

Costruendo un Nazionale Ricerca sull'IA Risorsa:

Un progetto per
la Nazionale
Nuvola di ricerca

CARTA BIANCA

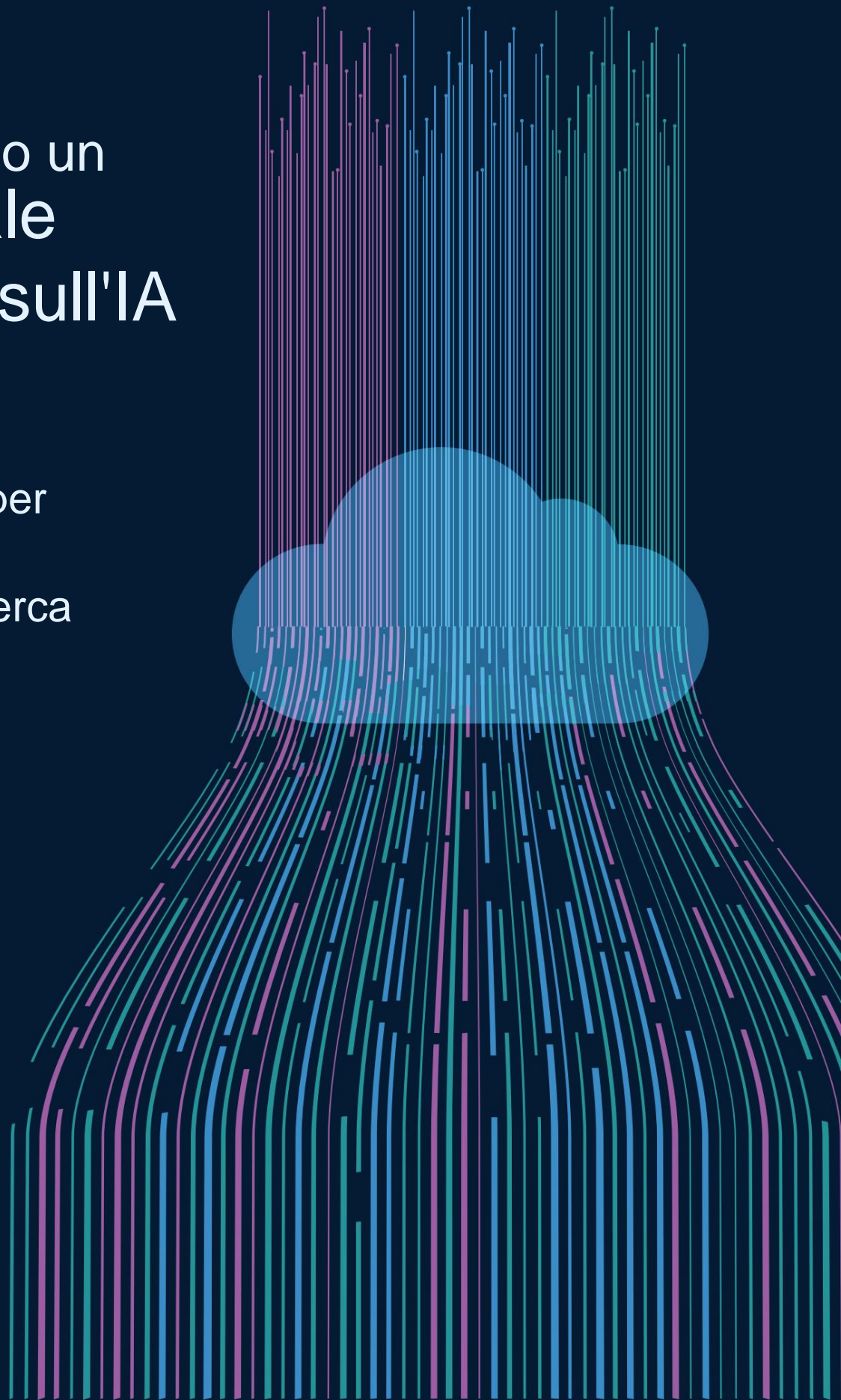
Daniel E. Ho
Jennifer King
Russell C. Wald
Christopher Wan

OTTOBRE 2021



Stanford University
Human-Centered
Artificial Intelligence

Stanford
Law School



Principali Autori

Daniel E. Ho, JD, Ph.D., è William Benjamin Scott e Luna M. Scott Professor of Law, Professore di Scienze Politiche e Senior Fellow presso lo Stanford Institute for Economic Policy Research presso la Stanford University. Dirige il Regulation, Evaluation, and Governance Lab (RegLab) a Stanford ed è Faculty Fellow presso il Center for Advanced Study in the Behavioral Sciences e Associate Director dello Stanford Institute for Human-Centered Artificial

Intelligenza (HAI). Ha conseguito il dottorato in giurisprudenza presso la Yale Law School e il dottorato di ricerca presso l'Università di Harvard e cancelliere per il giudice Stephen F. Williams presso la Corte d'Appello degli Stati Uniti per il Circuito del Distretto di Columbia.

Jennifer King, Ph.D., è la Privacy and Data Policy Fellow presso la Stanford HAI. Ha completato il suo dottorato in gestione e sistemi dell'informazione (scienza dell'informazione) presso l'Università della California, Berkeley School of Information. Prima di entrare in HAI, è stata direttrice della privacy dei consumatori presso il Center for Internet and Society della Stanford Law School dal 2018 al 2020.

Russell C. Wald è il direttore delle politiche per la Stanford HAI, a capo del team che promuove l'impegno di HAI con i governi e le organizzazioni della società civile. Dal 2013 ha ricoperto vari ruoli in affari governativi in rappresentanza della Stanford University. È membro a termine del Council on Foreign Relations, Visiting Fellow del National Security Institute della George Mason University e partner del Truman National Security Project. Si è laureato all'UCLA.

Christopher Wan è un candidato JD/MBA presso la Stanford University ed è stato assistente didattico per lo Stanford Policy Practicum: Creazione di un National Research Cloud. Serve anche come assistente di ricerca per la Stanford HAI e come investitore presso Bessemer Venture Partners. Ha ricevuto il suo

Laureato in informatica presso la Yale University, ha lavorato come ingegnere del software presso Facebook e come investitore di rischio presso In-Q-Tel e Tusk Ventures.

Lo Stanford Institute for Human-Centered Artificial Intelligence

Cordura Hall, 210 Panama Street, Stanford, CA 94305-4101

Ottobre 2021, V1.0

Contributori

Molte persone dedicate hanno contribuito a questo libro bianco. Per riconoscere questi contributi, elenchiamo qui i contributori per ogni capitolo e sezione.

Riepilogo esecutivo e introduzione Daniel E.

Ho, Tina Huang, Jennifer King, Marisa Lowe, Diego Núñez, Russell Wald, Christopher Wan, Daniel Zhang

La teoria per una nuvola di ricerca nazionale

Nathan Calvin, Shushman Choudhury, Tina Huang, Daniel E. Ho, Bitten Narayan, Diego Nunez, Frieda Rong, Russell Wald, Christopher Wan

Idoneità, allocazione e infrastruttura per l'informatica Daniel E. Ho,

Krithika Iyer, Tyler Robbins, Jasmine Shao, Russell Wald, Daniel Zhang

Protezione dell'accesso ai

dati Nathan Calvin, Shushman Choudhury, Daniel E. Ho, Ananya Karthik, Jennifer King, Christopher Wan

Design organizzativo

Sabina Beleuz, Drew Edwards, Daniel E. Ho, Jennifer King, Christopher Wan

Conformità alla privacy dei dati

Simran Arora, Neel Guha, Jennifer King, Sahaana Suri, Christopher Wan, Sadiki Wiltshire

Privacy tecnica e stanze sicure per dati virtuali

Neel Guha, Jennifer King, Christopher Wan

Salvaguardie per la ricerca etica Daniel

E. Ho, Jennifer King, Diego Núñez, Russell Wald, Daniel Zhang

Gestione dei rischi di sicurezza informatica

Neel Guha, Diego Núñez, Frieda Rong, Russell Wald

Proprietà intellettuale

Sabina Beleuz, Daniel E. Ho, Ananya Karthik, Diego Núñez, Christopher Wan

Casi studio

Daniel E. Ho, Krithika Iyer, Jennifer King, Marisa Lowe, Kanishka Narayan, Tyler Robbins

Vorremmo anche ringraziare Jeanina Casusi, Celia Clark, Shana Lynch, Kaci Peel, Stacy Peña, Mike Sellitto, Eun Sze e Michi Turner per il loro aiuto nella preparazione di questo White Paper.

Partecipanti esterni

Docenti ospiti e intervistati

Erik Brynjolfsson
Università di Stanford

Eric Horvitz
Microsoft

Brenda Leong
Il futuro di
Foro Privacy

Isabella Chu
Salute della popolazione
Scienze
Università di Stanford

Sara Giordano
Il futuro di
Foro Privacy

Amy O'Hara
Federale di Georgetown
Ricerca statistica
Banca dati
Università di Georgetown

Jack Clark
Antropico

Vince Kellen
UC San Diego,
CloudBank

Wade Shen
Attuare l'innovazione

Giovanni Etchemendy
Università di Stanford

Ed Lazowska
Università di
Washington

Suzanne Tallone
Calcola il Canada

Fei-Fei Li
Università di Stanford

Noemi Lefkowitz
Istituto Nazionale di
Norme e
Tecnologia (NIST)

Lee Tierich
Covington & Burling LLP

Marco Groman
Groman Consulting
Gruppo S.r.l

Eva Bianco
Laboratorio di politica della California
Università di Berkeley

Nel nostro processo, abbiamo anche coinvolto molti leader e sostenitori della società civile che hanno espresso molte prospettive sulla costruzione di un National Research Cloud. Abbiamo incorporato il loro feedback ove possibile e siamo grati per i loro pensieri condivisi e per averci aiutato a dare forma a un Libro bianco migliore.

Revisori

Ci siamo affidati a revisori esperti esterni straordinari per feedback e indicazioni. Ringraziamo Leisel Bogan, Jack Clark, John Etchemendy, Mark Krass, Marietje Schaake e Christine Tsang per la loro ponderata revisione del Libro bianco completo e ringraziamo Isabella Chu, Kathleen Creel, Luciana Herman, Sara Jordan, Vince Kellan, Brenda Leong, Ruth Marinshaw, Amy O'Hara e Lisa Ouellette per la loro competenza in materia su capitoli specifici.

Partecipanti al workshop

Il 5 agosto 2021, i coautori hanno ospitato una sessione di feedback per ascoltare una varietà di parti interessate del mondo accademico, della società civile, del governo e dell'industria. Siamo grati per il tempo e per i consigli utili che i partecipanti hanno offerto. Le affiliazioni dei partecipanti al workshop sono elencate solo a scopo identificativo. Al seminario hanno partecipato anche individui di Microsoft e AI Now, ma non hanno voluto essere identificati personalmente.

Capra Taka Responsabilità del governo Ufficio	Ciro Hodes La società futura Sara Giordano Forum sul futuro della privacy Vince Kellen Università di San Diego, CloudBank Michael Kratsios Scale AI Samanta Lai Istituto Brookling Brenda Leong Forum sul futuro della privacy	Asad Ramzanali Camera degli Stati Uniti di Rappresentanti Davide Robinson Ripresa Saiph Selvaggio Università del Nordest Michele Sellitto Stanford HAI Ishan Sharma Federazione di Scienziati americani Brittany Smith Dati e società John Smith IBM Brittany Smith Dati e società Victor Storchan JP Morgan Chase Keith Strier Attività di calcolo AI Organizzazione della forza per l'economia Cooperazione e Sviluppo (OCSE)
Katie Baxter Salesforce		
Leisel Bogan Centro Belfer Università di Harvard		
Jeffrey Brown IBM		
Miles Brundage IA aperta		
L. Jean Camp Università dell'Indiana a Bloomington	Ruth Marinshaw Ricerca di Stanford Centro di calcolo	
Dakota Cari Centro per la sicurezza e Tecnologia emergente Università di Georgetown	Giosuè Meltzer Istituto Brookling Sam Mulopulous al Senato degli Stati Uniti	
Shikai Chern Tecnologie Veritas		
Isabella Chu Scienze della salute della popolazione Università di Stanford	Dewey Murdick Centro per la sicurezza e Tecnologia emergente Università di Georgetown	
Jack Clark Antropico	Ricco Omar Centro per l'innovazione dei dati	
Meaghan inglese Patrick J. McGovern Fondazione	Calton Pu Georgia Tech	Lee Tierich Covington & Burling LLP Eva Bianco Laboratorio di politica della California Università di Berkeley

A proposito di HAI

L'Institute for Human-Centered Artificial Intelligence (HAI) della Stanford University applica analisi e ricerche rigorose a pressanti questioni politiche sull'intelligenza artificiale. Un pilastro di HAI è informare i responsabili politici, i leader del settore e la società civile diffondendo borse di studio a un vasto pubblico.

HAI è un istituto di ricerca apartitico, che rappresenta una gamma di voci. Le opinioni espresse in questo libro bianco riflettono le opinioni degli autori.

Informazioni su SLS Policy Lab

Il Policy Lab della Stanford Law School offre agli studenti un'esperienza coinvolgente nella ricerca di soluzioni ad alcuni dei problemi più urgenti del mondo sotto la direzione della facoltà e dei ricercatori di Stanford.

Diretto dall'ex SLS Dean Paul Brest, il Policy Lab riflette la convinzione della scuola che l'esame sistematico dei problemi della società, informato da una ricerca rigorosa, possa generare soluzioni ai problemi pubblici più impegnativi della società.

Indipendenza accademica

Questo libro bianco è stato sviluppato in modo indipendente dal gruppo di ricerca. Sebbene abbiamo sollecitato il feedback da un'ampia gamma di parti interessate, nessun donatore, società o altra parte interessata HAI ha avuto alcun coinvolgimento con la ricerca e la produzione di questo Libro bianco. Secondo la politica HAI, "i donatori non possono dettare argomenti di ricerca perseguiti dai ricercatori HAI" né "controllare il permesso di pubblicare i risultati della ricerca". Per ulteriori informazioni, consultare [la politica di HAI: https://hai.stanford.edu/about/fundraising-policy](https://hai.stanford.edu/about/fundraising-policy).

Divulgazioni

La Stanford University si è attivamente impegnata e ha fatto pressioni sul Congresso per approvare il *National Artificial Intelligence Research Resource Task Force Act*, lavorando con una coalizione di parti interessate del mondo accademico, della società civile e del settore. Il coautore Russell Wald ha fornito supporto agli sforzi di advocacy.

Il co-direttore di HAI Fei-Fei Li, che ha servito come docente ospite in classe, è stato uno dei primi sostenitori di una task force per studiare il National Research Cloud. Il Dr. Li è stato nominato membro della Task Force della National Artificial Intelligence Research Resource (NAIRR).

Il coautore Daniel Ho dirige lo Stanford RegLab, che ha ricevuto supporto di calcolo dal programma di crediti cloud di HAI (AWS e GCP), dal programma di concessione di crediti di calcolo AI for Earth Azure di Microsoft e dalla concessione di crediti cloud di Google per la ricerca COVID-19.

La co-autrice Jennifer King ha ricevuto finanziamenti illimitati per la ricerca da Mozilla, Facebook e Accenture nel suo precedente ruolo presso il Center for Internet and Society.

Lo Stanford Institute for Human-Centered Artificial Intelligence (HAI) riceve supporto finanziario e di cloud computing da A121 Labs, Amazon Web Services, Google, IBM, Microsoft e OpenAI.

Sommario

EXECUTIVE SUMMARY: Creare un National Research Cloud	9
INTRODUZIONE	15
CAPITOLO 1: Una teoria per un National Research Cloud	17
CAPITOLO 2: Ammissibilità, allocazione e infrastruttura per l'informatica	22
CAPITOLO 3: Protezione dell'accesso ai dati	35
CAPITOLO 4: Progettazione organizzativa	48
CAPITOLO 5: Conformità alla privacy dei dati	53
CAPITOLO 6: Privacy tecnica e stanze sicure per dati virtuali	61
CAPITOLO 7: Salvaguardie per la ricerca etica	66
CAPITOLO 8: Gestione dei rischi per la sicurezza informatica	70
CAPITOLO 9: Proprietà intellettuale	76
GLOSSARIO DEGLI ACRONIMI	82
APPENDICE	84
NOTE DI CHIUSURA	90

Casi studio

MODELLI DI CALCOLO

NSF CloudBank	27
XSEDE NSF	29
Fugaku	32
Calcola il Canada	34

MODELLI DI DATI

Iniziativa Coleridge	42
Scienze della salute della popolazione di Stanford	43
La legge sulle prove	46

MODELLI ORGANIZZATIVI

Istituto di politica scientifica e tecnologica	50
Partnership dati dell'Alberta	51

ALTRI MODELLI

Ricerca sui dati amministrativi nel Regno Unito	58
Laboratorio di politica della California	64

Sintesi:

Creazione di un cloud di ricerca nazionale

L'intelligenza artificiale (AI) sembra pronta a trasformare l'economia in settori che vanno dalla sanità alla finanza al commercio al dettaglio e all'istruzione. Quella che alcuni hanno coniato la "Quarta rivoluzione industriale"¹ è guidata da tre tendenze chiave: maggiore disponibilità di dati, aumento della potenza di calcolo e miglioramento della progettazione degli algoritmi. In primo luogo, quantità sempre maggiori di dati hanno alimentato la capacità di apprendimento dei computer, ad esempio addestrando un modello di linguaggio algoritmico su tutta Wikipedia. modelli che erano inimmaginabili solo 10 anni fa, che a volte abbracciano miliardi di parametri (un aumento esponenziale della portata rispetto ai modelli precedenti). per sconfiggere il campione del mondo nel gioco da tavolo Go.⁴

Storicamente, le partnership tra governi, università e industrie hanno ancorato l'ecosistema dell'innovazione statunitense. Il governo federale ha svolto un ruolo fondamentale nel sovvenzionare la ricerca di base, consentendo alle università di intraprendere ricerche ad alto rischio che possono richiedere decenni per essere commercializzate. Questo approccio ha catalizzato la tecnologia radar, Internet e i dispositivi GPS. Come affermano gli economisti Ben Jones e Larry Summers, "[e] anche sotto ipotesi molto prudenti, è difficile trovare un rendimento medio inferiore a \$ 4 per \$ 1 speso" per l'innovazione, e i rendimenti sociali potrebbero essere più vicini a \$ 20 per ogni dollaro esaurito.⁵ L'industria, a sua volta, ridimensiona e commercializza le applicazioni.

SFIDE ALL'ECOSISTEMA DELL'INNOVAZIONE AI

Tuttavia, questo ecosistema dell'innovazione deve affrontare serie sfide potenziali. La potenza di calcolo è diventata fondamentale per il progresso dell'IA, ma l'alto costo del calcolo ha posto la ricerca IA all'avanguardia in una posizione accessibile solo ai principali attori del settore e a una manciata di università d'élite.⁶ Accesso ai dati: gli ingredienti grezzi utilizzati per addestrare la maggior parte dei modelli di intelligenza artificiale — è sempre più limitato al settore privato e alle grandi piattaforme⁷, dal momento che le fonti di dati governative rimangono in gran parte inaccessibili alla comunità di ricerca sull'IA. l'industria dell'intelligenza artificiale minaccia la competitività tecnologica degli Stati Uniti.⁹

Quattro sfide correlate illustrano questa scoperta: in primo luogo, stiamo assistendo a una significativa fuga di cervelli di ricercatori che abbandonano le università. I .D vanno nell'industria e meno di un quarto nel mondo accademico.¹² In secondo luogo, queste tendenze indicano che molti ricercatori universitari faticano a impegnarsi nella scienza d'avanguardia, prosciugando il campo del variegato insieme di voci di ricerca di cui ha bisogno. In terzo luogo, la ricerca fondamentale che garantirebbe agli Stati Uniti di rimanere al timone dell'innovazione dell'IA viene messa da parte. Secondo una stima, l'82% degli algoritmi utilizzati oggi proviene da organizzazioni non profit e università finanziate a livello federale, ma "la leadership degli Stati Uniti è svanita negli ultimi decenni". il costo delle funzioni di governance fondamentali e il miglioramento della capacità interna del governo di sviluppare, testare e ritenere responsabili i sistemi di intelligenza artificiale.¹⁵ In breve, un crescente squilibrio nell'innovazione dell'IA si inclina verso l'industria, lasciando indietro la ricerca accademica e non commerciale.

Dato il ruolo di lunga data della ricerca accademica e non commerciale nell'innovazione, questo spostamento ha sostanziali conseguenze negative per l'ecosistema della ricerca americana.

LA RICERCA NAZIONALE AI RISORSA TASK FORCE ACT

In risposta a queste sfide, il Congresso ha promulgato il National AI Research Resource Task Force Act as parte del National Defense Authorization Act (NDAA) nel gennaio 2021.¹⁶ La legge fa parte della National Artificial Intelligence Initiative, che identifica ulteriori passi per aumentare gli investimenti nella ricerca, stabilire standard tecnici e costruire una forza lavoro IA più forte. La legge ha creato una Task Force – la cui composizione è stata annunciata il 10 giugno 2021¹⁷ – per studiare e pianificare l'implementazione di una "National Artificial Intelligence Research Resource" (NAIRR), vale a dire "un sistema che fornisce a ricercatori e studenti tutta la campi e discipline con accesso a risorse di calcolo, collocate insieme a set di dati governativi e non governativi pronti per l'intelligenza artificiale disponibili al pubblico."¹⁸ Questa risorsa di ricerca è stata anche denominata National Research Cloud (NRC) ed è stata fortemente sostenuta dall'NSCAI, che ha scritto che l'NRC "rafforzerà le basi dell'innovazione americana dell'IA sostenendo una crescita più equa del settore, espandendo le competenze dell'IA in tutto il paese e applicando l'IA a una gamma più ampia di campi".¹⁹

L'NRC dirige più risorse

verso lo sviluppo dell'IA nell'interesse pubblico e aiuta a garantire la leadership a lungo termine degli Stati Uniti nel settore, sostenendo il tipo di ricerca pura e di base che il settore privato non può intraprendere da solo.

Mentre altre iniziative hanno cercato di migliorare l'accesso al calcolo o ai dati in isolamento,²⁰ l'NRC genererà externalità positive distinte integrando il calcolo e i dati, i due colli di bottiglia per la ricerca sull'IA di alta qualità. In particolare, l'NRC fornirà un accesso conveniente a risorse computazionali di fascia alta, set di dati governativi su larga scala in un ambiente cloud sicuro e le competenze necessarie per beneficiare di questa risorsa attraverso una stretta collaborazione tra il mondo accademico, il governo e l'industria. Ampliando l'accesso a queste risorse critiche nella ricerca sull'IA, l'NRC sosterrà la ricerca scientifica di base sull'IA, la democratizzazione dell'innovazione dell'IA e la promozione della leadership statunitense nell'IA.

TEMI

Il programma Policy Lab della Stanford Law School ha riunito un gruppo di ricerca multidisciplinare di studenti laureati, personale e docenti provenienti dalle scuole di economia, diritto e ingegneria di Stanford per studiare la fattibilità e le considerazioni per la progettazione dell'NRC. Negli ultimi sei mesi, questo gruppo ha studiato i modelli esistenti per le risorse di calcolo e i dati governativi, ha intervistato un'ampia gamma di esperti governativi, informatici e politici ed ha esaminato i requisiti tecnici, aziendali, legali e politici. Questo White Paper è stato commissionato dall'Institute for Human-Centered Artificial Intelligence (HAI) di Stanford, che ha originato la proposta per l'NRC in collaborazione con altre 21 università di ricerca.²¹

Nel corso della nostra ricerca, abbiamo osservato tre temi principali che attraversano tutte le aree della nostra indagine. Abbiamo integrato questi temi in ogni sezione del nostro White Paper e ne ha tratto spunto per spiegare i nostri risultati.

- Complementarità tra calcolo e dati. Durante la valutazione degli ecosistemi di elaborazione e condivisione dei dati esistenti, una delle sfide sistemiche che abbiamo osservato è stata il disaccoppiamento delle risorse di elaborazione dalle infrastrutture di dati.

L'elaborazione ad alte prestazioni può essere inutile senza dati e uno dei principali ostacoli alla condivisione dei dati, in particolare per i dati governativi di alto valore, risiede nei requisiti per un ambiente informatico sicuro e che protegga la privacy.

- Riequilibrare la ricerca sull'IA verso la ricerca a lungo termine, accademica e non commerciale. Attualmente, l'innovazione dell'IA dipende in modo sproporzionato dal settore privato. Gli investimenti pubblici nelle infrastrutture di IA di base possono sia sostenere l'innovazione nell'interesse pubblico sia integrare gli sforzi di innovazione privati. L'NRC indirizza più risorse verso lo sviluppo dell'IA nell'interesse pubblico e aiuta a garantire la leadership a lungo termine degli Stati Uniti nel settore, sostenendo il tipo di ricerca di base pura che il settore privato non può intraprendere da solo.
- Coordinare gli approcci a breve e a lungo termine per la creazione dell'NRC. La nostra ricerca prende in considerazione molti percorsi a breve termine per sostenere una versione funzionante dell'NRC spiegando come lavorare all'interno dei vincoli esistenti. Identifichiamo inoltre le sfide strutturali, legali e politiche da affrontare a lungo termine per realizzare la visione completa dell'NRC.

Riassumiamo qui le nostre raccomandazioni principali.

MODELLO DI CALCOLO

- **La decisione "Make or Buy".** La scelta politica principale sarà se costruire un'infrastruttura informatica pubblica o acquistare servizi da fornitori di cloud commerciali esistenti.
 - È risaputo che, basandosi esclusivamente sui costi dell'hardware, è più conveniente possedere un'infrastruttura quando la domanda informatica è quasi continua.²² Il governo ha anche esperienza nella creazione di cluster di elaborazione ad alte prestazioni, generalmente costruiti da appaltatori e gestiti dai laboratori nazionali.²³ Anche la National Science Foundation (NSF) ha sostenuto molte iniziative di supercalcolo presso le istituzioni accademiche.²⁴
 - Le principali preoccupazioni compensative sono che i fornitori di cloud commerciali esistenti dispongono di stack software e usabilità che i ricercatori di intelligenza artificiale hanno ampiamente adottato e potrebbero considerare una piattaforma più user-friendly. I fornitori di servizi cloud commerciali offrono un modo per espandere rapidamente la capacità, sebbene la scalabilità e la disponibilità saranno ancora limitate dalla disponibilità delle attuali risorse di elaborazione dell'unità di elaborazione grafica (GPU).
 - Consigliamo una duplice strategia di investimento:
 - In primo luogo, il modello di calcolo dell'NRC può essere lanciato rapidamente sovvenzionando e negoziando il cloud computing per i ricercatori di intelligenza artificiale con i fornitori esistenti, espandendo le iniziative esistenti come il progetto CloudBank dell'NSF.²⁵
 - In secondo luogo, l'NRC dovrebbe investire in un progetto pilota per infrastrutture pubbliche per valutare la capacità di fornire risorse simili nel lungo periodo. Tale infrastruttura di proprietà pubblica sarebbe ancora costruita sotto contratto o

Una delle sfide sistemiche [per la ricerca di base sull'IA è] il disaccoppiamento delle risorse di calcolo dalle infrastrutture di dati. . . . [Un] ambiente informatico sicuro e che protegga la privacy [sarà fondamentale]

sovvenzione, ma potrebbe essere gestito in modo molto simile ai laboratori nazionali (ad esempio, Sandia National Laboratories, Oak Ridge National Laboratory) che possiedono sofisticate strutture di supercalcolo o strutture accademiche di supercalcolo.

- **Idoneità del ricercatore.** Mentre alcuni hanno sostenuto che l'NRC dovrebbe essere aperto all'accesso commerciale, ai fini di questo Libro bianco, abbiamo aderito allo spirito della legislazione che forma la NAIIR Task Force e abbiamo rivisto solo l'uso di un NRC per la ricerca sull'IA accademica e senza scopo di lucro. Raccomandiamo che l'idoneità NRC inizi con accademici che detengono lo status di "Principal Investigator" (PI) (ovvero, la maggior parte dei docenti) presso college e università statunitensi, nonché "agenzie governative affiliate" disposti a contribuire con set di dati di alto valore inediti a l'NRC in cambio di risorse di calcolo sovvenzionate. Lo stato PI dovrebbe essere interpretato in modo espansivo per comprendere tutti i campi di applicazione dell'IA. Gli studenti che lavorano con i PI dovrebbero presumibilmente ottenere l'accesso all'NRC. Ridimensionare l'NRC per soddisfare la domanda di tutti gli studenti negli Stati Uniti può essere impegnativo, ma raccomandiamo anche la creazione di programmi educativi come parte della nuova risorsa per aiutare a formare la prossima generazione di ricercatori di intelligenza artificiale.
- **Meccanismo.** Per mantenere bassi i costi di elaborazione dei premi, consigliamo un livello base di accesso al calcolo per soddisfare la maggior parte delle esigenze informatiche dei ricercatori. L'accesso di livello base evita un sovraccarico elevato per l'amministrazione delle sovvenzioni e può soddisfare le esigenze di elaborazione per la stragrande maggioranza dei ricercatori. Per i ricercatori con esigenze eccezionali, consigliamo un processo di sovvenzione semplificato per un ulteriore accesso al calcolo.

MODELLO DI ACCESSO AI DATI

- **Focus sui dati del governo.** Concentriamo le nostre raccomandazioni per la fornitura di dati/l'accesso ai dati governativi perché: (1) esiste già un'ampia gamma di piattaforme per la condivisione di dati privati,²⁶ e (2) la distribuzione da parte dell'NRC di set di dati privati solleverebbe un groviglio di spinose questioni di proprietà intellettuale. Raccomandiamo che i ricercatori siano autorizzati a calcolare su tutti i set di dati che essi stessi contribuiscono, a condizione che certifichino di avere i diritti su tali dati e che l'uso di tali dati sia per scopi di ricerca accademica.
- **Accesso a più livelli.** Consigliamo un modello di accesso a più livelli: per impostazione predefinita, i ricercatori avranno accesso ai dati del governo che sono già pubblici; i ricercatori possono quindi presentare domanda attraverso un processo semplificato per ottenere l'accesso a livelli di sicurezza più elevati in base a un progetto specifico. Sarà fondamentale per l'NRC alla fine sostituire l'attuale approccio relazionale frammentato, agenzia per agenzia. Fornendo ambienti virtuali sicuri e armonizzando gli standard di sicurezza (ad esempio, Federal Risk and Authorization Management Program (FedRAMP)²⁷), l'NRC può collaborare con proposte per un National Secure Data Service²⁸ per fornire un modello per accelerare la ricerca sull'IA, proteggendo al contempo la riservatezza dei dati e priorità alla sicurezza dei dati.
- **Incentivi di agenzia.** Per incentivare le agenzie federali a condividere i dati con l'NRC e migliorare lo stato della tecnologia del settore pubblico, consigliamo all'NRC di consentire al personale dell'agenzia federale di utilizzare le risorse di calcolo dell'NRC. In linea con le pratiche dei programmi di condivisione dei dati esistenti, come la Coleridge Initiative,²⁹ raccomandiamo inoltre che l'NRC fornisca formazione e supporto per lavorare con le agenzie per modernizzare e armonizzare i loro standard di dati.
- **Investimento strategico per le fonti di dati.** A breve termine, raccomandiamo che l'NRC concentri i suoi sforzi sulla messa a disposizione di set di dati governativi non sensibili, a rischio da basso a moderato, piuttosto che dati governativi sensibili (ad esempio, dati su individui) o dati del settore privato, a causa di privacy dei dati e problemi di proprietà intellettuale. I ricercatori possono ancora utilizzare le risorse di calcolo NRC sui dati privati, ma dovrebbero fare affidamento sui meccanismi esistenti per acquisire i dati per i propri bucket privati sull'NRC. Ad esempio, le immagini prese dai satelliti di osservazione della Terra, come

Le immagini Landsat forniscono un promettente set di dati governativi a basso rischio e ad alto rendimento, poiché rendere tali immagini satellitari liberamente disponibili ai ricercatori ha generato benefici economici annuali stimati in 3-4 miliardi di dollari, in particolare se combinati con il calcolo ad alte prestazioni.³⁰ Agenzie come la National Oceanic and Atmospheric Administration, l'US Geological Survey, il Census Bureau, l'Administrative Office of the US Courts e il Bureau of Labor Statistics, ad esempio, hanno anche ricchi set di dati che possono essere distribuiti più facilmente. A lungo termine, l'accesso a set di dati ad alto rischio, come quelli di proprietà dell'Internal Revenue Service (IRS) e del Department of Veterans Affairs (VA), dipenderà dal modello di accesso a più livelli.

FORMA ORGANIZZATIVA

Dove collocare istituzionalmente l'NRC pone un compromesso tra la facilità di coordinamento per ottenere il calcolo e la facilità di accesso ai dati. Ad esempio, l'ubicazione dell'NRC all'interno di una singola agenzia renderebbe più facile il coordinamento con i fornitori di servizi informatici, ma renderebbe più difficile l'accesso ai dati tra le agenzie, in assenza di ulteriore autorità statutaria. Molti sforzi per rendere più facile l'accesso ai dati del governo, in particolare il Foundations for Evidence-Based Policymaking Act del 2018, si sono rivelati tra le sfide più scoraggianti della modernizzazione del governo. Istituito come centro di ricerca e sviluppo finanziato a livello federale (FFRDC) nel breve periodo e partenariato pubblico-privato (PPP) nel lungo periodo.

- **FFRDC.** Gli FFRDC presso le agenzie governative affiliate ridurrebbero i costi significativi per proteggere i dati da quelli agenzie ospitanti. Questo approccio sarà anche coerente con la maggiore dipendenza dai crediti cloud commerciali nel breve periodo, rendendo meno centrale il coordinamento del calcolo e dei dati. A lungo termine, tuttavia, il coordinamento semplificato tra dati e calcolo può essere più difficile con FFRDC ospitati presso agenzie specifiche quando (1) l'NRC si allontana dai crediti cloud commerciali e si sposta verso il proprio cluster di calcolo ad alte prestazioni e (2) un diventa disponibile un numero maggiore di set di dati interagenzia.
- **PPP.** A lungo termine, raccomandiamo la creazione di un modello PPP, governato da funzionari di agenzie governative affiliate, ricercatori accademici e rappresentanti del settore tecnologico, che possa ospitare sia risorse di calcolo che di dati.

ULTERIORI CONSIDERAZIONI

- **Privacy dei dati.** Come prima cosa, un NRC da cui provengono dati amministrativi sensibili o identificabili individualmente più agenzie vengono utilizzate per costruire e addestrare modelli di intelligenza artificiale dovranno affrontare le sfide del Privacy Act del 1974.³² La legge ha lo scopo di porre un controllo sulla condivisione dei dati tra agenzie e sulla divulgazione di dati sensibili senza consenso.
 - ° Al fine di evitare conflitti con la condivisione non consensuale dei dati tra agenzie, raccomandiamo che l'NRC non sia istituito come propria agenzia federale, né che al personale dell'agenzia federale sia consentito l'accesso ai dati tra agenzie.
 - ° Per evitare conflitti con il requisito della legge "nessuna divulgazione senza consenso", qualsiasi dato rilasciato all'NRC non deve essere identificabile individualmente. Nonostante questi vincoli, la maggior parte della ricerca sull'intelligenza artificiale rientrerà probabilmente nell'eccezione di ricerca statistica della legge, subordinata alle proposte in linea con lo scopo principale di un'agenzia.
 - ° Date le preoccupazioni sui potenziali rischi per la privacy, le agenzie federali potrebbero desiderare di condividere i dati, subordinatamente all'uso di misure tecniche per la privacy (ad esempio, privacy differenziale). Sebbene utile in molti casi, tecnico

gli approcci non sono una panacea e non dovrebbero sostituire le politiche di accesso ai dati.

° L'NRC dovrebbe esplorare la progettazione di "stanze sicure per i dati" virtuali che consentano ai ricercatori di accedere ai dati in un ambiente sicuro, monitorato e basato su cloud.

° Ulteriori interventi legislativi potrebbero anche facilitare la condivisione dei dati con l'NRC (ad esempio, richiedendo la modernizzazione dell'IT per includere piani di condivisione dei dati con l'NRC).

- **Etica.** La rapida innovazione nella ricerca sull'IA solleva una serie di potenziali sfide etiche. Data la portata dell'NRC, non sarà possibile esaminare ogni singola proposta di ricerca per potenziali violazioni etiche, soprattutto perché gli standard etici sono ancora in evoluzione. L'NRC dovrebbe adottare un duplice approccio.

° In primo luogo, per l'accesso predefinito delle PI ai dati e al calcolo di base, l'NRC dovrebbe istituire un processo di revisione ex post per le accuse di violazioni della ricerca etica. L'accesso può essere revocato quando si dimostri che la ricerca viola in modo manifesto e grave gli standard etici. Sottolineiamo che lo standard elevato per una violazione dovrebbe essere informato dalle implicazioni del discorso accademico e dalle potenziali conseguenze politiche del coinvolgimento del governo nell'amministrazione dell'NRC e nella determinazione delle direzioni della ricerca accademica.

° In secondo luogo, per le applicazioni che richiedono l'accesso a set di dati limitati o risorse oltre il calcolo predefinito, che saranno necessariamente sottoposte a revisione, i ricercatori dovrebbero essere tenuti a fornire una dichiarazione sull'impatto etico. Uno dei vantaggi di iniziare con i PI è che i docenti universitari sono responsabili in base agli IRB esistenti per la ricerca sui soggetti umani, nonché per i principi della revisione tra pari.

° Esortiamo le parti non NRC (ad esempio, le università) a esplorare una serie di misure per affrontare le preoccupazioni etiche nel calcolo dell'IA (ad esempio, un processo di revisione etica³³ o incorporare esperti di etica nei progetti³⁴).

- **Sicurezza.** Raccomandiamo che l'NRC prenda l'iniziativa nella definizione di classificazioni e protocolli di sicurezza, in parte per contrastare un sistema di sicurezza balcanizzato tra le agenzie federali che ostacolerebbe la capacità di ospitare set di dati. L'NRC dovrebbe utilizzare personale di sicurezza dedicato per lavorare con le agenzie governative affiliate e i rappresentanti delle università per armonizzare e modernizzare gli standard di sicurezza delle agenzie.
- **Proprietà Intellettuale (IP).** Sebbene le prove sugli incentivi IP ottimali per l'innovazione siano contrastanti, raccomandiamo che l'NRC adotti lo stesso approccio all'assegnazione di diritti di brevetto, diritti d'autore e diritti sui dati agli utenti dell'NRC che si applicano agli accordi di finanziamento federali. L'NRC dovrebbe inoltre prendere in considerazione le condizioni per richiedere ai ricercatori dell'NRC di divulgare o condividere i loro risultati di ricerca con una licenza ad accesso aperto.
- **Risorse umane.** Data la sua ambizione, per rendere l'NRC un successo saranno necessarie risorse umane significative, dagli ingegneri di sistema ai responsabili dei dati, dagli amministratori delle sovvenzioni al personale per la privacy, l'etica e la sicurezza informatica.

Data la sua ambizione, per rendere
l'NRC un successo saranno necessarie
risorse umane significative, dagli
ingegneri di sistema ai responsabili
dei dati, dagli amministratori delle
sovvenzioni al personale per la privacy,
l'etica e la sicurezza informatica.

introduzione

Nel marzo 2020, l'Institute for Human-Centered Artificial Intelligence (HAI) di Stanford ha pubblicato una lettera aperta, co-firmata dai presidenti e dai rettori delle 22 migliori università del paese, al presidente degli Stati Uniti e al Congresso degli Stati Uniti sollecitando l'adozione di un National Research Cloud (NRC).¹ La proposta NRC mira a colmare un divario significativo nell'accesso all'informatica e ai dati che, sostengono i sostenitori, ha distorto la traiettoria a lungo termine della ricerca sull'intelligenza artificiale (AI).² Senza accesso a tali risorse critiche, la ricerca sull'intelligenza artificiale può essere dominata da interessi commerciali a breve termine e minare l'ecosistema dell'innovazione storica in cui la ricerca di base, fondamentale e non commerciale ha gettato le basi per applicazioni che potrebbero essere lontane decenni, non ancora commerciabili o promuovere l'interesse pubblico.

Nel gennaio 2021, il Congresso ha promulgato il National Artificial Intelligence Research Resource Task Force Act (NAIRR), costituendo una task force per esaminare il progetto dell'NRC.³ La task force è stata annunciata nel giugno di quest'anno e include uno dei sostenitori originali dell'NRC e co-direttore di HAI (Fei-Fei Li).⁴

Questo White Paper è il culmine di un praticantato politico indipendente di due quarti presso il programma Policy Lab della Stanford Law School, che è stato co-insegnato da tre di noi (Ho, King, Wald) e un assistente didattico (Wan) e ha riunito giurisprudenza, affari e studenti di ingegneria per contemplare le dimensioni chiave del progetto dell'NRC. Abbiamo intervistato e convocato un'ampia gamma di parti interessate, inclusi avvocati per la privacy, tecnici di cloud computing, esperti di dati governativi, professionisti della sicurezza informatica, potenziali utenti e gruppi di interesse pubblico. Gli studenti hanno studiato le disposizioni legali in materia, le opzioni politiche e le strade per la progettazione istituzionale dell'NRC. Il team del praticantato ha lavorato in modo indipendente per dare forma alle sue raccomandazioni.

La proposta per un NRC è ambiziosa e questo Libro bianco copre molti argomenti. Iniziamo con la domanda fondamentale - perché costruire l'NRC (capitolo 1)? - e spieghiamo quella che consideriamo una teoria convincente dell'impatto. Poi esaminiamo chi dovrebbe avere accesso all'NRC (capitolo 2), cosa comprende l'NRC (capitolo 2), come l'accesso ai dati riservati può (o meno) essere concesso (capitolo 3) e dove dovrebbe essere localizzato l'NRC (Capitolo 4). Dedichiamo molto tempo alla parte relativa all'accesso ai dati (capitoli 3, 5 e 6), a causa delle complessità della condivisione dei dati da parte del governo ai sensi del Privacy Act del 1974.⁵ Come notiamo in questi capitoli, la parte relativa ai dati dell'NRC è complementare a sforzi di lunga data per consentire un maggiore accesso della ricerca ai dati amministrativi ai sensi, ad esempio, del Foundations for Evidence-Based Policymaking Act del 2018⁶ e della proposta del National Secure Data Service Act.⁷ Tale condivisione deve essere effettuata in modo sicuro e in modo da proteggere la privacy. Consideriamo anche questioni di standard etici (capitolo 7), sicurezza informatica (capitolo 8) e proprietà intellettuale (capitolo 9) che informano la progettazione dell'NRC.

Riconosciamo la complessità dell'impresa e che ci sono molte domande a cui non viene data risposta. La portata contemplata dell'NRC potrebbe essere per l'intelligenza artificiale ciò che il Progetto genoma umano è stato per la genomica (o ciò che gli acceleratori di particelle sono stati per la fisica): investimenti pubblici per una ricerca scientifica fondamentale ambiziosa e non commerciale per garantire lo sviluppo a lungo termine di un'area critica di innovazione per gli Stati Uniti. Ci sono molte aree in cui avremmo voluto avere l'opportunità di impegnarci in ricerche più approfondite. Ci auguriamo tuttavia che questo White Paper fornisca un utile contributo alla task force NAIRR, al Congresso, alla Casa Bianca e a tutti coloro che sono interessati all'ecosistema dell'innovazione AI.

Dobbiamo gratitudine alle molte persone che hanno contribuito con tempo, feedback e approfondimenti. Soprattutto, ringraziamo gli straordinari studenti che hanno dato forma a questo White Paper: Simran Arora, Sabina Beleuz, Nathan Calvin, Shushman Choudhury, Drew Edwards, Neel Guha, Krithika Iyer, Ananya Karthik, Kanishka Narayan, Tyler Robbins, Frieda Rong, Jasmine Shao e Sadiki Wiltshire. Abbiamo beneficiato di troppe persone da nominare, ma un ringraziamento speciale va a Taka Ariga, Kathy Baxter, Miles Brundage, Jean Camp, Shikai Chern, Bella Chu, Jack Clark, Kathleen Creel, John Etchemendy, Deep Ganguli, Eric Horvitz, Sara Jordan, Vince Kellen, Mark Krass, Sebastien Krier, Ed Lazowska, Brenda Leong, Fei-Fei Li, Ruth Marinshaw, Michelle Mello, Amy O'Hara, Hodan Omaar, Saiph Savage, Marietje Schaake, Mike Sellitto, Wade Shen, Keith Strier, Suzanne Talon, Lee Tiedrich, Christine Tsang ed Evan White per utili approfondimenti e feedback. Il personale HAI e gli assistenti di ricerca che sono stati essenziali nell'aiutarci durante le fasi finali della redazione e della compilazione del Libro bianco includono Tina Huang, Marisa Lowe, Diego Núñez e Daniel Zhang.

Come spieghiamo in questo Libro bianco, l'NRC è un'idea che vale la pena prendere sul serio. Vale la pena chiarire, tuttavia, cosa risolverebbe e non risolverebbe. L'NRC consentirebbe un accesso molto maggiore - e in tal senso, democratizzerebbe - forme di intelligenza artificiale e ricerca sull'IA che sono aumentate nelle richieste computazionali, ma non impedirebbe o sposterebbe categoricamente la centralizzazione del potere all'interno dell'industria tecnologica. L'NRC sposterebbe l'attenzione degli attuali sforzi di IA su dimensioni più pubbliche e socialmente guidate fornendo l'accesso a set di dati governativi precedentemente limitati, affrontando gli sforzi di lunga data per migliorare l'accesso a dati del settore pubblico di alto valore, ma non creerebbe un sistema per impedire a tutti usi non etici dell'intelligenza artificiale. L'NRC faciliterebbe gli audit di modelli, set di dati e sistemi di intelligenza artificiale su larga scala per violazioni della privacy e pregiudizi, ma non equivarrebbe a un requisito normativo per valutazioni di equità e responsabilità. Non è né uno strumento dell'antitrust né un organismo di certificazione per algoritmi etici, che sono aree che meritano di essere prese seriamente in proposte politiche indipendenti . progettazione dell'NRC.

Sebbene da solo non possa risolvere tutto ciò che affligge l'IA, l'NRC promette di fare un importante passo avanti affermativo.

Capitolo 1:

La teoria per a Nube di ricerca nazionale

Questo capitolo articola una teoria dell'impatto per l'NRC. Nella politica convenzionale termini analitici, 1 quale problema (o fallimento del mercato) affronta l'NRC? Da un certo punto di vista, l'innovazione dell'IA è vivace negli Stati Uniti, con importanti progressi nel linguaggio, nella visione, nei dati strutturati e nelle applicazioni che si sviluppano in tutti i settori. Tuttavia, da un'altra prospettiva, l'attuale commercializzazione dell'innovazione passata maschera un sistematico sottoinvestimento nella ricerca di IA di base e non commerciale che potrebbe garantire la salute a lungo termine dell'innovazione tecnologica in questo paese.

L'attuale commercializzazione dell'innovazione passata maschera un sistematico sottoinvestimento nella ricerca di base sull'IA non commerciale che potrebbe garantire la salute a lungo termine dell'innovazione tecnologica in questo paese.

Il nostro caso per l'NRC è fondato sia sull'efficienza che su ragioni distributive. In primo luogo, l'NRC può produrre externalità positive, in particolare nel tempo, sostenendo gli investimenti nella ricerca di base che potrebbero essere commercializzati decenni dopo. In secondo luogo, può aiutare a livellare il campo di gioco ampliando l'accesso dei ricercatori sia al calcolo che ai dati, garantendo che la ricerca sull'IA sia fattibile non solo per le istituzioni accademiche più elite o per le grandi aziende tecnologiche. Data la portata della trasformazione economica che l'IA dovrebbe avviare nei prossimi decenni, la posta in gioco è potenzialmente significativa. Mentre i più grandi interessi privati come le società di tecnologia delle piattaforme e alcune istituzioni accademiche d'élite continuano a progettare, sviluppare e implementare sistemi di intelligenza artificiale che possono essere prontamente commercializzati, una storia diversa si sta svolgendo per il settore pubblico e la stragrande maggioranza delle istituzioni accademiche, che mancano accesso agli input fondamentali della ricerca sull'IA. I costi crescenti associati allo svolgimento di ricerca e sviluppo stanno esacerbando la disconnessione tra gli attuali vincitori e vinti nello spazio dell'IA.

Questo capitolo procede in tre parti. Innanzitutto, esaminiamo l'attuale panorama della ricerca sull'IA. In secondo luogo, articoliamo le tendenze mutevoli nella ricerca sull'IA e l'equilibrio del settore accademico. In terzo luogo, esplicitiamo i rischi dell'inazione federale e i vantaggi di una strategia di investimento che abbinia dati e risorse di calcolo.

CHIAVE ASPORTO

- Il governo federale svolgerà un ruolo centrale nel plasmare, coordinare e consentire lo sviluppo

dell'IA.

- Ricerca sull'intelligenza artificiale e

lo sviluppo dipende sempre più dall'accesso a computer e dati su larga scala, causando la migrazione di talenti dell'IA dal settore accademico a quello privato e limitando la gamma di voci in grado di

contribuire alla ricerca sull'IA.

- Non commerciale e di base

La ricerca sull'intelligenza artificiale è fondamentale per il

salute a lungo termine dell'ecosistema dell'innovazione.

- Un NRC che fornisce l'accesso ai dati e al calcolo contribuirà a promuovere il lungo termine salute nazionale del ecosistema AI e mitigare i rischi di un aumento delle disuguaglianze nel panorama AI della nazione.

IL PAESAGGIO DELLA RICERCA AI

Il campo della ricerca sull'IA, così come lo consideriamo in questo Libro bianco, è interpretato in senso lato. Comprende non solo gli accademici che si identificano come ricercatori nell'intelligenza artificiale o nell'apprendimento automatico, ma anche la più ampia comunità di ricercatori che utilizzano l'IA applicata nel loro lavoro, nonché coloro che ne esaminano gli impatti sulla società e sull'ambiente.

Molti credono, coerente con la legislazione che richiede la Task Force NAIIR, che l'IA avrà un impatto drammatico sulla società. Nove delle 10 aziende attualmente più grandi al mondo per capitalizzazione di mercato sono aziende tecnologiche che pongono l'IA al centro dei loro modelli di business.² I dati recenti dell'Indice AI dimostrano la quantità crescente di investimenti che le aziende di intelligenza artificiale hanno attirato. La più recente iterazione dell'Indice del 2021 descrive in dettaglio come gli investimenti privati globali nell'intelligenza artificiale siano cresciuti del 40% dal 2019 fino a un totale di 67,9 miliardi di dollari, con i soli Stati Uniti che rappresentano oltre 23,6 miliardi di dollari.³ Mentre numerose previsioni del settore privato sull'impatto economico dell'intelligenza artificiale sottolineano il potenziale dell'intelligenza artificiale di guidare una crescita economica significativa attraverso un forte aumento della produttività del lavoro, altri si preoccupano del ritmo del cambiamento strutturale nel mercato del lavoro e della dislocazione economica per i lavoratori automatizzati senza lavoro o colpiti dalla gig economy.⁴

Tali impatti sono previsti in tutti i domini. L'intelligenza artificiale ha una promessa sostanziale di trasformare l'assistenza sanitaria e la ricerca scientifica: i progressi relativi all'IA nel campo del ripiegamento delle proteine sono destinati ad accelerare notevolmente lo sviluppo di vaccini e farmaci.⁵ L'integrazione dei sistemi correlati all'IA in agricoltura può migliorare i raccolti attraverso un uso mirato dei pesticidi e del monitoraggio del suolo.⁶ E gli esperti di sicurezza nazionale hanno identificato l'IA come un motore chiave di nuove capacità di difesa,⁷ tra cui la guerra informatica e la raccolta di informazioni.

Molti paesi hanno riconosciuto l'importanza dell'IA come motore del progresso nella sicurezza economica, scientifica e nazionale, rilasciando piani nazionali che coordinano gli investimenti per il continuo progresso nell'IA.⁸ Il piano nazionale della Cina ha annunciato finanziamenti per miliardi di dollari volti a rendere il paese il paese globale leader nell'IA entro il 2030.⁹ Il governo giapponese ha collaborato con Fujitsu per costruire

il supercomputer più veloce del mondo (Fugaku).¹⁰ Compute Canada ha fornito allo stesso modo l'accesso ai computer di ricerca agli accademici di tutto il paese. La risorsa informatica di fascia alta nazionale del Regno Unito, HECToR, è stata lanciata nel 2007 al costo di 118 milioni di dollari e utilizzata da quasi 2.500 ricercatori di oltre 250 organizzazioni separate che hanno prodotto oltre 800 pubblicazioni accademiche.¹¹

Il governo degli Stati Uniti ha inizialmente presentato un altro approccio decentralizzato, fornendo supporto per lo sviluppo dell'IA attraverso le sovvenzioni della National Science Foundation e la spesa per la difesa, ma si è astenuto dal rilasciare un piano nazionale unificato per coordinare le risorse tra governo, industria privata e università.¹² La creazione di un National AI Initiative Office,¹³ l'aggiornamento della National Strategic Computing Initiative¹⁴ e la pubblicazione del rapporto finale della National Security Commission on Artificial Intelligence (NSCAI)¹⁵ hanno introdotto un approccio più completo e coordinato. Negli Stati Uniti, il modello più vicino all'NRC potrebbe essere il consorzio HPC COVID-19, che ha rapidamente fornito il calcolo di 50.000 GPU e 6,8 milioni di core per quasi 100 progetti in 43 membri del consorzio accademico, industriale e del governo federale uniti dal obiettivo comune di combattere la pandemia di COVID-19.¹⁶

Storicamente, le partnership tra governo, le università e l'industria hanno ancorato l'ecosistema dell'innovazione statunitense. Il governo federale ha svolto un ruolo fondamentale nel sovvenzionare la ricerca di base, consentendo alle università di intraprendere ricerche ad alto rischio che possono richiedere decenni per essere commercializzate. Questo approccio ha catalizzato la tecnologia radar,¹⁷ Internet, ¹⁸ e i dispositivi GPS.¹⁹ Questa storia ha ispirato la raccomandazione dell'NSCAI per nuovi investimenti nella ricerca e sviluppo dell'IA creando un'infrastruttura nazionale di ricerca sull'IA che democratizzi l'accesso al risorse che alimentano l'intelligenza artificiale. Molti politici ritengono che nel prossimo futuro saranno necessari investimenti sostanziali diversi anni per sostenere questi sforzi, mentre i ritorni su tali investimenti potrebbero potenzialmente trasformare l'economia, la società e la sicurezza nazionale americana.²⁰

A dire il vero, alcuni potrebbero contestare la teoria dell'impatto. In primo luogo, alcuni studi contestano la premessa che l'IA sarà economicamente trasformativa. Alcuni economisti lo sostengono

molte delle valutazioni ottimistiche non tengono conto di quanto possa essere limitata l'adozione dell'innovazione dell'IA a causa dell'incapacità dell'IA di modificare compiti essenziali ma difficili da migliorare.²¹ Altri criticano allo stesso modo le prove di una quarta rivoluzione industriale.²² In secondo luogo, alcuni suggeriscono che l'approvvigionamento dell'NRC può rafforzare la posizione delle grandi aziende tecnologiche di piattaforma (il che ovviamente provoca dibattiti sull'antitrust nel settore tecnologico²³), poiché l'NRC può essere difficile da avviare senza un coinvolgimento dell'hardware o fornitori di servizi cloud nel processo di approvvigionamento. In terzo luogo, alcuni sostengono che l'NRC genererebbe grandi esternalità negative sotto forma di impronte energetiche. Ad esempio, uno studio ha rilevato che la quantità di energia necessaria per addestrare GPT-3, un modello leader di elaborazione del linguaggio naturale (NLP), richiedeva l'equivalente di emissioni di gas serra di 552,1 tonnellate di anidride carbonica,²⁴ circa 35 volte le emissioni annuali di una media Americano.²⁵ Ampliare l'accesso al calcolo senza controlli appropriati può contribuire allo spreco informatico.²⁶ Infine, alcuni critici sostengono che qualsiasi progresso nell'IA sia intrinsecamente troppo rischioso per ulteriori investimenti,²⁷ dati i rischi ampiamente documentati di parzialità,²⁸ conseguenze indesiderate²⁹ e danni.³⁰

Siamo consapevoli di queste critiche e le prendiamo sul serio. Questo libro bianco procede sulla premessa operativa che anima la legislazione NRC: che sarà importante per il paese mantenere la leadership nell'IA, compreso un rigoroso interrogatorio sui suoi usi, limiti e promesse, e che ciò richiede il supporto dell'accesso al calcolo e ai dati. Gli investimenti pubblici nella ricerca sull'IA per scopi non commerciali possono aiutare ad affrontare alcuni dei problemi di danno sociale che vediamo attualmente in contesti commerciali³¹, nonché contribuire a spostare l'attenzione più ampia del campo verso la tecnologia sviluppata nell'interesse pubblico dal settore pubblico e società civile, incluso il mondo accademico. Le considerazioni precedenti, tuttavia, hanno modellato le nostre opinioni su aspetti chiave, come la strategia di investimento sequenziale, data l'incertezza del potenziale dell'IA; la seria considerazione delle infrastrutture di proprietà pubblica; le disposizioni per la revisione etica del calcolo e dell'accesso ai dati; e, cosa più importante, l'abilitazione di un'indagine accademica indipendente sui potenziali danni dei sistemi di intelligenza artificiale. L'NRC non è un'approvazione dell'adozione cieca e ingenua dell'IA su tutta la linea; è un meccanismo per garantire che una gamma più ampia di voci abbia accesso agli elementi di base della ricerca sull'IA.

L'NRC non è un'approvazione dell'adozione cieca e ingenua dell'IA su tutta la linea; è un meccanismo per garantire che una gamma più ampia di voci abbia accesso all'essenziale elementi della ricerca sull'IA.

SPOSTAMENTO FONTI DI RICERCA AI

Ora spieghiamo come e perché la ricerca sull'IA ha migrato dalla ricerca di base a lungo termine alle applicazioni commerciali a breve termine.

In primo luogo, molti progressi attuali alimentati da modelli su larga scala sono costosi da addestrare, rispetto alle dimensioni dei budget accademici tipici. Ad esempio, il costo stimato per l'addestramento dell'algorithm AlphaGo Zero della controllata Alphabet DeepMind, in grado di battere il campione mondiale umano del gioco Go, è stato di oltre 25 milioni di dollari.³² Per riferimento, il budget totale annuale 2019 per il Robotics Institute della Carnegie Mellon University, dei principali istituti di ricerca accademici della nazione, era di 90 milioni di dollari.³³ Un white paper del Bipartisan Policy Center³⁴ e del Center for a New American Security ha rilevato che il budget per l'anno fiscale 2020 per la ricerca e lo sviluppo dell'IA non per la difesa annunciato dalla Casa Bianca era di 973 milioni di dollari. Al contrario, la spesa combinata in ricerca e sviluppo nel 2018 da parte di cinque delle principali società di piattaforme tecnologiche è stata di 80 miliardi di dollari. In sintesi, le università di ricerca non possono tenere il passo con le risorse informatiche del settore privato. Questo non vuol dire che il calcolo su larga scala sia necessario per tutta la ricerca accademica sull'IA, o che la ricerca accademica sia in concorrenza con la ricerca industriale, ma illustra perché alcuni settori della ricerca sull'IA non sono più accessibili al ricerca

In secondo luogo, il divario tra mondo accademico e industria maschera disparità significative tra le istituzioni accademiche. Utilizzando il QS World University Rankings dal 2012, le aziende tecnologiche Fortune 500 e le 50 migliori università hanno pubblicato cinque volte più documenti all'anno per conferenza sull'IA rispetto alle università classificate tra 200 e 500.³⁵ Collaborano anche aziende private

sei volte di più con le prime 50 università che con quelle classificate tra 301 e 500.³⁶ Questo divario di calcolo interno tra le università pone sfide significative per chi è al tavolo.

In terzo luogo, la ricerca di base sull'IA ha perso capitale umano.³⁷ Quando questo è combinato con un accesso ridotto a

calcolo e dati nell'accademia, la prospettiva di condurre ricerca di base nelle università diventa meno allettante. I migliori talenti nel campo dell'intelligenza artificiale ora ricevono stipendi del settore privato di gran lunga superiori a quelli accademici.³⁸ Il

la partenza della facoltà di intelligenza artificiale dalle università americane ha portato a quello che alcuni analisti hanno soprannominato la fuga di cervelli di intelligenza artificiale: mentre i dottorati di ricerca in intelligenza artificiale nel 2011 avevano all'incirca le stesse probabilità di entrare nell'industria rispetto al mondo accademico, due terzi dei dottorati di ricerca in intelligenza artificiale ora entrano nell'industria e meno di un quarto entra nel mondo accademico.³⁹ Uno studio suggerisce che l'abbandono della facoltà di intelligenza artificiale ha anche un effetto negativo sulla formazione di startup da parte degli studenti.⁴⁰

Mentre i dottorati di ricerca in intelligenza artificiale nel 2011 lo erano

all'incirca con la stessa probabilità di entrare

nell'industria come nel mondo accademico, due terzi dei

dottorati di ricerca in intelligenza artificiale ora

entrano nell'industria e meno di un quarto nel

mondo accademico.

In quarto luogo, mentre la ricerca sull'IA su larga scala migra verso l'industria, il focus della ricerca si sposta inevitabilmente. Mentre i ricercatori accademici nel campo dell'intelligenza artificiale possono non avere accesso al volume di dati necessari per addestrare i modelli di intelligenza artificiale,⁴¹ le aziende di grandi piattaforme hanno accesso a vasti set di dati, compresi quelli relativi o creati dai loro clienti. Questa divisione dei dati a sua volta distorce la ricerca sull'intelligenza artificiale verso applicazioni incentrate sul profitto privato, piuttosto che sul beneficio pubblico. L'NRC può svolgere un ruolo chiave nello sbloccare l'accesso ai dati del settore pubblico, il che può aiutare a riorientare l'attenzione della ricerca sull'IA lontano dai set di dati del settore privato.⁴⁴

Lo svuotamento della capacità di IA accademica può essere visto nell'analisi di OpenAI della relazione tra il calcolo e 15 "scoperte" relativamente note nell'IA tra il 2012 e il 2018.⁴⁵ Sebbene l'analisi avesse lo scopo di enfatizzare il ruolo della potenza di calcolo, illustra anche un divario emergente tra il settore privato e i contributi accademici nel tempo. Dei 15 sviluppi esaminati, 11 sono stati realizzati da aziende private mentre solo quattro provenivano da istituzioni accademiche. Inoltre, questo squilibrio aumenta nel tempo: sebbene la ricerca del settore privato abbia continuato ad accelerare dal 2012, la produzione accademica è rimasta stagnante. L'ultima delle principali scoperte ad alta intensità di calcolo nell'analisi di OpenAI derivanti dal mondo accademico è stata il rilascio nel 2014 di Oxford del suo programma di riconoscimento delle immagini VGG; Il lavoro della NYU sulle reti neurali convoluzionali risale al 2013. Dal 2015 al 2018, tutte le otto scoperte incluse nell'analisi di OpenAI provenivano da società private. Nel loro insieme, questo porta gli osservatori a sostenere che i ricercatori accademici sono sempre più incapaci di competere alla frontiera della ricerca sull'IA. Con meno scoperte accademiche ad alta intensità di calcolo, le innovazioni dell'IA si sono concentrate su interessi privati (ad esempio, la pubblicità online) piuttosto che su vantaggi non commerciali a lungo termine. A dire il vero, il settore privato è stato, ovviamente, centrale per la ricerca sull'IA, ma la preoccupazione riguarda l'equilibrio a lungo termine dell'ecosistema dell'innovazione dell'IA.

SCOPING INTERVENTO FEDERALE IN DATI E CALCOLO

Come possiamo raggiungere un approccio più equilibrato nei confronti della ricerca e dello sviluppo? Per prima cosa consideriamo i rischi dell'inazione federale e discutiamo alcuni dei vantaggi unici dell'indirizzare dati e calcolare insieme.

Rischi di inerzia federale

I rischi dell'inazione federale sono duplici. In primo luogo, la ricerca di base sull'IA che fino ad oggi ha spianato la strada ai progressi nell'IA e l'apprendimento automatico rallenterà. Secondo uno studio recente, circa l'82% degli algoritmi utilizzati oggi proviene da gruppi no profit e università sostenuti dalla spesa pubblica.⁴⁷ Anche quando la ricerca industriale ha successo, è tipicamente focalizzata sul prodotto o incrementale, più difficile da riprodurre e potrebbe non essere pubblicato o open-source. Un caso interessante risiede nelle recenti scoperte nel ripiegamento delle proteine. Alla fine del 2020, DeepMind, sussidiaria di Alphabet, ha annunciato di aver sviluppato un programma chiamato AlphaFold, un sistema guidato dall'intelligenza artificiale in grado di prevedere con precisione la struttura di un vasto numero di proteine, utilizzando solo la sequenza di nucleotidi contenuta nel suo DNA. Sia per preoccupazione per la privatizzazione o per accelerare l'adozione di sistemi correlati, un consorzio di accademici, guidato da scienziati dell'Università di Washington, ha sviluppato un concorrente open source chiamato RoseTTaFold.⁴⁸ DeepMind ha reso AlphaFold disponibile a un vasto pubblico, ma le preoccupazioni illustrano i rischi della scienza posta dalla ricerca sull'intelligenza artificiale esclusivamente privata, che ricorda la corsa al sequenziamento del genoma umano, in cui l'investimento pubblico nel progetto genoma umano ha anticipato le preoccupazioni per un'azienda privata che brevettava il genoma umano.⁴⁹

In secondo luogo, l'inerzia federale potrebbe aumentare notevolmente disuguaglianze nel panorama dell'IA. Senza un maggiore accesso all'informatica, all'istruzione e alla formazione, gran parte dell'economia potrebbe non essere in grado di adattarsi, che si tratti di servizi finanziari, assistenza sanitaria, istruzione o governo. Diversificare la gamma della ricerca sull'IA può anche promuovere il progresso e la produttività. Uno studio suggerisce che la diversità delle traiettorie di ricerca sull'IA, ovvero le domande, gli argomenti e i problemi specifici che i ricercatori scelgono

indagare—è diventato più limitato negli ultimi anni e che la ricerca sull'IA nel settore privato è meno diversificata rispetto alla ricerca accademica.⁵⁰ Gruppi accademici più piccoli con una minore collaborazione del settore privato sembrano rafforzare la diversità della ricerca sull'IA.⁵¹ Dal punto di vista delle vie di ricerca sottosviluppate, come l'etica e la responsabilità nell'IA, l'ampliamento della gamma di argomenti e metodi di ricerca nel campo aumenta la probabilità di trovare scoperte che rendano possibili ulteriori progressi a lungo termine.⁵² Prove recenti suggeriscono che tra il 2005 e il 2017, solo cinque aree metropolitane in gli Stati Uniti hanno rappresentato il 90 per cento della crescita dei posti di lavoro nel settore dell'innovazione.⁵³ Secondo l'economista di Stanford Erik Brynjolfsson, il probabile impatto della concentrazione geografica è che “ci sono un sacco di persone — centinaia di milioni negli Stati Uniti e miliardi in tutto il mondo — che potrebbero essere innovativi e chi non lo è perché non ha accesso alle competenze informatiche di base, o alle infrastrutture, o al capitale, o persino alla cultura e agli incentivi per farlo . . .⁵⁵ Ampliare l'insieme delle voci che possono interrogare tali sistemi sarà fondamentale per un futuro inclusivo ed equo.

In sintesi, investimenti federali nell'infrastruttura pubblica di intelligenza artificiale può promuovere una distribuzione più equa della partecipazione e dei guadagni per l'innovazione dell'IA in generale, rafforzare la competitività degli Stati Uniti e sostenere la ricerca fondamentale nelle applicazioni non commerciali e del settore pubblico.

Capitolo 2:

Idoneità, allocazione e infrastruttura per l'informatica

Questo capitolo tratta l'idoneità, l'allocazione delle risorse e l'elaborazione infrastruttura per l'NRC: chi dovrebbe avere accesso a cosa e come?

In primo luogo, quando si determina chi dovrebbe ottenere l'accesso, è fondamentale tenere a mente gli obiettivi generali dell'NRC. Come discusso nel capitolo 1, c'è un grande divario di risorse nel mondo accademico rispetto all'industria privata. Nell'interesse di sostenere la ricerca di base e democratizzare il settore, questa sezione si concentrerà sull'identificazione di un gruppo target per l'ammissibilità. Come articoliamo di seguito, ci asteniamo dal prendere in considerazione l'espansione a un insieme più ampio di parti commerciali e non accademiche a causa dell'attenzione dell'NRC sulla ricerca scientifica fondamentale a lungo termine. Uno degli approcci più ristretti sarebbe un modello di facoltà specialistica che si rivolgerebbe ai ricercatori impegnati nel lavoro di base sull'IA. Tuttavia, le difficoltà nel definire l'IA e i domini in rapida espansione in cui viene applicata l'IA rendono questo modello troppo vincolato per realizzare il pieno impatto dell'NRC. Raccomandiamo invece di tracciare il criterio più comune per il finanziamento della ricerca federale e sostenere che l'idoneità dipenda dallo status di "Principal Investigator" (PI) presso le università statunitensi.¹ Uno dei compromessi è che i PI possono essere meno diversificati rispetto a un segmento più ampio di ricercatori,² quindi un'espansione a lungo termine potrebbe considerare di andare oltre questo gruppo. Sebbene l'NRC miri a formare la prossima generazione di ricercatori di intelligenza artificiale, avvertiamo che un'espansione immediata a tutti gli studenti laureati e universitari porrebbe notevoli sfide nel ridimensionamento. Pertanto, raccomandiamo che gli studenti ottengano l'accesso principalmente partecipando alla ricerca sull'IA sponsorizzata dalla facoltà, invece dell'accesso generale degli studenti, e che ottengano una formazione attraverso la creazione di programmi educativi.

In secondo luogo, discutiamo tre modelli per l'allocazione del credito informatico: sviluppo di un nuovo processo di sovvenzione, delega delle sovvenzioni per blocchi di calcolo alle università per l'allocazione interna tra i docenti o accesso universale. Ciascuno di questi modelli compromette la facilità di amministrazione con la personalizzazione per specifici obiettivi NRC. Raccomandiamo un approccio utilizzato da altri cloud di ricerca nazionali, vale a dire un approccio ibrido di accesso predefinito universale per la maggior parte dei ricercatori, con un processo di sovvenzione per il calcolo in eccesso oltre l'allocazione predefinita. Un tale approccio manterrebbe bassi i costi amministrativi per la stragrande maggioranza dei ricercatori, consentendo al contempo l'adattamento attraverso un processo di sovvenzione competitivo per gli utenti più bisognosi.

CHIAVE ASPORTO

- Idoneità del ricercatore per l'accesso NRC dovrebbe iniziare con lo status di "Principal Investigator" negli Stati Uniti università.
- L'NRC dovrebbe adottare un approccio ibrido di accesso predefinito universale per la maggior parte dei ricercatori e un processo di sovvenzione quando le richieste di calcolo o i dati superano i livelli di base.
- L'NRC dovrebbe adottare una doppia strategia di investimento sviluppando programmi per espandere l'accesso ai servizi cloud esistenti e pilotare la capacità di fornire servizi di proprietà pubblica risorse.

In terzo luogo, consideriamo la decisione "make-or-buy" per il CNR. Un'opzione sarebbe che l'NRC fornisse borse di ricerca per l'uso di servizi cloud commerciali su cui molti ricercatori fanno già affidamento (la decisione di "acquisto").

In alternativa, l'NRC potrebbe creare e fornire l'accesso a un cluster informatico pubblico ad alte prestazioni (la decisione "make"). È risaputo che, basandosi esclusivamente sui costi dell'hardware, è più conveniente possedere un'infrastruttura quando la domanda di elaborazione è quasi continua. D'altra parte, i provider di cloud commerciali esistenti hanno sviluppato stack software altamente utilizzabili che i ricercatori di intelligenza artificiale hanno ampiamente adottato. I fornitori di servizi cloud commerciali offrono un modo per espandere rapidamente la capacità. Raccomandiamo quindi una doppia strategia di investimento per (a) lanciare rapidamente l'NRC sovvenzionando e negoziando il cloud computing per i ricercatori di intelligenza artificiale con i fornitori esistenti, espandendo le iniziative esistenti come il progetto CloudBank della National Science Foundation; e (b) investire in un progetto pilota

per le infrastrutture pubbliche per valutare la capacità di fornire risorse simili nel lungo periodo. Tale infrastruttura di proprietà pubblica sarebbe probabilmente costruita sotto contratto o sovvenzione, ma potrebbe essere gestita in modo molto simile ai laboratori nazionali che possiedono sofisticate strutture di supercalcolo, come nel caso di altre risorse di ricerca nazionali (ad esempio, Compute Canada, il giapponese Fugaku).

Le nostre raccomandazioni si basano su una serie di studi di casi presentati in questo capitolo e nel resto del Libro bianco.

La tabella 1 riassume il confronto tra i modelli esistenti e le tre decisioni chiave di progettazione. All'inizio, notiamo che poche iniziative esistenti hanno tentato di fornire potenza di calcolo alla scala dell'NRC. Allo stesso tempo, consideriamo l'NRC complementare alle aree più tradizionali dell'informatica scientifica.³

	ELEGGIBILITÀ			ASSEGNAZIONE				PROPRIETÀ	
Esistente Programma	PI Soltanto	Qualunque Facoltà	Studenti	Esistente Concessione Processi	Università Allocazione	Nuovo Processi	Predefinito Accesso con livelli	Privato	Pubblico
CloudBank	X		X	X					X
Stanford HAI-AWS Programma nuvola		X			X			X	
Stanford Sherlock Grappolo	X						X	X	
Google Co		X	X				X	X	
Calcolare Canada	X						X		X
Fugaku		X				X			X
XSEDE	X	X					X		X
INCITE	X					X			X

Tabella 1: Principali differenze di progettazione tra casi di studio informatici. "Altra facoltà" indica un insieme di ammissibilità per la facoltà diversa dallo stato PI (ad esempio, che richiede l'affiliazione a Stanford per il cluster Sherlock) e "nuovo processo" viene utilizzato per indicare la creazione di qualsiasi processo diverso da quelli attualmente elencati (ad esempio, Fugaku è attualmente sollecitazione di proposte alle strutture di ricerca).

Eleggibilità

Il primo compito è identificare quali ricercatori dovrebbero essere ammessi al CNR. Il capitolo 1 ha discusso la necessità di sostenere l'innovazione dell'IA nelle università. Pertanto, questa sezione esaminerà l'ammissibilità all'interno del mondo accademico analizzando i compromessi accesso-risorse in linea con gli obiettivi NRC.

All'inizio, notiamo che non analizziamo l'ammissibilità in profondità al di là dei ricercatori accademici. La legislazione che costituisce la task force dell'NRC contempla specificatamente "l'accesso alle risorse informatiche per i ricercatori di tutto il paese".⁴ L'NRC è definito come "un sistema che fornisce a ricercatori e studenti di

discipline con accesso a risorse informatiche."⁵ L'interpretazione più naturale di questo linguaggio suggerisce un focus centrale sulla ricerca scientifica e accademica.⁶

Introduzione dell'accesso commerciale all'NRC, in particolare per le imprese con risorse insufficienti come le piccole imprese e le startup, potrebbero benissimo avvantaggiare l'ecosistema dell'innovazione statunitense. Ma le sfide per incorporare l'accesso commerciale all'NRC sono enormemente complesse. In primo luogo, includere gli sviluppatori di software presso le startup come "ricercatori" ai sensi dell'NDAA solleverebbe un'ampia gamma di questioni di confine che l'NRC potrebbe non essere adeguatamente attrezzato per giudicare. Secondo la Small Business Administration (SBA), ci sono oltre 31 milioni di piccole imprese negli Stati Uniti.⁷ Ogni anno vengono aperte oltre 627.000 imprese.⁸ Tutte queste imprese dovrebbero essere idonee a calcolare sulla NRC? Come si potrebbe evitare l'ammissibilità al gioco (ad es. consociate strategiche/spin-off)? E come farebbe questo a far avanzare la missione scientifica dell'NRC? In secondo luogo, sebbene potenzialmente preziosa, è meno chiaro in che modo l'inclusione di startup e piccole imprese soddisfi la teoria dell'impatto dell'NRC. Come attualmente interpretato, la preoccupazione che anima l'NRC risiede nell'importanza della ricerca fondamentale a lungo termine e non commerciale che può garantire la leadership dell'IA per i decenni a venire. La commercializzazione non è l'elemento

dell'ecosistema dell'innovazione AI che affronta le sfide strutturali articolate nel capitolo 1. Infine, ridimensionare l'NRC per consentire un accesso commerciale significativo porrebbe serie sfide pratiche. Perché anche la Task Force deve farlo

Considerata la fattibilità dell'NRC, non abbiamo considerato in profondità una concezione che estenderebbe il termine "ricercatore" per comprendere ampie porzioni del settore privato commerciale. L'espansione a organizzazioni non accademiche e senza scopo di lucro può essere una considerazione più ragionevole, in quanto l'obiettivo di alcune entità (ad esempio, giornalismo investigativo senza scopo di lucro, organizzazioni della società civile) potrebbe essere più vicino al nucleo della missione dell'NRC di responsabilizzare a lungo termine ricerca benefica che attualmente non può verificarsi.⁹ A lungo termine, l'NRC dovrebbe considerare i compromessi di una tale espansione.

Anche se l'NRC adotta un modello informatico più ampio lungo la strada, riteniamo che concentrarsi sui ricercatori accademici sia un importante punto di partenza in quanto illumina alcune delle principali considerazioni operative per NRC accesso.

MODELLO DI FACOLTÀ DI SPECIALITÀ

Uno degli approcci più ristretti all'ammissibilità dell'NRC sarebbe quello di limitarlo ai docenti impegnati nella ricerca sull'IA. In base a questo approccio, i responsabili politici indirizzerebbero le risorse di calcolo esclusivamente verso i docenti che lavorano su progetti di intelligenza artificiale identificabili, che spesso richiedono grandi quantità di potenza di calcolo. Un vantaggio di questo approccio è che la familiarità dei ricercatori con l'infrastruttura significherebbe probabilmente che meno fondi sarebbero destinati al servizio cloud formazione per utenti inesperti.

Eppure l'insieme di docenti di intelligenza artificiale di base autoidentificati sono pochi e concentrato in un piccolo numero di università, che hanno già maggiori probabilità di ottenere l'accesso all'informatica su larga scala. Limitare l'accesso alla facoltà di intelligenza artificiale di base minerebbe quindi la missione di democratizzare la ricerca sull'IA. Inoltre, l'applicazione dell'IA si sta espandendo rapidamente in tutti i domini. La ricerca interdisciplinare che dispiega l'IA in nuovi domini sarà vitale per mantenere la leadership americana nell'IA, nonché per animare le domande di ricerca di base. Limitare l'ammissibilità alla facoltà di IA di base (comunque definita) potrebbe compromettere la capacità dei ricercatori di tutte le discipline accademiche (ad esempio, nelle scienze fisiche, nelle scienze sociali e umanistiche) di contribuire alla realizzazione del pieno potenziale dell'IA.

MODELLO DI FACOLTA' GENERALE

Un punto di partenza più naturale per l'ammissibilità NRC è con Principal Investigators (PIs) presso college e università statunitensi, il criterio più comunemente utilizzato per le sovvenzioni federali. I requisiti per lo status di PI sono stabiliti dalle singole università e includono un'ampia gamma di ricercatori certificati dalla loro università come qualificati per condurre grandi progetti di ricerca. È soggetto ai processi di formazione e certificazione del proprio istituto, che a loro volta chiariscono le responsabilità di un ricercatore in merito alla gestione e all'esecuzione delle proprie proposte di ricerca. I programmi esistenti per l'allocazione della potenza di calcolo in genere stabiliscono l'ammissibilità in base allo stato PI in quanto garantisce che il ricercatore disponga dell'infrastruttura

realizzare un progetto di ricerca su vasta scala. CloudBank, un programma NSF che distribuisce fondi per risorse di cloud computing commerciale, assegna sovvenzioni a PI, che possono distribuire fondi ad altri ricercatori e studenti sul

project.¹¹ Compute Canada consente a tutti i docenti a cui è stato concesso lo status di PI dalla propria università di ricevere automaticamente un numero prestabilito di crediti informatici e richiedere ulteriori crediti se necessario. Il PI può quindi sponsorizzare altri per accedere al credito.¹²

Riconosciamo che lo stato PI non include tutti i ricercatori affiliati all'università. Nel 2013, degli oltre 200.000 ricercatori accademici autoidentificatisi, poco meno di 60.000 erano occupati in un ruolo diverso dal corpo docente a tempo pieno, una posizione che potrebbe non essere qualificata per lo status di PI.¹³ Dal 1973 al 2013, la percentuale di ricercatori a tempo pieno la facoltà a tempo tra i titolari di dottorato in ingegneria è diminuita del 2%, mentre la percentuale di "altri" lavori accademici (inclusi gli associati alla ricerca) è aumentata del 12%. I ricercatori su un progetto e la capacità amministrativa pesa fortemente a favore della coerenza con gli attuali criteri di ammissibilità delle sovvenzioni.

STUDENTI

Gli studenti laureati e laureandi dovrebbero poter accedere all'NRC? Una delle sfide principali

qui sta nella scala e nell'amministrabilità. Una stima è che ci sono quasi 20 milioni di studenti universitari negli Stati Uniti.¹⁵ In secondo luogo, l'idoneità orientata ai PI non preclude agli studenti universitari l'accesso alle risorse per intraprendere ricerche sull'IA sotto la direzione dei PI. IL

Il modello Compute Canada, ad esempio, limita l'idoneità ai docenti, ma consente ai docenti di sponsorizzare collaboratori, incluso qualsiasi studente ricercatore. Un modello di accesso per l'NRC che consente ai PI di sponsorizzare gli studenti fornisce ulteriori opportunità di ricerca e formazione per gli studenti. In terzo luogo, numerosi servizi cloud esistenti forniscono già un accesso limitato ai crediti di calcolo per scopi didattici.

Google Colaboratory, ad esempio, fornisce accesso gratuito, ma non garantito in modo affidabile, ai servizi cloud.¹⁶ Amazon Web Services offre gratuitamente fino a 35 dollari di crediti AWS a tutti i docenti universitari e agli studenti. Nonostante le risorse esistenti, gli studenti potrebbero aver bisogno di più risorse. La consociata di Google e la community online Kaggle, ad esempio, fornisce gratuitamente 30 ore di accesso alla GPU a settimana e ha rilevato che il 15% degli utenti ha superato il limite.¹⁷

Sebbene l'esatta portata delle esigenze di potenza di calcolo degli studenti non sia chiara, consigliamo di finanziare una risorsa educativa una volta che il ricercatore ha bisogno e le risorse sono limitate sono stati misurati. Attualmente, CloudBank della NSF sta testando una risorsa per la comunità e l'istruzione per destinare un piccolo set di crediti a scopi educativi.¹⁸ Questa risorsa consente a un professore universitario di richiedere un numero limitato di crediti per corsi degli studenti o ricerca.

Indipendentemente dal modello di ammissibilità adottato dall'NRC, ci sarà anche un bisogno significativo di personale di supporto, documentazione per la formazione e materiali didattici in modo che i ricercatori possano utilizzare efficacemente le risorse di calcolo e dati (vedi Appendice D). Il motivo per cui alcuni studenti e ricercatori potrebbero non sfruttare tutti i crediti cloud disponibili potrebbe, ad esempio, derivare dalla difficoltà nell'utilizzo delle piattaforme cloud. Se l'NRC serve accademici di una vasta gamma di discipline, questa questione del capitale umano sarà particolarmente rilevante per servire diversi modelli di ricerca. Un solido programma di formazione per gli utenti dell'NRC garantirà la facilità d'uso e incoraggerà un utilizzo appropriato del cloud.

Assegnazione delle risorse Modelli

Consideriamo ora tre modelli di allocazione delle risorse: (1) un nuovo processo di sovvenzione; (2) assegnazione di sovvenzioni in blocco alle università; e (3) accesso universale, ma potenzialmente a più livelli.

PROCESSO DI CONCESSIONE NRC

Definizione di un nuovo processo di concessione per l'accesso al calcolo avrebbe un vantaggio principale. Il programma potrebbe essere creato appositamente per lo scopo della ricerca sull'IA, con revisori che hanno familiarità con concetti, pratiche e tendenze dell'IA. Un tale processo potrebbe quindi consentire migliori decisioni di allocazione e fornire all'NRC un maggiore controllo sui propri investimenti.

Detto questo, stabilire un processo di revisione tra pari per tutte le domande richiederebbe un uso intensivo delle risorse, richiedendo l'istituzione di un programma di amministrazione delle sovvenzioni simile a quelli della National Science Foundation (NSF) o dei National Institutes of Health (NIH). Ad esempio, per implementare la revisione tra pari richiesta per il processo di revisione del merito, la NSF ha bisogno ogni anno di una comunità sufficientemente ampia da condurre quasi 240.000 revisioni all'anno.¹⁹ Poiché la portata prevista è ampia, siamo consapevoli di aggiungere un onere di servizio significativo per i docenti esperti in AI per ogni applicazione per l'accesso al computer. La revisione tra pari per l'accesso al calcolo richiederebbe un sovraccarico significativo e ritardi nell'allocazione del calcolo.

ACCESSO UNIVERSITARIO

Per ridurre i costi amministrativi, uno schema alternativo sarebbe quello di assegnare crediti alle università sulla base del numero di ricercatori ammissibili. L'NRC potrebbe allocare risorse alle università come sovvenzioni in blocco e, a sua volta, fare affidamento sull'università per distribuire l'accesso al computer. (Ad esempio, l'NRC potrebbe acquistare quantità significative di elaborazione dai fornitori di servizi cloud, creare crediti virtuali convertibili in risorse cloud appropriate e delegare l'allocazione alle università).

competenza per la revisione e la distribuzione delle risorse. Tuttavia, porterebbe a un processo altamente decentralizzato, fornendo poca supervisione per comprendere la distribuzione dell'utilizzo e dando all'NRC scarso controllo sull'allocazione delle risorse. Sebbene non raccomandiamo questo percorso come principale schema di allocazione, riteniamo che una certa allocazione ai team di supporto IT universitari possa essere giustificata per supportare i ricercatori nell'utilizzo dell'NRC. Il programma "Campus Champions" di XSEDE, ad esempio, fornisce ai dipendenti dell'università l'accesso al sistema per supportare la transizione computazionale.²⁰

ACCESSO UNIVERSALE

L'ultimo modello potenziale fornirebbe l'accesso universale al calcolo di livello base a tutti i PI idonei. Il modello più vicino è il cloud di ricerca nazionale di Compute Canada, che fornisce accesso al calcolo di livello base a tutte le facoltà in Canada. Ciò ridurrebbe in modo significativo il carico amministrativo, sia per un istituto che gestisce il processo di revisione, sia per gli accademici che richiedono l'accesso all'NRC. Lo svantaggio principale è che il calcolo di livello base potrebbe non essere sufficiente per esigenze specializzate.

Si consiglia di combinare un modello di base universale con un processo di concessione per le esigenze di calcolo oltre l'accesso a livello di base. La ridotta complessità nell'amministrazione di un modello di calcolo di accesso di base universale lo rende un'opzione interessante per l'NRC nell'allocazione delle risorse di calcolo, in particolare rispetto all'obiettivo dell'NRC di aprire l'accesso alle risorse di calcolo.²¹ XSEDE, ad esempio, utilizza un modello simile di semplificato "Startup Allocations" (rilasciato per periodi di un anno, in genere entro due settimane dalla domanda) e "Research Allocations" per richieste di calcolo più significative. Compute Canada fornisce l'accesso al 15% dei PI a una maggiore capacità di elaborazione sulla base di una competizione di merito. Una domanda critica sarà, ovviamente, il livello di calcolo di base che determinerà i costi complessivi, i requisiti di spazio fisico e simili. Per valutare questo aspetto, raccomandiamo uno studio approfondito delle esigenze informatiche previste, basato sui centri di calcolo accademici esistenti.²²

Il processo di concessione per il calcolo aggiuntivo potrebbe assumere più forme; ad esempio, mentre si potrebbe consentire ai singoli PI di presentare domanda direttamente all'NRC per eccesso

compute, l'NRC potrebbe anche allocare "blocchi" di risorse a livello universitario e consentire alle università di supervisionare la loro amministrazione. In ogni caso, a causa dell'entità di tali richieste, le revisioni delle sovvenzioni dovrebbero essere condotte in base al merito e gestite da una combinazione di personale NRC e un comitato consultivo esterno di docenti universitari. Nel 2021, Compute Canada, ad esempio, ha completato la revisione di 650 proposte di ricerca in circa

revisori volontari delle istituzioni accademiche canadesi per valutare il valore scientifico della proposta.²³ Al fine di evitare conflitti di interesse, raccomandiamo vivamente di non far partecipare docenti o consulenti del settore privato che abbiano conflitti di interesse con fornitori che forniscono servizi all'NRC. Idealmente, la revisione della proposta dovrebbe essere indipendente, cieca e basata su meriti scientifici nella misura del possibile.

CASO DI STUDIO: CLOUDBANK

Nel 2018, il Directorate for Computer and Information Science and Engineering (CISE) della National Science Foundation (NSF) ha creato la Cloud Access Solicitation per fornire finanziamenti per le attività di ricerca relative all'IA. cloud, CloudBank è un caso di studio interessante per esplorare i modelli di allocazione delle risorse. Accessibile tramite un portale, CloudBank aiuta i ricercatori a utilizzare appieno le risorse cloud facilitando il processo di "gestione dei costi, traduzione e aggiornamento degli ambienti informatici nel cloud e apprendimento delle tecnologie basate su cloud".²⁴

CloudBank è un progetto di collaborazione stabilito tramite un accordo di cooperazione NSF con il San Diego Supercomputer Center (SDSC) e la divisione Information Technology Services dell'UC San Diego, l'Istituto di eScience dell'Università di Washington e la divisione di scienza dei dati e informazioni dell'UC Berkeley.²⁵ Ciascuno di queste istituzioni gestisce un'area, in base al suo vantaggio comparativo.²⁶ Ad esempio, SDSC è responsabile della costruzione del portale online e UC San Diego è responsabile della gestione degli account degli utenti.²⁷

CloudBank mira anche a ridurre il costo del cloud computing: utilizza sia gli sconti in corso con i fornitori di cloud dell'Università della California sia gli sconti che derivano dall'acquisto di cloud in blocco dalla società di consulenza per l'approvvigionamento di cloud Strategic Blue, che collabora regolarmente con artisti del calibro di AWS, Microsoft e Google.²⁸ Inoltre, non vi sono costi generali associati alle allocazioni cloud tramite CloudBank, poiché i termini dell'accordo cooperativo NSF vietano i costi indiretti.²⁹ Con questi meccanismi di risparmio sui costi, i ricercatori possono permettersi maggiori capacità di calcolo da una varietà di importanti fornitori di cloud.

Richiedendo l'uso di CloudBank durante la loro candidatura ai progetti NSF selezionati,³⁰ i ricercatori possono ottenere l'accesso non solo a varie risorse hardware avanzate, ma anche a una varietà di servizi per rendere il processo più supportato e monitorato.³¹ CloudBank offre anche alla comunità di ricerca l'accesso dei membri alle informazioni relative all'istruzione e alla formazione.³²

PUNTI CHIAVE

- **Integrato nel processo di sovvenzione esistente:**
I ricercatori ammissibili per alcuni esistenti
Le sovvenzioni NSF possono semplicemente richiedere l'accesso a CloudBank attraverso la stessa domanda di contributo.
- **Singolo punto di accesso per l'accesso al computer:**
La banca delle nuvole
portale fornisce un unico punto di accesso per i ricercatori a accedere ai fondi per utilizzare su qualsiasi nuvola commerciale
fornitore che preferiscono.
- **Riduzione dei costi:** no
i costi generali sono associato all'utilizzo di CloudBank.
- **Accesso studenti:**
Sono fissati fondi limitati da parte per borse di studio a studenti e classi.

Informatica Infrastruttura

Gli ambienti di cloud computing collegano i dispositivi informatici locali come i computer desktop a server di grandi dimensioni, generalmente distribuiti geograficamente, contenenti hardware fisico. Questo hardware, a sua volta, è responsabile dell'archiviazione dei dati e dell'esecuzione del calcolo sulle reti di computer, il tutto mediato da una raccolta di servizi software. Questo modello centralizza il solito gestione operativa per coloro che utilizzano la rete e fornisce unità regolabili di calcolo e memorizzazione dei dati per consentire le fluttuazioni della domanda. Gli utenti interagiscono con il cloud avviando connessioni virtuali al server, istanze cloud, ed eseguendo processi containerizzati da remoto. Queste operazioni sono gestite dal cloud e disponibili per il monitoraggio tramite dashboard. Il cloud computing può essere servito tramite cluster locali, tramite fornitori esterni o una combinazione di questi, e accessibile tramite reti con sicurezza e connettività variabili, da regioni accessibili a Internet a regioni con air gap.

L'infrastruttura per l'NRC potrebbe essere sviluppata con due approcci generali: (1) l'NRC potrebbe utilizzare piattaforme cloud commerciali come spina dorsale dell'infrastruttura; o (2) il governo federale potrebbe assumere un appaltatore per costruire una struttura pubblica di calcolo ad alte prestazioni (HPC) specificatamente per l'NRC. Questa sezione affronta alcuni vantaggi e svantaggi di entrambi. (Forniamo un confronto dei costi stimati di questi due approcci nell'Appendice A.) I due approcci discussi qui non si escludono a vicenda e, in ultima analisi, raccomandiamo una strategia di investimento ibrida. A breve termine, l'NRC dovrebbe ampliare i programmi di credito cloud (simile al programma CloudBank di NSF) per fornire sia un accesso di livello base semplificato sia una revisione del merito per le applicazioni che vanno oltre l'accesso di livello base. A lungo termine, l'NRC dovrebbe investire in un progetto pilota per sviluppare un'infrastruttura informatica pubblica. Anche con l'infrastruttura pubblica, sarà fondamentale soddisfare la "domanda esplosiva" (per espandere le risorse quando la domanda di elaborazione raggiunge il picco). Il successo degli investimenti iniziali dovrebbe guidare il modello prospettico in merito all'opportunità di affidarsi a infrastrutture di proprietà pubblica o privata a lungo termine. Notiamo che per scalare correttamente a entrambe le risorse

richiederà la costruzione di capacità istituzionali presso le istituzioni accademiche.

NUVOLA COMMERCIALE

Il più grande vantaggio di utilizzare il cloud commerciale servizi per l'NRC è che esiste già un'infrastruttura significativa.³³ Secondo questo modello, l'NRC sovvenzionerebbe semplicemente crediti per l'utilizzo di servizi cloud commerciali (simile al programma CloudBank di NSF). Pertanto, invece di passare anni a costruire nuove risorse informatiche, i responsabili politici potrebbero lanciare l'NRC subito dopo aver determinato i dettagli amministrativi del programma. (Notiamo, tuttavia, che potrebbero esserci ancora significative carenze di GPU nel breve periodo; con la scala contemplata dell'NRC, sarebbe necessario costruire un'infrastruttura significativa.) Poiché molti ricercatori utilizzano già servizi cloud commerciali per la loro ricerca sull'IA, il la transizione al programma NRC potrebbe essere relativamente agevole. Inoltre, le piattaforme cloud commerciali offrono all'NRC una maggiore flessibilità per modificare le dimensioni e l'ambito del programma. Le piattaforme cloud commerciali fanno pagare la quantità di elaborazione effettivamente utilizzata.³⁴ Pertanto, la dimensione dell'NRC potrebbe espandersi o ridursi in linea con il mutare della domanda. Al contrario, un sistema HPC dedicato avrebbe una determinata quantità di hardware che costa lo stesso, indipendentemente dall'efficacia con cui viene utilizzato.

Lavorare direttamente anche con provider di cloud commerciali offre diversi vantaggi per l'NRC. Il mercato dei servizi cloud commerciali è altamente competitivo e presenta numerosi provider in grado di soddisfare le esigenze dell'NRC. L'NRC avrebbe la possibilità di utilizzare uno o più fornitori. Se si sceglie di utilizzare un solo fornitore, il potere contrattuale del governo potrebbe essere al massimo nel contribuire a ridurre i prezzi per l'NRC. In alternativa, l'utilizzo di più fornitori offre all'NRC una maggiore flessibilità nei servizi e nell'hardware disponibili. In ogni caso, i responsabili politici avrebbero l'opportunità di negoziare contratti e prezzi con fornitori di cloud commerciali ogni pochi anni, il che sarà fondamentale per il contenimento dei costi.³⁵ L'NRC

inoltre, non essere vincolato all'utilizzo dello stesso fornitore o gruppo di fornitori per la durata del programma. Piuttosto, il personale dell'NRC potrebbe rivalutare quale infrastruttura del fornitore di cloud commerciale soddisferebbe al meglio le esigenze dell'NRC al momento inizio di ogni nuovo contratto.

CASO DI STUDIO: XSEDE

L'Extreme Science and Engineering Discovery Environment (XSEDE) è un'organizzazione finanziata dalla NSF che integra e coordina la condivisione di servizi digitali avanzati come supercomputer e risorse di visualizzazione e analisi dei dati di fascia alta.³⁶ XSEDE è una partnership collaborativa di 19 istituzioni, o "Fornitori di servizi", molti dei quali sono centri di supercalcolo o senza scopo di lucro presso le università e forniscono strutture informatiche per i ricercatori XSEDE. e le arti.³⁸ Le allocazioni XSEDE sono disponibili per qualsiasi ricercatore o educatore presso un istituto accademico, di ricerca senza scopo di lucro o educativo degli Stati Uniti, esclusi gli studenti.³⁹ Tuttavia, i ricercatori possono condividere le loro allocazioni creando account utente con altri collaboratori, inclusi gli studenti.⁴⁰

I ricercatori hanno due percorsi diversi per richiedere le allocazioni: Startup Allocation e Research Allocation. Le allocazioni di avvio ripartiscono le risorse XSEDE per attività computazionali su piccola scala.⁴¹ Le allocazioni di avvio sono uno dei modi più rapidi per ottenere l'accesso e iniziare a utilizzare le risorse XSEDE, poiché le richieste vengono generalmente esaminate e assegnate entro due settimane.⁴² Le richieste di allocazione di avvio richiedono anche una documentazione minima: l'abstract del progetto e il curriculum vitae (CV) dei ricercatori.⁴³ Gli Startup Assegnazioni durano tipicamente un anno, ma le richieste supportate da borse di studio valutate nel merito possono richiedere assegnazioni che durano fino a tre anni. I ricercatori possono anche presentare richieste di rinnovo se il loro lavoro necessita di risorse continue di basso livello.⁴⁴

Per esigenze di ricerca che vanno oltre il computazionale limiti previsti da un'allocazione di avvio, i ricercatori devono presentare una richiesta di allocazione di ricerca.⁴⁵ XSEDE incoraggia vivamente i propri utenti a richiedere un'allocazione di avvio prima di richiedere un'allocazione di ricerca, al fine di ottenere risultati di riferimento e documentare in modo più accurato le proprie esigenze di ricerca nell'allocazione di ricerca. ⁴⁶ Le richieste di allocazione della ricerca devono includere una serie di documenti, come un piano di utilizzo delle risorse, una relazione sullo stato di avanzamento, calcoli delle prestazioni del codice, CV e referenze.⁴⁷ Le richieste vengono accettate e riviste trimestralmente dall'XSEDE Resource Allocations Committee (XRAC), che valuta l'adeguatezza metodologica delle proposte, l'adeguatezza del piano di ricerca, l'uso efficiente delle risorse e il merito intellettuale.⁴⁸

XSEDE si attiene a una "regola del progetto unico", in base alla quale ogni ricercatore dispone di una sola assegnazione XSEDE per le proprie attività di ricerca.⁴⁹ Ad esempio, se un ricercatore ha diverse borse di studio che richiedono supporto computazionale, quelle linee di lavoro

devono essere riuniti in un'unica richiesta di assegnazione. Ciò riduce al minimo lo sforzo richiesto dal ricercatore per inviare le richieste e riduce il sovraccarico nella revisione di tali richieste.

XSEDE usa anche un "Programma Campus Champion" per semplificare l'accesso alle risorse.⁵⁰ Il Programma Campus Champion è un gruppo di oltre 700 Campus Champions che sono dipendenti o affiliati di oltre 300 college, università e istituti di ricerca statunitensi.⁵¹ Questi Campus Champions facilitano e supportano utilizzo delle risorse assegnate da XSEDE da parte di ricercatori, educatori e studenti nei loro campus. Ad esempio, i Campus Champions ospitano sessioni di sensibilizzazione e seminari di formazione per i ricercatori delle loro istituzioni, raccogliendo anche informazioni su problemi e sfide che devono essere affrontati dai proprietari delle risorse XSEDE.⁵²

Infine, XSEDE accoglie con favore le opportunità di collaborazione con altri membri della comunità scientifica e di ricerca.⁵³ Ad esempio, XSEDE assiste altre organizzazioni nell'acquisizione e nella gestione delle risorse informatiche e aiuta ad allocare e gestire l'accesso a tali risorse. Recentemente, XSEDE ha collaborato con accademici e industria privata per formare il COVID-19 High Performance Computing Consortium, che fornisce ai ricercatori potenti risorse informatiche per comprendere meglio COVID-19 e sviluppare trattamenti per affrontare le infezioni.⁵⁴

PUNTI CHIAVE

■ **Infrastruttura finanziata dalla Confederazione:**

XSEDE è un NSF iniziativa finanziata che integra e coordina condiviso risorse di supercalcolo e analisi dei dati con i ricercatori.

■ **Accesso a più livelli a**

elaborazione: per l'accesso di base all'elaborazione, XSEDE sfrutta un processo di revisione rapido e con pochi ostacoli. Per l'accesso oltre la linea di base, XSEDE ha il proprio allocazioni delle risorse commissione che esamina le domande ogni trimestre.

■ **"Campioni Campus Programma: "** XSEDE collabora con dipendenti e affiliati di college, università e istituti di ricerca

per aiutare i ricercatori ad accedere al computer risorse.

■ **Collaborazione:** XSEDE collabora con il settore privato nell'acquisizione, funzionamento e gestione delle risorse di calcolo.

Le piattaforme cloud commerciali forniscono anche altro vantaggi per il CNR. Il lavoro di gestione, manutenzione e aggiornamento dell'hardware alla base dell'NRC verrebbe gestito da parti private che hanno già esperienza nella gestione di servizi cloud su larga scala e hanno investito miliardi di dollari per farlo. Questa disposizione consente ai ricercatori di accedere a una maggiore varietà di hardware che viene costantemente ampliato e aggiornato.⁵⁵ Con un forte incentivo economico a continuare a migliorare le offerte cloud, i servizi cloud commerciali offrono un assortimento di tipi di istanza, ovvero le varie permutazioni e combinazioni di GPU /CPU, memoria, storage e specifiche di rete che costituiscono un'istanza di calcolo, con hardware diverso in una gamma di fasce di prezzo. Pertanto, i ricercatori avrebbero la flessibilità di scegliere sia l'hardware che meglio si adatta alle esigenze di

i loro progetti e il modo migliore per allocare i loro crediti cloud limitati. I ricercatori potrebbero anche avere accesso a tecnologie all'avanguardia appositamente progettate per la ricerca sull'IA, come chip ottimizzati per l'addestramento e l'inferenza, sviluppati e utilizzati esclusivamente da fornitori di cloud commerciali.

Tuttavia, l'utilizzo di servizi cloud commerciali per l'NRC comporta notevoli compromessi. Mentre i costi iniziali per sovvenzionare i crediti cloud potrebbero essere inferiori rispetto alla costruzione di un'infrastruttura pubblica, molti studi dimostrano che affidarsi a servizi cloud commerciali sarebbe probabilmente molto più costoso a lungo termine.⁵⁶ Ad esempio, uno studio del Community Cluster Program della Purdue University mostra che il costo ammortizzato del suo cluster on-premise in cinque anni è 2,73 volte più economico rispetto all'utilizzo di AWS, 3,24 volte più economico rispetto all'utilizzo di Azure e 5,54 volte più economico rispetto all'utilizzo di Google Cloud.⁵⁷ Uno studio simile dell'Università dell'Indiana stima che l'investimento totale nel suo II supercomputer di proprietà locale, Big Red II, è di circa 10,1 milioni di dollari, mentre il costo totale di una prenotazione di tre anni su AWS è di circa 24,9 milioni di dollari.⁵⁸ I confronti dei costi in altri studi sono ancora più drammatici. Ad esempio, uno studio sui cluster Advanced Research Computing presso Virginia Tech mostra che il costo quinquennale per il suo cloud locale è di circa \$ 15,5 milioni, mentre il costo quinquennale per le istanze AWS riservate che utilizzano gli stessi carichi di lavoro sarebbe di circa \$ 136,3 milioni.⁵⁹

Cosa spiega queste disparità di costo? Stime che confrontano i servizi cloud commerciali con un HPC dedicato

Mentre i costi iniziali per sovvenzionare i crediti cloud potrebbero essere inferiori rispetto alla costruzione di un'infrastruttura pubblica, molti studi dimostrano che affidarsi a servizi cloud commerciali probabilmente sarà molto più costoso a lungo termine.

cluster mostrano che i servizi cloud commerciali sono più costoso per ciclo di elaborazione.⁶⁰ Almeno in parte, ciò è dovuto al fatto che i servizi commerciali sono ottimizzati per le applicazioni commerciali. Compute Canada, ad esempio, ha scoperto che costruire la propria infrastruttura era più economico rispetto all'utilizzo di servizi commerciali, perché non avevano le stesse esigenze di utilizzo di base dei clienti commerciali, un compromesso che ha fatto guadagnare al proprio sistema più potenza di calcolo a scapito della disponibilità.⁶¹ Sebbene l'analisi è stata pubblicata nel 2016, il benchmarking dei costi di Compute Canada ha concluso:

Attualmente, è molto più conveniente per la federazione Compute Canada procurarsi e gestire un'infrastruttura informatica interna piuttosto che esternalizzare a fornitori di cloud commerciali. . . . I costi basati sul cloud variavano da 4 a 10 volte in più rispetto al costo di proprietà e gestione dei nostri cluster. Alcuni componenti erano notevolmente più costosi, in particolare lo storage persistente che era 40 volte il costo dello storage di Compute Canada.⁶²

In definitiva, la differenza di costo tra servizi cloud commerciali e sistemi HPC dipende dalla frequenza e dall'efficienza con cui viene utilizzato il sistema HPC. Forniamo un calcolo dei costi che aggiorniamo Compute Canada di seguito, arrivando a differenze di costo di grandezza comparabile. Istanze cloud commerciali con hardware comparabile in uso costante, anche con sconti sostanziali,

sarebbe significativamente più costoso nel tempo per l'NRC rispetto a un sistema HPC dedicato. Portare il costo dei servizi cloud commerciali al di sotto di quello di un HPC

Il sistema richiederebbe ai responsabili politici di negoziare sconti eccezionalmente elevati con i fornitori di cloud commerciali o di fare grandi sacrifici nella velocità dell'hardware o nella scala complessiva dell'NRC. È anche un calcolo dei costi simile

cosa ha portato la Stanford University a investire contemporaneamente sia in hardware on-premise che in una soluzione commerciale basata su cloud per la sua iniziativa Population Health Sciences (vedere il case study del riquadro nel capitolo 3). La pratica più comune nei centri NSF, come l'iniziativa XSEDE (vedere il case study del riquadro di seguito), è anche quella di costruire infrastrutture invece di fare affidamento su crediti cloud commerciali, a causa di queste considerazioni sui costi.

Infine, affidarsi al cloud commerciale può sollevare interrogativi sul consolidamento del settore. Ci sono due risposte principali a questa domanda. Uno è che la creazione di cluster HPC dedicati e di proprietà pubblica richiederebbe l'acquisto di hardware sofisticato da attori del settore esistenti, che esistono anche in settori concentrati. In altre parole, è difficile immaginare che nessuna delle due opzioni coinvolga l'industria privata. Un altro grande vincolo risiede nel tempo: un'infrastruttura pubblica completamente matura NRC non potrebbe essere messa in piedi dall'oggi al domani. Inoltre, un cloud di proprietà pubblica richiederebbe comunque a una grande azienda tecnologica di costruire l'infrastruttura sotto contratto, come nel caso di National Labs, o utilizzando una sovvenzione, come nel caso di XSEDE.

INFRASTRUTTURE PUBBLICHE

Costruire un nuovo cluster HPC sarebbe una soluzione su misura, su misura per soddisfare le esigenze di calcolo specifiche dell'NRC. Questo approccio sarebbe un territorio relativamente ben esplorato per il governo federale.⁶³ Il Dipartimento dell'Energia degli Stati Uniti (DOE) e il Dipartimento della Difesa degli Stati Uniti (DOD) già stipulano regolarmente contratti con una manciata di aziende per costruire cluster HPC ogni pochi anni.⁶⁴ Lo stesso DOE utilizza già due dei tre cluster HPC più veloci al mondo e ha recentemente finanziato lo sviluppo di due nuovi supercomputer che, una volta completati, saranno i più veloci del mondo con un margine significativo.⁶⁵ La National Science Foundation concede comunemente sovvenzioni per la costruzione di infrastruttura di calcolo ad alte prestazioni.⁶⁶ Detto questo

familiarità, i responsabili politici avrebbero stime ragionevoli di quanto costerebbe un nuovo cluster HPC per l'NRC

e avrebbe già rapporti con le imprese che presenterebbero le offerte per l'appalto.

Il costo dell'hardware per una tale scala di elaborazione è, ovviamente, notevole.⁶⁷ Ad esempio, il supercomputer IBM utilizzato presso l'Oak Ridge National Laboratory (ORNL), noto come "Summit", costava 200 milioni di dollari.⁶⁸ Al momento del suo completamento nel 2018, Summit è stato il supercomputer più veloce del mondo e, a partire dal 2020, è ancora il secondo più veloce.⁶⁹ Frontier, il nuovo supercomputer Cray in costruzione all'ORNL nel 2021, è costato 500 milioni di dollari. Una volta completato, si prevede che sarà il supercomputer più veloce al mondo, "fino a 50 volte" più veloce del Summit. esigenze.

Un tale sistema sarebbe più efficiente in termini di costo per ciclo a lungo termine rispetto a sovvenzionare i servizi cloud commerciali. L'NRC potrebbe anche espandere e aggiornare più cluster nel tempo per soddisfare le mutevoli esigenze e l'ambito del programma.

Inoltre, un cluster dedicato per l'NRC ha il vantaggio di dare al governo federale un maggiore controllo sulle risorse computazionali (ad esempio, riducendo l'incertezza sui prodotti e sulle piattaforme, come l'improvvisa deprecazione delle API richieste). Questo livello di controllo sull'hardware consente inoltre ai responsabili politici una maggiore flessibilità con le operazioni NRC. Adottare l'approccio dell'infrastruttura pubblica (vale a dire, "fare" non "acquistare") comporta diversi compromessi significativi da pesare rispetto agli obiettivi politici dell'NRC. In primo luogo, la costruzione di un nuovo cluster HPC richiederebbe circa due anni, oltre al tempo necessario per sollecitare e valutare le proposte di potenziali appaltatori.⁷¹ Se l'NRC spera di stimolare rapidamente e aiutare a democratizzare la ricerca sull'IA negli Stati Uniti, il programma non sarebbe l'ideale, data la rapidità con cui avanzano le scoperte dell'IA. Naturalmente, anche stipulare contratti con fornitori di servizi cloud o concedere sovvenzioni per la costruzione di supercomputer richiederebbe un processo. Tuttavia, la creazione di un cluster potrebbe sollevare problemi contrattuali più impegnativi, come il superamento del budget e i ritardi del progetto.⁷² L'esperienza degli appaltatori nella creazione di questo tipo di hardware può aiutare a mitigare alcune di queste preoccupazioni, così come il loro interesse

preso in considerazione per futuri contratti governativi. Ma i rischi sono comunque ancora presenti.

In secondo luogo, l'usabilità e il set di funzionalità del software stack per le infrastrutture pubbliche non è affatto provato. Uno degli ostacoli più comuni all'adozione del cloud computing da parte dei ricercatori risiede nell'usabilità dei sistemi,⁷³ e l'infrastruttura pubblica ha meno esperienza nell'agevolare il percorso di onboarding alla scala contemplata. Questo è il motivo per cui raccomandiamo un progetto pilota per valutare se un centro HPC nazionale può essere amministrato in modo da garantire la facilità di transizione al cloud e lo stack software che i ricercatori sono abituati a utilizzare con fornitori privati.

In terzo luogo, anche i responsabili politici dovrebbero renderne conto costi di manutenzione e amministrazione del sistema.⁷⁴ Avrebbero bisogno di trovare strutture per alloggiare e gestire l'hardware e per tenere conto degli elevati costi energetici della gestione di un cluster HPC, oltre alla prevenzione dei disastri

e costo di ripristino.⁷⁵ Questi costi sono significativi. Nel 2021, l'Oak Ridge Leadership Computing Facility ha richiesto 225 milioni di dollari per far funzionare tutti i suoi sistemi.⁷⁶ L'Argonne Leadership Computing Facility, a sua volta, ha richiesto 155 milioni di dollari.⁷⁷ Inoltre, il ciclo di vita dei sistemi DOE HPC è tradizionalmente di circa sette anni, dopo quali nuovi sistemi vengono costruiti e quali vecchi vengono dismessi.⁷⁸ Sebbene non sia chiaro quale sarà la durata di vita dei nuovi sistemi, questa cifra di sette anni ci porterebbe a sostenere che l'NRC dovrebbe aspettarsi di aggiornare i suoi sistemi o di costruirne di nuovi con alcuni grado di regolarità.

Infine, dando al governo federale un maggiore controllo sulle risorse informatiche non renderebbe immediatamente l'NRC al sicuro dagli attacchi.⁷⁹ Come per l'utilizzo di un'infrastruttura cloud commerciale, la sicurezza dipenderà principalmente dal modello di accesso ai dati implementato dall'NRC.⁸⁰ Discutiamo approfonditamente i problemi di sicurezza nel capitolo 8.

CASO DI STUDIO: FUGAKU

Nel 2014, il Ministero giapponese dell'istruzione, della cultura, dello sport, della scienza e della tecnologia ha lanciato una partnership pubblico-privata tra il Riken Institute, l'Organizzazione di ricerca per la scienza e la tecnologia dell'informazione (RIST), finanziata dal governo, e Fujitsu per creare il supercomputer successore del K computer che supporta un'ampia gamma di applicazioni scientifiche e sociali.⁸¹ Il risultato è stato Fugaku, che è stato nominato il supercomputer più veloce del mondo nel 2020.⁸²

L'obiettivo tecnico di Fugaku era di essere 100 volte più veloce del precedente computer K, con una prestazione di 442 petaFLOPS nel benchmark LINPACK ad alte prestazioni FP64 del TOP500 . è composto da poco più di 150.000 CPU collegate, con ciascuna CPU che utilizza chip per computer con licenza ARM.⁸⁵ Nonostante avesse circa 1,9 volte più parti rispetto al suo predecessore del computer K, Fugaku è stato completato in tre mesi in meno.⁸⁶ Il budget di sei anni per Fugaku era circa \$ 1 miliardo.⁸⁷

RIST ha sollecitato proposte di utilizzo attraverso il "Programma di Promozione Ricerca sul supercomputer Fugaku. Nell'ambito del programma, Fugaku è già stato utilizzato per studiare l'effetto di maschere e goccioline respiratorie al fine di informare la politica giapponese durante la pandemia di COVID-19.⁸⁸ Per l'anno fiscale 2021, sono stati selezionati 74 progetti pubblici e industriali per l'accesso su larga scala a Fugaku .⁸⁹ Attualmente, RIST sta ancora richiedendo proposte che rientrino in specifiche categorie di utilizzo, e qualsiasi ricercatore interessato può fare domanda.⁹⁰

PUNTI CHIAVE

- **Potenza di calcolo significativa:**
Fugaku è stato il più veloce supercomputer nel 2020.
- **Nuovo processo di applicazione per la potenza di calcolo:**
Le applicazioni sono state sollecitate per testare il supercomputer su una serie di compiti e avere il controllo su chi ha ricevuto potenza di calcolo.

CONFRONTO DEI COSTI

Per concludere questo capitolo, forniamo un costo approssimativo confronto tra un servizio cloud commerciale leader e un sistema HPC governativo dedicato (IBM Summit) (vedere l'Appendice A per i dettagli). Rimandiamo il lettore al lavoro sostanziale che è stato pubblicato sull'economia del cloud computing per un'analisi più completa, gran parte del quale sottolinea la varianza nella domanda di computer.⁹¹

Si prevede che la costruzione di un'infrastruttura pubblica autonoma sarà meno costosa rispetto all'implementazione dell'NRC attraverso un accordo di contratto con i fornitori nell'arco di cinque anni. Con uno sconto del 10% sulle tariffe standard per cinque anni e in condizioni di utilizzo costante, l'opzione di cloud computing più potente di AWS (nota come istanze P3) potrebbe costare 7,5 volte i costi totali stimati di Summit, utilizzando hardware comparabile. Utilizziamo uno sconto del 10% che è stato negoziato da un'importante università di ricerca con un fornitore di servizi cloud commerciali. Al contrario, il governo dovrebbe negoziare uno sconto dell'88% affinché AWS sia competitivo in termini di costi con un cluster HPC dedicato a lungo termine. Anche in uno scenario in cui l'utilizzo di NRC oscilla notevolmente, il cloud computing commerciale potrebbe costare 2,8 volte il costo stimato di Summit. (Sebbene la variabilità nell'utilizzo influisca pesantemente su queste stime, l'uso di pianificatori può contribuire a livellare la domanda.⁹²)

Queste stime dei costi hanno limitazioni importanti. Primo, governo può essere in grado di negoziare il costo verso il basso. Abbiamo utilizzato come punto di riferimento l'accordo aziendale di una delle principali università con AWS, che offre uno sconto del 10% rispetto alle tariffe di mercato. Ma, a meno che lo sconto negoziato non riguardi ordini di entità maggiore, il cloud commerciale rimarrà significativamente più costoso. In secondo luogo, queste stime dei costi si concentrano principalmente sull'informatica.⁹³ Come ha mostrato l'analisi di Compute Canada, la differenza di costo nello storage era ancora maggiore. In terzo luogo, l'uso di tariffe commerciali è probabilmente più favorevole per i fornitori di cloud, poiché gli standard di sicurezza del governo in genere aumentano le tariffe a causa dei requisiti normativi. Ad esempio, un requisito di "sovranità dei dati" per dati e hardware che risiedono negli Stati Uniti, o requisiti di cloud privato per determinati set di dati di agenzie, possono aumentare significativamente il costo del cloud computing commerciale. In quarto luogo, questo semplice confronto dei costi

è statico e non riflette i cambiamenti nei costi dell'hardware e nelle strutture dei prezzi che potrebbero verificarsi in un periodo di cinque anni in condizioni di mercato in rapida evoluzione.

Ma, se l'NRC si espandesse effettivamente, i sistemi verrebbero acquistati in modo incrementale nel tempo, aggiornando le risorse disponibili e fornendo opzioni a prezzi diversi, simili alle attuali opzioni commerciali. Infine, come notato sopra, queste stime dei costi tengono conto della manutenzione prevista per il Summit, ma potrebbero non tenere conto di tutti questi costi non hardware, motivo per cui raccomandiamo un progetto pilota per esplorare la possibilità di aprire le strutture informatiche del governo a utenti NRC.

In breve, offriamo questo semplice confronto per evidenziare alcune delle considerazioni di costo salienti per la decisione make-or-buy, che arriva a una conclusione molto simile all'analisi fatta da Compute Canada.

CASO DI STUDIO: COMPUTE CANADA

Compute Canada si è formata nel 2006 come partnership tra le regioni canadesi organizzazioni accademiche HPC di condividere le infrastrutture in tutto il Canada.⁹⁴ La missione dichiarata dell'organizzazione è quella di "consentire l'eccellenza nella ricerca e nell'innovazione a beneficio del Canada, implementando in modo efficace, efficiente e sostenibile una rete informatica di ricerca avanzata all'avanguardia supportata da competenza di livello mondiale."⁹⁵

L'infrastruttura di Compute Canada include cinque sistemi HPC ospitati presso università di ricerca in tutto il Canada.⁹⁶ Dal 2015 al 2019, Compute Canada ha utilizzato circa 125 milioni di dollari canadesi (CAD) in finanziamenti per costruire quattro di questi sistemi.⁹⁷ Hanno anche studiato l'utilizzo di risorse cloud commerciali invece di costruire questi nuovi sistemi.⁹⁸ Tuttavia, alla fine hanno concluso che affidarsi a fornitori di cloud commerciali sarebbe molto più costoso e non potrebbe fornire la latenza desiderata per la ricerca su larga scala e ad alta intensità di dati . CAD) per il funzionamento dei suoi sistemi HPC e \$ 20 milioni (CAD) per supporto, formazione e sensibilizzazione.¹⁰⁰ La domanda di risorse HPC di Compute Canada supera di gran lunga l'attuale capacità dell'infrastruttura e si prevede che continuerà a crescere.¹⁰¹ Nel 2018, Compute Canada ha stimato che avrebbero bisogno di circa \$ 90 milioni (CAD) all'anno per cinque anni da investire nell'espansione dell'infrastruttura fino al punto in cui potrebbe soddisfare le richieste previste.¹⁰²

Circa 16.000 ricercatori di tutte le discipline scientifiche utilizzano Compute Canada infrastrutture per supportare il loro lavoro.¹⁰³ Compute Canada distribuisce le sue risorse in due modi. In primo luogo, i Principal Investigator e gli utenti sponsorizzati possono richiedere un'allocazione di risorse senza priorità pianificata per il proprio gruppo di ricerca.¹⁰⁴ Compute Canada ritiene che molti gruppi di ricerca possano soddisfare le proprie esigenze di calcolo in questo modo.¹⁰⁵ In alternativa, i ricercatori che necessitano di risorse maggiori o prioritarie possono presentare un progetto proposto ai "Research Allocation Competitions" annuali.¹⁰⁶ Le proposte presentate vengono sottoposte a una revisione scientifica tra pari e a una revisione del personale tecnico per valutarne i meriti.¹⁰⁷ La revisione scientifica esamina l'eccellenza scientifica e

la fattibilità del progetto di ricerca specifico, l'adeguatezza delle risorse richieste per raggiungere gli obiettivi del progetto e la probabilità che le risorse richieste vengano utilizzate in modo efficiente.¹⁰⁸ Questa revisione è condotta su base volontaria da 80 esperti specifici della disciplina provenienti da istituzioni accademiche canadesi .¹⁰⁹ La revisione tecnica è condotta dallo stesso personale di Compute Canada, che verifica l'accuratezza delle risorse computazionali necessarie per ciascun progetto, sulla base dei requisiti tecnici delineati nella domanda, e formula raccomandazioni su quali risorse dovrebbero essere allocate per soddisfare le esigenze del progetto. ¹¹⁰ Nel 2021, Compute Canada ha ricevuto 651 domande per la Research Allocation Competition e ha esaminato tutte le domande nell'arco di cinque mesi.¹¹¹

PUNTI CHIAVE

■ **Accesso predefinito** **con livelli:** tutti

I PI possono accedere a un pianificatore allocazione delle risorse di calcolo senza priorità con un processo di applicazione integrato per richiedere di più. La maggior parte dei ricercatori trovare l'impostazione predefinita allocazione sufficiente per i loro bisogni.

■ **Domanda** **ampiamente** **utilizzata e in aumento:** Compute del Canada infrastruttura è ampiamente utilizzato in ambito accademico discipline, con una domanda in costante aumento risorse. Compute Canada intende farlo investire pesantemente in infrastrutture per soddisfare le crescenti richieste.

Capitolo 3:

Protezione dell'accesso ai dati

Dopo aver calcolato le risorse, la prossima decisione di progettazione critica per l'NRC è come archiviare e fornire ai suoi utenti l'accesso ai set di dati: l'obiettivo di "accesso ai dati" dell'NRC.

In effetti, come articolato nell'invito all'azione originale dell'NRC, le agenzie governative dovrebbero "raddoppiare i loro sforzi per rendere disponibili gratuitamente dati più numerosi e di migliore qualità per la ricerca pubblica", in quanto "alimentano" scoperte uniche nella ricerca. I problemi socialmente più significativi dipendono da set di dati ampi ma inaccessibili nel settore pubblico. Dai dati sul clima ospitati dalla National Oceanic and Atmospheric Administration (NOAA), dai dati sanitari del più grande sistema sanitario integrato del paese nel Department of Veterans Affairs (VA), o dai dati sull'occupazione nel Department of Labor (DOL), tali dati potrebbero alimentare sia la ricerca fondamentale che utilizza l'intelligenza artificiale sia il riorientare gli sforzi dai progetti incentrati sul consumatore (ad esempio, l'ottimizzazione della pubblicità) verso argomenti più socialmente urgenti (ad esempio, il cambiamento climatico).

Come notato nell'accusa del Congresso, facilitare un ampio accesso ai dati è fondamentale pilastro del NRC. È importante sottolineare che, come discuteremo di seguito, limitiamo l'ambito delle nostre raccomandazioni a facilitare l'accesso ai dati del governo del settore pubblico, che come condizione per l'accesso ai dati amministrativi del governo, i ricercatori NRC dovrebbero utilizzare solo per scopi di ricerca accademica. Gli utenti di NRC dovrebbero anche essere in grado di eseguire calcoli su qualsiasi set di dati privato a loro disposizione. Esistono meccanismi disponibili per la condivisione di tali set di dati, ma identifichiamo la principale sfida dell'NRC nel fornire l'accesso a dati governativi precedentemente non disponibili.

I dati del governo sono intenzionalmente decentralizzati. In base alla progettazione del Privacy Act del 1974, non esiste un archivio centralizzato per i dati del governo degli Stati Uniti o un metodo fondamentale per collegare i dati tra le agenzie governative.² Il risultato è un'infrastruttura di dati tentacolare e decentralizzata con livelli molto diversi di finanziamento, esperienza, applicazione di standard e accesso e condivisione delle policy. Pertanto, l'NRC dovrà sviluppare una strategia dei dati unificata che possa funzionare con un'ampia gamma di agenzie, standard di sicurezza adottati in modo non uniforme e all'interno della legislazione sulla privacy dei dati esistente.

Gli sforzi precedenti hanno cercato di migliorare l'accesso e la condivisione dei dati federali, sia tra agenzie che con ricercatori esterni, ma esistono ancora ostacoli significativi per consentire l'accesso alla ricerca sull'IA del tipo richiesto dall'NRC.³ Collegando le politiche di governance dei dati con l'accesso a compute, basandosi su modelli di successo esistenti e collaborando con le agenzie per creare sistemi interoperabili che soddisfino i problemi di sicurezza e privacy, l'NRC può consentire un maggiore accesso ai dati che aiuteranno i ricercatori di intelligenza artificiale a rispondere a pressanti domande scientifiche e sociali e ad aumentare l'innovazione dell'IA.⁴

CHIAVE ASPORTO

- L'NRC dovrebbe adottare un modello a più livelli per l'accesso e l'archiviazione dei set di dati delle agenzie federali. I livelli dovrebbero corrispondere alla sensibilità dei dati.
- L'NRC può aiutare ad armonizzare il panorama frammentato della condivisione dei dati federali.
- L'NRC dovrebbe prendere in considerazione incentivare la partecipazione delle agenzie concedendo alle agenzie che forniscono dati il diritto di utilizzare il calcolo NRC risorse.
- L'NRC dovrebbe sequenziare strategicamente l'acquisizione dei dati concentrandosi prima sui set di dati a rischio da basso a moderato che lo sono attualmente inaccessibile.
- A causa di vincoli legali e molte opzioni esterne, il NRC dovrebbe concentrare i suoi sforzi sulla semplificazione dell'accesso ai set di dati del governo. I ricercatori dovrebbero ancora essere autorizzati a utilizzare le risorse di calcolo NRC su set di dati privati, a condizione che i ricercatori certifichino di avere i diritti per utilizzare tali dati.

Per prima cosa spiegheremo perché l'NRC dovrebbe concentrare i suoi sforzi sull'agevolazione della condivisione dei dati del governo federale piuttosto che della condivisione dei dati del settore privato. Esaminiamo quindi come e perché lo status quo per la condivisione dei dati federali non riesce a realizzare l'enorme potenziale dei dati del governo. Sebbene il concetto di centralizzare fonti di dati disparate per sbloccare gli approfondimenti della ricerca non sia nuovo,⁵ ci sono sfide uniche per farlo nel contesto dell'NRC. Discuteremo anche gli elementi chiave del nostro modello proposto: (1) l'uso di FedRAMP come sistema per classificare i set di dati in base alla loro sensibilità e per modificare l'accesso ad essi tramite credenziali a più livelli per gli utenti NRC; (2) promozione degli sforzi di standardizzazione e armonizzazione interagenzia per modernizzare le pratiche di condivisione dei dati; e (3) considerazioni strategiche su come mettere in sequenza gli sforzi per semplificare l'accesso a particolari set di dati.

I casi di studio inclusi in questo White I documenti sono stati scelti come esempi di iniziative di condivisione dei dati di successo⁶ e per illustrare la gamma di decisioni di progettazione disponibili. Sebbene ogni caso di studio fornisca uno sguardo unico su approcci diversi, emergono alcuni temi comuni. Innanzitutto, molte delle entità di condivisione dei dati che abbiamo studiato non solo hanno un unico punto di accesso per i ricercatori per richiedere l'accesso, ma consentono anche alle agenzie governative di mantenere un certo controllo sui requisiti di accesso ai propri dati. Come discuteremo di seguito, questa concezione dell'NRC come intermediario di dati fornirebbe vantaggi reali nella razionalizzazione dell'accesso ai dati pur mantenendo la fiducia tra le agenzie che desiderano proteggere i propri dati. In secondo luogo, alcune iniziative utilizzano i finanziamenti e la formazione del personale come carote per incentivare le agenzie a impegnarsi nella condivisione dei dati. L'NRC può imparare da queste iniziative nel formulare la propria serie di incentivi per le agenzie.

CONDIVISIONE DEI DATI PRIVATI

L'NRC dovrebbe facilitare affermativamente il set di dati privato condivisione? Sebbene ci siano indubbi vantaggi nel fornire ai ricercatori l'accesso ai dati privati,⁷ l'NRC avrà il suo maggiore impatto concentrando i suoi sforzi prima sui meccanismi per accedere e condividere i dati del governo.

Come questione iniziale, una varietà di meccanismi per la condivisione generale dei dati esiste già.⁸ Settore privato

le parti interessate, inoltre, possono e hanno spesso costruito le proprie piattaforme interne per consentire l'accesso a set di dati approvati riducendo al minimo i problemi di proprietà intellettuale⁹ o fornire l'accesso alle proprie interfacce di programmazione delle applicazioni (API) per rendere i dati open source più facilmente accessibili.¹⁰ Concentrandosi sulla fornitura di accesso ai dati del settore pubblico, in particolare ai dati amministrativi che sono tradizionalmente inaccessibili alla maggior parte dei ricercatori,¹¹ l'NRC svolgerebbe un ruolo unico e pertinente per i ricercatori di tutte le discipline senza doversi occupare di complesse preoccupazioni relative ai dati del settore privato o la necessità di incentivare la partecipazione di attori non governativi.

Dall'NRC deriverebbero complessi problemi di proprietà intellettuale consentire, facilitare o persino richiedere alle parti interessate del settore privato e ai ricercatori indipendenti di condividere liberamente i propri dati privati insieme ai dati del settore pubblico.

Complesse preoccupazioni in materia di proprietà intellettuale deriverebbero dal fatto che l'NRC consenta, faciliti o addirittura richieda alle parti interessate del settore privato e ai ricercatori indipendenti di condividere liberamente i propri dati privati insieme ai dati del settore pubblico. In primo luogo, ciò comporterebbe domande complesse su quali licenze dovrebbero essere disponibili o obbligatorie

utenti di NRC al fine di incoraggiare la condivisione dei dati, nonostante i timori su come tale condivisione possa influire sulla redditività e sulla commercializzazione future. Mentre l'obbligo di una licenza open-source (ad esempio, Creative Commons) avvantaggerebbe maggiormente i ricercatori fornendo il più ampio accesso ai dati e avvantaggerebbe gli amministratori NRC rimuovendo alcune possibili preoccupazioni di violazione della proprietà intellettuale, le parti interessate del settore privato potrebbero sentirsi scoraggiate dal caricamento di conseguenza.

Al contrario, se gli utenti possono scegliere di adottare una licenza che

consente loro di preservare i propri diritti di proprietà intellettuale, le parti interessate del settore privato potrebbero sentirsi più a proprio agio nel condividere i propri dati, ma ciò sposterebbe alcune responsabilità sugli utenti, o sullo stesso NRC, facendo affidamento sugli utenti per rispettare la licenza. Ciò comporterebbe un'enfasi sull'applicazione, che va dalle spiegazioni e dalle esclusioni di responsabilità degli utenti allo standard del settore di un sistema completo di notifica e rimozione.

I proprietari di dati potrebbero voler impedire il caricamento di opere protette da copyright, ad esempio, chiedendo allo stesso NRC di valutare se i dati privati sono già protetti da copyright. Gli standard di settore per condurre la diligenza dei dati, utilizzando strumenti manuali o automatizzati, richiederebbero molta manodopera¹² o costi proibitivi.

Nessuno dei precedenti impedirebbe ai ricercatori di farlo utilizzando le risorse di calcolo NRC sui propri set di dati privati. Come gli attuali fornitori di servizi cloud, l'NRC può stipulare in un accordo di licenza con l'utente finale (EULA) che i ricercatori devono accettare di possedere i diritti di proprietà intellettuale sui dati che stanno utilizzando.¹⁴ Questo EULA può anche attribuire la responsabilità all'utente finale, piuttosto che al NRC, per qualsiasi utilizzo dei dati gravato dalle disposizioni vigenti in materia di PI. Inoltre, la discussione di cui sopra riguarda se i ricercatori dovrebbero essere tenuti a condividere i loro dati privati, non se i ricercatori dovrebbero essere tenuti a condividere i risultati delle loro ricerche condotte sull'NRC. Quest'ultimo punto è discusso nel capitolo 9.

L'ATTUALE PATCHWORK SISTEMA PER L'ACCESSO DATI FEDERALI

L'NRC potrebbe svolgere un ruolo fondamentale semplificando l'accesso ai dati governativi in un sistema attualmente decentralizzato.¹⁵ In alcuni casi, le agenzie potrebbero semplicemente non disporre di un metodo standardizzato per la condivisione dei dati.¹⁶ A causa di vincoli legali percepiti, rischi o problemi di sicurezza, le agenzie spesso hanno scarso incentivo pratico a condividere i propri dati.¹⁷ Esempi di successo di ricercatori che ottengono l'accesso a dati governativi da singole agenzie spesso si basano sul fatto che i ricercatori abbiano relazioni personali

con gli amministratori e la disponibilità da parte dell'amministratore a contrastare questi vincoli al servizio del progetto di ricerca.¹⁸ Sebbene questo processo basato sulle relazioni abbia prodotto alcuni successi,¹⁹ il risultato molto più comune è che i dati semplicemente non vengono condivisi o cui hanno accesso i ricercatori.²⁰ In effetti, un funzionario governativo ha indicato che superare gli ostacoli alla messa a disposizione di alcuni dati governativi per la ricerca è stata la sfida più grande in una lunga carriera.

Un funzionario del governo ha indicato che superare gli ostacoli per rendere disponibili determinati dati del governo per la ricerca è stata la sfida più grande in una

Le agenzie in genere richiedono al destinatario dei dati di rispettare un accordo sull'utilizzo dei dati (DUA). Questi DUA prescrivono tali limitazioni sull'utilizzo dei dati come la durata dell'uso, lo scopo dell'uso e le garanzie sulla privacy e la sicurezza dei dati.²¹ Tuttavia, i DUA soffrono di un problema centrale: il processo di negoziazione dei DUA è altamente frammentato e incoerente tra agenzie governative, aumentando drasticamente la complessità nell'ottenere le approvazioni per loro.²² Alcune agenzie hanno un ufficio o un processo designato per gestire i DUA, ma altre agenzie si affidano a processi estemporanei e accordi ad hoc, quid pro quo.²³ Uno di questi esempi è il Research Data Assistance Center, un'unità centralizzata all'interno dei Centers for Medicare & Medicaid Services (CMS) dedicata a supportare le richieste di accesso ai dati.²⁴ Al contrario, i DUA all'interno del Department of Housing and Urban Development e del Department of Education sono gestiti in unità aziendali decentralizzate, ciascuna con diversi canali di instradamento e team legali, che possono confondere i revisori quando sono multipli

le richieste di dati tra le stesse parti vengono instradate simultaneamente ma separatamente.²⁵ In effetti, i negozianti di DUA universitari in un sondaggio si sono lamentati del fatto che il processo fosse un gioco di "patate bollenti burocratiche" e si sono chiesti: "Perché non esiste un solo modello per tutto?"²⁶ In definitiva, la mancanza di standardizzazione significa che i DUA spesso richiedono un'ampia revisione e revisione, creando notevoli ritardi.

Anche i requisiti agenzia per agenzia impediscono la condivisione dei dati. Questi requisiti possono variare dall'obbligo per i ricercatori di accedere ai dati solo presso una struttura in loco, utilizzando apparecchiature autorizzate dal governo, al limite della quantità di cicli computazionali che possono essere utilizzati per analizzare i dati o alla limitazione della quantità di dati disponibili contemporaneamente.²⁷ Queste restrizioni sono particolarmente problematico, dato che i moderni modelli di intelligenza artificiale possono richiedere enormi quantità di dati e calcoli per essere più efficaci.

In generale, le ragioni di questa disfunzione vanno da valide preoccupazioni sulla sicurezza e responsabilità al banale e prosaico. I sistemi informatici all'interno di alcune agenzie operano letteralmente decenni indietro rispetto alla frontiera tecnologica; un rapporto del 2016 del Government Accountability Office (GAO) descriveva esempi dettagliati di questi sistemi legacy, discutendo di come diverse agenzie dipendessero da hardware e software che non erano più aggiornabili e richiedessero personale specializzato per la manutenzione.²⁸ Una mancanza di incentivi, un rischio- Anche la cultura avversa e l'autorità statutaria di un'agenzia svolgono un ruolo importante nel consentire o ostacolare la condivisione dei dati.²⁹

Non siamo affatto i primi osservatori a notare questi problemi. I sostenitori hanno lavorato per anni per standardizzare e modernizzare le pratiche governative in materia di dati e tecnologia. L'approccio solido e integrato all'utilizzo dei dati per fornire in missione, servire il pubblico e amministrare le risorse dell'apprendimento automatico e dell'intelligenza artificiale in mente.

ACCESSO E STOCCAGGIO DEI DATI A LIVELLI

La natura decentralizzata dei dati del governo ha implicazioni a cascata su molti aspetti dell'ecosistema dei dati del governo. Un'area chiave che influirà sull'NRC è la mancanza di protocolli coerenti di accesso all'archiviazione e all'autenticazione tra le agenzie governative.

Perché molti set di dati governativi contengono dati sensibili (ad esempio, alto rischio dovuto a problemi di privacy individuale),³² una componente cruciale del modello di dati dell'NRC consisterà in una tassonomia di archiviazione a più livelli che distingue tra set di dati in base alla loro sensibilità e restringe di conseguenza l'accesso a diversi gruppi di ricerca. Interpretando l'archiviazione a più livelli e l'accesso come due facce della stessa medaglia, facciamo riferimento a modelli esistenti basati sui livelli di rischio del set di dati e proponiamo un quadro per l'NRC che mira a raggiungere il duplice obiettivo di semplificare il processo di consentire l'accesso della ricerca ai dati del governo mentre mantenendo la privacy e la sicurezza.

FedRAMP: un framework a più livelli per l'archiviazione dei dati nel cloud

Un tipo di tassonomia dello storage su più livelli esiste già per i servizi cloud governativi di terze parti in uno dei principali framework di sicurezza informatica del governo federale, il Federal Risk and Authorization Management Program (FedRAMP).³³ Entrato in vigore nel 2011, il framework è stato progettato per governare tutte le implementazioni, con alcune eccezioni dettagliate nel capitolo 8 di questo libro bianco. FedRAMP offre due percorsi per i fornitori di servizi cloud per ricevere l'autorizzazione federale. In primo luogo, una singola agenzia può emettere ciò che è noto come autorità di operare (ATO) a un fornitore di servizi cloud dopo che il pacchetto di autorizzazione di sicurezza del fornitore è stato esaminato dal personale dell'agenzia e l'agenzia ha identificato eventuali carenze che devono essere affrontate.³⁴ Questi tipi di ATO sono validi per ciascun fornitore in più agenzie, in quanto altre agenzie possono riutilizzare il pacchetto di sicurezza di un'agenzia iniziale nella concessione di ATO. La seconda opzione disponibile per i fornitori di servizi cloud consiste nell'ottenere un'autorizzazione provvisoria dal FedRAMP Joint Authorization Board, composto da rappresentanti del Dipartimento della Difesa (DOD), del Dipartimento della Sicurezza Nazionale (DHS) e del

Amministrazione dei servizi generali (GSA). Queste ATO provvisorie offrono garanzie alle agenzie che DHS, DOD e GSA hanno esaminato le considerazioni sulla sicurezza, ma prima che un'agenzia specifica possa utilizzare i servizi di un fornitore, tale agenzia deve emettere la propria ATO.²⁸ Sia nel primo che nel secondo caso , FedRAMP classifica i sistemi in livelli di impatto basso, moderato o alto (vedere Tabella 2).

Perché i requisiti FedRAMP si applicano a tutti i federali agenzie quando i dati federali vengono raccolti, mantenuti, elaborati, diffusi o eliminati sul cloud, l'NRC stesso dovrà essere conforme agli standard di sicurezza FedRAMP indipendentemente dalla forma organizzativa che assume.³⁵ Ogni set di dati portato all'NRC dovrebbe da rivedere sotto FedRAMP con livelli di accesso appropriati. Se un servizio cloud è già stato valutato in FedRAMP perché è stato utilizzato in passato per ospitare dati federali, il servizio può ereditare lo stesso livello di conformità FedRAMP nell'NRC senza un'ulteriore valutazione.³⁶

Oltre a classificare i set di dati, l'altra funzione di FedRAMP è identificare un set completo di "controlli", ovvero requisiti e meccanismi che i fornitori di servizi cloud devono implementare prima che il set di dati del governo possa essere ospitato su di loro.³⁷ Si basano sul National Institute of Standards and Technology (NIST) Special Publication 800-53, che fornisce standard e requisiti di sicurezza per i sistemi informativi utilizzati dal governo federale.³⁸

Questi controlli variano ampiamente e includono i requisiti come garantire che l'organizzazione che richiede la certificazione "disabiliti automaticamente gli account inattivi", "stabilisca e amministri account utente con privilegi in conformità con uno schema di accesso basato sui ruoli che organizza l'accesso al sistema e i privilegi in ruoli", "fornisce formazione sulla consapevolezza della sicurezza sul riconoscimento e la segnalazione di potenziali indicatori di minaccia interna" o sviluppa piani di sicurezza regolari in caso di violazione.³⁹ I requisiti diventano più faticosi per i dati "ad alto impatto" di FedRAMP (ad esempio, la creazione di air-gap a livello di sistema per proteggere i dati sensibili).⁴⁰

LIVELLO	TIPO DI DATI	IMPATTO DELLA VIOLAZIONE DEI DATI	NUMERO DI CONTROLLI
Rischio a basso impatto - Linea di base bassa - SaaS a basso impatto	Dati destinati all'uso pubblico	Effetti avversi limitati; preserva la sicurezza, le finanze, la reputazione o la missione di un'agenzia	125
Rischio di impatto moderato - Ad esempio, informazioni di identificazione personale	Dati non classificati controllati non disponibili al pubblico	Può danneggiare le operazioni di un'agenzia	325
Rischio di impatto elevato - Ad esempio, forze dell'ordine, assistenza sanitaria, servizi di emergenza	Informazioni federali sensibili	Impatti catastrofici come la chiusura delle operazioni di un'agenzia, la rovina finanziaria o la minaccia alla proprietà o alla vita	421

Tabella 2: i livelli FedRAMP sono designati in base al grado di rischio associato alla violazione di un sistema informativo. I livelli di base di sicurezza si basano su riservatezza, disponibilità e integrità, come definito nel Federal Information Processing Standard 199.41

L'ottenimento di queste certificazioni e la creazione di piani di conformità possono comportare costi significativi, anche se le specifiche tecniche sottostanti possono essere affrontate o esistono già. Una questione chiave per strutturare l'NRC è che gli oneri principali per garantire la conformità FedRAMP dovrebbero ricadere sul personale istituzionale dell'NRC, non sulle agenzie di origine o sui singoli ricercatori accademici. Nell'ambito del processo di certificazione FedRAMP, il personale NRC dovrà considerare come concedere l'accesso ai PI in conformità con le regole FedRAMP, ma tale processo può e dovrebbe evitare di richiedere alle agenzie di origine o alle singole università di sostenere spese sostanziali associate all'assunzione di consulenti e avvocati per certificare la conformità FedRAMP.⁴²

Mentre FedRAMP stabilisce standard comuni per l'archiviazione cloud dei dati governativi all'interno delle agenzie,⁴³ è un'eccezione a un panorama di standard federali di condivisione dei dati altrimenti balcanizzato,⁴⁴ sebbene non faciliti lo scambio di dati. L'NRC deve mantenere la conformità non solo ai requisiti FedRAMP, ma anche ai requisiti di qualsiasi agenzia con cui sta collaborando per l'accesso ai dati.⁴⁵ I sostenitori interessati ad aumentare la disponibilità dei dati del governo hanno lottato a lungo per stabilire un equivalente FedRAMP universale tra diverse agenzie che fornisca standard condivisi per la condivisione dei dati basata sulla sensibilità dei dati.⁴⁶ Come discusso nel capitolo 8, la definizione di tali standard di sicurezza universali e "centralizzati" non solo garantisce l'uniformità interna, ma rimuove anche le barriere alla condivisione dei dati.

L'implementazione degli standard FedRAMP da parte dell'NRC può anche fornire alle agenzie partner un'importante opportunità per riesaminare i propri standard e condividere tra loro le migliori pratiche⁴⁷. esigenze di ricerca. Il CNR

può trarre ispirazione dalle migliori pratiche delle agenzie, nonché da FedRAMP per sviluppare uno standard NRC comune per determinare i dati ad alto, moderato o basso rischio, nonché quali conseguenze dovrebbero derivare da tale valutazione.

Nella sezione successiva sulle considerazioni strategiche, discutiamo come abilitare questo processo incentivando le agenzie a partecipare all'NRC e selezionando i set di dati che presentano un minor rischio per la privacy e la sicurezza.

Inoltre, data la diversità dei tipi di dati e fonti che potrebbero essere archiviate sulla piattaforma, la politica NRC dovrebbe garantire che esistano standard e protezioni per l'archiviazione dei dati nelle aree in cui FedRAMP presenta punti ciechi. FedRAMP è in parte animato da rischi da parte di attori malintenzionati come criminali informatici o governi stranieri avversari, ma come discuteremo nel capitolo 6, i rischi per la privacy possono sorgere anche per il caso d'uso previsto dell'analisi da parte dei ricercatori NRC. Di particolare interesse sono i casi in cui vengono combinati insieme di dati disparati, che possono consentire nuove inferenze che rendono i dati precedentemente anonimi identificabili individualmente, anche quando i dati stessi non contenevano informazioni identificabili.⁴⁹ Tali combinazioni può anche alterare il livello di rischio originale dei dati, creando un output che merita una classificazione di rischio più elevata. Inoltre, i modelli e le rappresentazioni di apprendimento automatico possono rivelare involontariamente le proprietà dei dati utilizzati per addestrarli⁵⁰ e la diffusione di questi modelli potrebbe comportare rischi per la privacy.

Questa non è una sfida unica per l'NRC; gli Stati Uniti Anche il Census Bureau e altre agenzie governative impegnate nel collegamento dei dati hanno dovuto sviluppare mezzi per affrontare questo problema.⁵¹ Una soluzione prevede l'applicazione di metodi di rumore aggiuntivo ai dati (privacy differenziale) al fine di offuscare i dati individuali preservando l'utilità dei dati per la ricerca. Ne discuteremo più dettagliatamente insieme ad altre tecnologie che migliorano la privacy nel capitolo 7. Le tecnologie non ⁵² Tuttavia, il miglioramento della privacy sono una panacea e, a seconda della natura del particolare set di dati, gli obiettivi di garantire l'anonimizzazione, consentendo anche ai ricercatori di accedere a dati granulari, possono entrare in conflitto.

L'NRC può attingere anche dai dati delle "Cinque Casseforti". framework di sicurezza utilizzato da UK Data Service⁵³, Federal Statistical Research Data Centers Network e Coleridge Initiative, un modello incentrato su dati, progetti, persone, impostazioni di accesso e output.⁵⁴ L'attuazione dell'Evidence Act del 2019 sta già utilizzando un quadro Five Safes simile per prendere decisioni in merito al collegamento dei dati.

contratti attraverso i quali sono stati raccolti. Allo stesso modo, l'NRC potrebbe controllare l'ambito di diffusione di modelli, codice e dati, a seconda della sensibilità. È meno probabile che l'identificabilità teorica costituisca un problema quando l'accesso e la diffusione sono limitati e i dati sono di minore importanza natura sensibile o non riguarda affatto gli individui.⁵⁶

Facilitare l'accesso dei ricercatori con un livello Modello di accesso

In che modo i ricercatori dovrebbero ottenere l'accesso a dati specifici risorse? Attualmente, l'approvazione procede a livello di agenzia per agenzia.⁵⁷ Così come il valore dell'NRC per il sostegno alla ricerca sull'IA dipenderà in parte dalla misura in cui sarà in grado di riunire set di dati di diverse agenzie, dipenderà anche dalla misura in cui che può semplificare il processo di accesso ai dati. Un modo per raggiungere questa semplificazione sarà attraverso un sistema di accesso a più livelli per gli utenti NRC, simile al sistema a più livelli di FedRAMP per l'archiviazione dei dati federali sul cloud, dove livelli più alti consentirebbero l'accesso a dati ad alto rischio, soggetti agli altri requisiti di calcolo e utilizzo dei dati. Discutiamo questo sistema di accesso in modo più approfondito nel Capitolo 7.

Il capitolo 2 ha affermato che l'accesso al calcolo dovrebbe iniziare con i PI nelle istituzioni accademiche. Questa autorizzazione può anche fungere da riferimento, in cui tutti i PI registrati nell'NRC possono accedere liberamente e utilizzare set di dati a basso rischio sull'NRC. Livelli aggiuntivi imporrebbero più requisiti, come cittadinanza, nulla osta di sicurezza, restrizioni di distribuzione o restrizioni di calcolo e di sistema. Questi livelli di accesso saranno simili a quelli utilizzati per determinare la classificazione FedRAMP per l'archiviazione dei dati, ma mentre l'accesso e la sensibilità dell'archiviazione possono richiamare considerazioni simili; potrebbero non essere necessariamente gli stessi.

CASO DI STUDIO: INIZIATIVA COLERIDGE (RICERCA DATI AMMINISTRATIVI FACILITÀ)

In collaborazione con il Census Bureau e il finanziamento dell'Ufficio of Management and Budget, la Coleridge Initiative, un'organizzazione senza scopo di lucro, ha lanciato l'Administrative Data Research Facility (ADRF), una piattaforma informatica sicura per le agenzie governative per condividere e lavorare con i microdati delle agenzie.⁵⁸ L'ADRF è disponibile sul Federal Risk and Budget Programma di gestione delle autorizzazioni (FedRAMP) Marketplace e ha una certificazione FedRAMP Moderate. Attualmente, la piattaforma supporta oltre 100 set di dati di 50 agenzie.⁵⁹

L'ADRF fornisce l'accesso gratuito ai ricercatori sponsorizzati dall'agenzia e ai ricercatori affiliati all'agenzia che seguono i programmi di formazione dell'ADRF. Negli ultimi tre anni, oltre 500 dipendenti di circa 100 agenzie hanno seguito i programmi di formazione ADRF.⁶⁰

L'ADRF fornisce uno spazio di lavoro condiviso per i progetti e il Data Explorer, uno strumento per visualizzare una panoramica e i metadati (nome, descrizione del campo e tipo di dati) dei set di dati disponibili sull'ADRF.⁶¹ Per accedere ai dati riservati, gli utenti devono soddisfare la revisione requisiti stabiliti dall'agenzia che fornisce i dati. Per esportare i dati, gli utenti devono passare attraverso un unico processo di "revisione dell'esportazione".⁶² L'ADRF ha un processo di revisione predefinito molto complesso, che richiede ai ricercatori di inviare tutto il codice e l'output per il progetto per l'approvazione all'amministratore dei dati e genera costi aggiuntivi, se si richiede l'esportazione di più di 10 file.⁶³ L'agenzia che fornisce i dati può anche modificare il processo di revisione predefinito, se lo desidera.

Prima di trasferire i file di dati, l'ADRF fornisce un'applicazione per hashing dei dati per trasmettere dati in modo sicuro.⁶⁴ L'ADRF segue anche il modello di sicurezza "Five Safes" utilizzato da altre agenzie governative, come il servizio dati del Regno Unito.⁶⁵

La gestione dei dati per l'ADRF è definita in conformità al Titolo III dell'Evidence-Based Policymaking Act del 2018.⁶⁶ Una volta che un set di dati ristretto viene condiviso con l'ADRF, una persona all'interno dell'agenzia verrà assegnata al data steward per tutte le richieste di progetto. Da lì, le procedure vengono sviluppate con l'agenzia, in termini di aspettative su come i dati saranno protetti, utenti autorizzati e procedure di controllo per la continua conformità.

I data steward hanno accesso a un portale online nell'ADRF. Tutto le richieste di progetto relative a dati specifici vengono indirizzate all'amministratore dei dati tramite questa proposta. Una volta concesso l'accesso, l'amministratore dei dati ha anche la possibilità di monitorare la conformità del progetto.

PUNTI CHIAVE

- **Saldi che forniscono coerente limitato accesso ai dati con i requisiti dell'agenzia:** l'ADRF bilancia la creazione di ogni accesso limitato al set di dati ed esportare il modulo di revisione in modo coerente in il portale dei dati con i requisiti dell'agenzia per l'accesso e l'esportazione dei dati. Ciò consente alle agenzie di controllare l'accesso ai propri dati, fornendo al tempo stesso un unico punto di accesso per i ricercatori. Attualmente, la piattaforma supporta solo dati coerenti con una FedRAMP Moderate certificato.
- **Standardizzato per entrambi gli utenti e data steward:** insieme a ogni accesso ai dati per utenti, anche un punto di contatto presso l'agenzia che fornisce i dati ha accesso alla piattaforma. Ciò consente un facile accesso per approvare e tenere traccia dei progetti e lavorare con l'ADRF sull'accesso requisiti.
- **Il framework per la sicurezza dei dati di Five Safes:** l'ADRF garantisce la sicurezza dei dati concentrandosi su cinque aspetti: dati, progetti, persone, impostazioni di accesso e risultati.
- **La formazione come funzione centrale:** l'ADRF ospita workshop e forma dipendenti governativi e altri ricercatori sull'uso dei dati.

I modelli esistenti per l'accesso dei ricercatori a set di dati sensibili possono aiutare a tracciare un quadro di come l'NRC potrebbe mantenere e monitorare un sistema di accesso a più livelli. L'NRC può emulare sia la Coleridge Initiative che il Center for Population Health Sciences (PHS) di Stanford,⁶⁷ ad esempio, che fungono da intermediari di dati, facilitando l'accesso ai dati governativi. In effetti, questi intermediari sono stati documentati come mezzi efficaci per superare le barriere alla condivisione dei dati perché, in sostanza, negoziano e semplificano le relazioni tra i contributori di dati e gli utenti.⁶⁸ Ad esempio, in qualità di intermediario fidato, l'NRC potrebbe centralizzare l'assunzione di DUA processo promulgando un modulo standard universale per i DUA di agenzia.⁶⁹

Inoltre, analogamente all'esempio dell'iniziativa Coleridge, uno o più rappresentanti designati all'interno dell'agenzia potrebbero essere assegnati come amministratore dei dati per tutte le richieste di progetto per un determinato set di dati limitato. Qualsiasi progetto che richieda l'accesso ai dati in livelli superiori potrebbe iniziare solo dopo che la sua proposta è stata esaminata e approvata da un rappresentante competente. Poiché l'accesso agli NRC inizia con i PI, i ricercatori dovrebbero anche ottenere l'approvazione dai loro Institutional Review Boards (IRB) universitari, se necessario. Dopo l'approvazione del progetto e l'autorizzazione del ricercatore dell'NRC, i dati sarebbero resi disponibili attraverso il portale protetto dell'NRC. Eventuali violazioni dei termini di utilizzo o della privacy del soggetto potrebbero comportare sanzioni che vanno dalla retrocessione del livello di accesso alla rimozione dei privilegi NRC o sanzioni professionali, civili o penali, a seconda dei casi.

Il NRC può emulare sia il Iniziativa Coleridge e Stanford Centro per la popolazione Scienze della salute (PHS), per esempio, che fungono da dati intermediari, nel facilitare l'accesso ai dati governativi.

CASO DI STUDIO: STANFORD CENTRO PER LA POPOLAZIONE SCIENZE DELLA SALUTE

Lo Stanford Center per la salute della popolazione Sciences (PHS) fornisce una serie crescente di set di dati relativi alla salute della popolazione e metodi di accesso a Stanford ricercatori e affiliati.⁷⁰ I dati PHS L'ecosistema ospita set di dati di alto valore, collegamenti e filtri di dati e strumenti analitici per aiutare i ricercatori. Il PHS collabora con un'ampia gamma di enti pubblici, senza scopo di lucro e privati per concedere in licenza set di dati a livello di popolazione per ricercatori universitari, che vanno da set di dati pubblici a basso rischio a dati riservati contenenti informazioni sanitarie protette (PHI) e informazioni di identificazione personale (PII), come Medicare, reclami commerciali come Optum e MarketScan data⁷¹ e cartelle cliniche elettroniche.

PUNTI CHIAVE

- **Architettura dei dati mista, ma esperienza utente coerente:** PHS utilizza un mix di servizi dati locali e cloud, ma cerca comunque di fornire un'esperienza utente coerente.
- **L'accesso limitato ai dati ha un unico punto di ingresso:** il PHS Data Portal standardizza, centralizza e semplifica i requisiti di accesso ai dati e la formazione, invece di indirizzare gli utenti verso un processo che richiede tempo lavorando direttamente con ciascun data steward.
- **Riduce i costi e il tempo associati a approvvigionamento di dati:** PHS sfrutta le relazioni esistenti con le agenzie per consolidare i set di dati in un unico portale, risparmiando ai ricercatori il tempo e il denaro necessari per ottenere l'accesso attraverso le richieste delle singole agenzie.

CASO STUDIO: STANFORD CENTRO PER LA POPOLAZIONE SCIENZE DELLA SALUTE (CONTINUA)

Oltre all'archiviazione sicura dei dati e agli strumenti computazionali per i ricercatori, PHS fornisce standardizzati e protocolli di accesso e gestione dei dati ben documentati, che aumentano il comfort del proprietario dei dati con la condivisione dei dati. PHS ha anche personale a tempo pieno che coltiva e mantiene i rapporti con le organizzazioni che detengono i dati. Ciò consente a PHS di collaborare con questi gruppi per centralizzare l'hosting dei dati e fornire un accesso sicuro a un'ampia gamma di ricercatori.

Il PHS Data Portal è ospitato su una piattaforma di terze parti che consente la scoperta e l'esplorazione dei dati e passaggi standardizzati chiaramente delineati per l'accesso ai dati. La piattaforma di terze parti, Redivis, utilizza un sistema di accesso a quattro livelli: (1) panoramica dei dati e documentazione di base; (2) accesso ai metadati, incluse definizioni, descrizioni e caratteristiche; (3) un campione dell'1 o del 5% del set di dati; e (4) accesso completo ai dati.⁷²

Se i dati sono classificati come pubblici, i ricercatori possono accedervi utilizzando un software specializzato o semplicemente scaricarli direttamente.⁷³ Per i dati riservati, il portale dispone di moduli integrati per richiedere facilmente l'accesso.⁷⁴ Dopo aver identificato il set di dati, il ricercatore deve richiedere l'adesione a l'organizzazione che ospita i dati.⁷⁵ Un amministratore dell'organizzazione proprietaria del set di dati può stabilire i requisiti relativi ai membri e allo studio che devono essere soddisfatti, inclusa la formazione e la qualificazione istituzionale, per poter accedere ai dati. Le applicazioni dei membri possono essere impostate per l'approvazione automatica o richiedere l'approvazione amministrativa. Una volta concesso l'accesso a un set di dati, i ricercatori possono manipolare i dati utilizzando un software specializzato. Le restrizioni di utilizzo vengono anche specificate individualmente su ciascun set di dati per controllare se è possibile esportare output completo, parziale o nullo e quale livello di revisione è richiesto per l'esportazione. Tutte le richieste di dati ed export vengono gestite direttamente sulla piattaforma Data Explorer.

Attualmente, il PHS Data Portal è principalmente per la facoltà, il personale, gli studenti o altri affiliati di Stanford.⁷⁶ Addirittura con lo stato di affiliazione, alcuni set di dati commerciali potrebbero richiedere ulteriori accordi di data rider per l'accesso. I collaboratori non di Stanford devono soddisfare tutti gli stessi requisiti di accesso degli affiliati di Stanford, oltre a eventuali requisiti imposti dalla propria istituzione. Inoltre, è spesso necessario un accordo "data rider" sul DUA originale.⁷⁷

Per lavorare con dati riservati, il PHS fornisce due servizi informatici per dati ad alto rischio: (1) Nero, con entrambe le versioni della piattaforma on-premise e Google Cloud Platform (GCP); e (2) cluster PHS-Windows Server.⁷⁸ Entrambi sono gestiti dallo Stanford Research Computing Center (SRCC). Entrambi i servizi sono conformi a HIPAA.⁷⁹ I dati illimitati possono essere utilizzati su qualsiasi altro ambiente computazionale di Stanford (Sherlock, Oak) o semplicemente scaricati sulla macchina locale del ricercatore.

PROMUOVERE L'INTERAGENZIA ARMONIZZAZIONE E ADOZIONE DI MODERNO ACCESSO AI DATI STANDARD

Il panorama federale della condivisione dei dati ne risente standard e pratiche divergenti e le singole agenzie, lasciate sole, hanno tradizionalmente affrontato ostacoli elevati per armonizzare e modernizzare i loro standard di accesso ai dati. È problematico sia dal punto di vista dell'agenzia che della società. Come rileva un rapporto della Conferenza amministrativa degli Stati Uniti da un'indagine sull'uso dell'IA nel governo federale, quasi la metà delle agenzie ha sperimentato l'IA per migliorare le capacità decisionali e operative, ma spesso mancano dell'infrastruttura tecnica e della capacità di dati utilizzare le moderne tecniche e strumenti di intelligenza artificiale.⁸¹ La mancanza di uno standard moderno e uniforme per la condivisione dei dati nella ricerca sull'IA, pertanto, rende più difficile per le agenzie realizzare guadagni in termini di accuratezza, efficienza e responsabilità, che successivamente si ripercuotono sui cittadini a valle, che sono influenzati dalle decisioni dell'agenzia.⁸²

La mancanza di uno standard moderno e uniforme per la condivisione dei dati in

La ricerca sull'intelligenza artificiale lo rende più difficile

affinché le agenzie realizzino guadagni
in termini di accuratezza, efficienza e
responsabilità, che successivamente incidono sui
cittadini a valle, che sono influenzati dalle decisioni
dell'agenzia.

L'NRC può aiutare a superare la riluttanza delle agenzie a condividere i dati consentendo l'accesso alle agenzie per calcolare i propri dati. Questo risolverebbe almeno due cruciali problemi per gli enti pubblici. In primo luogo, l'accesso alle risorse informatiche collettive dell'NRC supererebbe alcune difficoltà che le agenzie hanno tradizionalmente affrontato nella creazione delle proprie risorse informatiche.⁸³

competenza del governo in materia di IA.⁸⁴ Dal punto di vista della società, ciò potrebbe aumentare le capacità del governo nell'adozione responsabile dell'IA, contribuire a ridurre il costo delle funzioni di governance fondamentali e aumentare l'efficienza, l'efficacia e la responsabilità dell'agenzia.⁸⁵

L'NRC può anche imparare dagli altri e allinearsi con gli altri iniziative per armonizzare e modernizzare gli standard. L'Evidence Act, che richiede alle agenzie di nominare responsabili dei dati e della valutazione, ne è un esempio. La legislazione che autorizza la creazione dell'NRC potrebbe prevedere un mandato federale per incoraggiare l'adozione di pratiche di condivisione.⁸⁶ Tuttavia, come si discute nel capitolo 5, un mandato federale da solo, senza ulteriori aiuti o incentivi, potrebbe non essere sufficiente per incentivare l'armonizzazione degli standard di accesso e condivisione dei dati.⁸⁷ La task force dovrebbe pertanto prendere in considerazione l'idea di unire il mandato con ulteriori vantaggi, come la fornitura di finanziamenti per assistere le agenzie nell'espansione delle loro capacità tecniche o del personale a sostegno dell'NRC e della strategia nazionale sull'IA. L'NRC è in linea con l'attuale caso bipartisan per il National Secure Data Service (NSDS) (descritto nel caso di studio di seguito), un servizio che faciliterebbe l'accesso dei ricercatori ai dati con maggiore privacy e trasparenza, raccomandato dalla Commission on Evidence-Based Definition delle politiche nel 2018. Sia l'NRC che l'NSDS sono iniziative complementari di condivisione dei dati che hanno il potenziale per migliorare notevolmente l'efficacia operativa del servizio pubblico. Elaboriamo ulteriormente la proposta NSDS nel capitolo 5. Infine, i programmi di formazione sono strade promettenti per aumentare l'adozione dell'NRC e il supporto dell'agenzia. Ad esempio, come descritto nel caso di studio sopra, la Coleridge Initiative ha ospitato workshop per formare oltre 500 dipendenti di circa 100 agenzie sull'uso dei dati negli ultimi tre anni.

CASO DI STUDIO: LA LEGGE SULLE PROVE

Nel perseguimento di un accesso più ampio e più sicuro e di un collegamento dei dati amministrativi del governo, nel marzo 2016 il Congresso ha istituito una Commissione bipartisan per l'elaborazione di politiche basate sulle prove. Il rapporto finale della commissione⁸⁸ includeva 22 raccomandazioni al governo federale affinché costruisse meccanismi,⁸⁹ e capacità istituzionale per fornire un accesso sicuro ai dati pubblici a fini statistici e di ricerca. Una raccomandazione era quella di creare un "National Secure Data Service" (NSDS) per facilitare l'accesso ai dati allo scopo di costruire prove, pur mantenendo la privacy e la trasparenza. Attraverso questo servizio, l'NSDS potrebbe aiutare i ricercatori collegando temporaneamente i dati esistenti e fornendo un accesso sicuro, senza creare essa stessa una stanza di smistamento dei dati.

Il Foundations for Evidence-Based Policymaking Act del 2018⁹⁰ ha creato parte della base legislativa per le raccomandazioni della commissione. In particolare, ha creato nuovi ruoli per i responsabili dei dati, della valutazione e della statistica e ha cercato di aumentare l'accesso e il collegamento di set di dati precedentemente nell'ambito del Confidential Information Protection and Statistical Efficiency Act (CIPSEA).⁹¹

Infine, la Strategia federale in materia di dati 2020 e l'azione associata Plan⁹² ha cercato di attuare tali disposizioni legislative. La strategia includeva piani per migliorare la governance dei dati, per rendere i dati più accessibili, per migliorare l'uso dei dati da parte del governo e per aumentare l'uso e la qualità degli inventari dei dati, dei metadati e della sensibilità dei dati.

Il passaggio centrale rimanente previsto dall'Evidence-Based iniziale Policymaking Commission è un National Secure Data Service (NSDS) modellato sul Data Service del Regno Unito.⁹³ Il Data Service del Regno Unito fornisce accesso a una serie di sondaggi pubblici, studi longitudinali, dati del censimento del Regno Unito, dati aggregati internazionali, dati aziendali e dati qualitativi. Oltre all'accesso, fornisce indicazioni e formazione per l'utilizzo dei dati, sviluppa best practice e standard per la privacy e dispone di personale specializzato che applica tecniche di controllo statistico per fornire l'accesso a dati troppo dettagliati, sensibili o riservati per essere resi disponibili con licenze standard.

PUNTI CHIAVE

- **Condivisione prioritaria dei dati in presenza di un caso operativo di servizio pubblico:** esiste un precedente in iniziative di condivisione di dati amministrativi su larga scala giustificate da miglioramenti dell'efficienza e dell'efficacia operativa del servizio pubblico.
- **Iniziativa National Secure Data Service (NSDS):** l'NSDS è sostenuto dal supporto bipartisan.
- **Modelli istituzionali che bilanciano talento esterno, innovativo e influenza interna e interagenzia:** A A livello federale Ricerca finanziata e Il modello Development Center (FFRDC), ospitato all'interno di un'agenzia esistente (NSF), può bilanciare la capacità di portare talenti esterni con l'influenza dell'agenzia interna.

INVESTIMENTO IN SEQUENZA IN PATRIMONIO DI DATI

Dati i notevoli ostacoli nella negoziazione dell'accesso ai dati, l'NRC dovrà sequenziare strategicamente su quali agenzie e set di dati concentrarsi per l'uso da parte dei ricercatori. Il governo federale raccoglie petabyte di dati,⁹⁴ ciascuno con vari gradi di restrizione o apertura. Nel valutare a quali set di dati dare la priorità, l'NRC può attingere dall'esempio di altre iniziative di condivisione dei dati, nonché concentrarsi su set di dati a breve termine che non pongono sfide complesse per quanto riguarda la privacy o la condivisione dei dati. Un esempio del settore privato è Google Earth Engine, che ha aggregato petabyte (circa 1 milione di gigabyte) di immagini satellitari e set di dati geospaziali, quindi ha collegato tale accesso ai servizi di cloud computing di Google per consentire agli scienziati di rispondere a una serie di domande di ricerca cruciali.⁹⁵ Questo processo di aggregazione di dati complessi e di ospitarli in un'infrastruttura informatica amichevole per facilitare la ricerca, dimostra il valore convincente dell'accoppiamento di calcolo e dati. Come altro esempio, ADR UK identifica specifiche aree di ricerca che sono di pressante interesse politico, come il "mondo del lavoro"⁹⁶ e dà la priorità all'accesso ai dati per i ricercatori che lavorano su tali argomenti. Il servizio dati del Regno Unito offre set di dati derivati da sondaggi, fonti amministrative e relative alle transazioni, inclusi dati sulla produttività dall'Annual Respondents Database,⁹⁷ dati sull'innovazione dall'UK Innovation Survey,⁹⁸ dati geospaziali dall'indagine sulla forza lavoro,⁹⁹ Understanding Society,¹⁰⁰ e dati sensibili sullo sviluppo dell'infanzia.¹⁰¹

Quando si assegna la priorità ai set di dati e alle agenzie per la partnership NRC, si consigliano i seguenti criteri:

- Dati preziosi per i ricercatori di intelligenza artificiale, ma attualmente non disponibili in una forma conveniente. Ad esempio, in una richiesta di commenti del luglio 2019, l'Office of Management and Budget (OMB) ha chiesto ai membri del pubblico di fornire un contributo sulle caratteristiche dei modelli che li rendono adatti alla ricerca e sviluppo dell'IA, quali dati sono attualmente limitati e come la liberazione di tali dati accelererebbe la ricerca e lo sviluppo di IA di alta qualità.¹⁰² In una risposta, la Data Coalition ha sostenuto che il rilascio controllato di informazioni private ma

i dati in data.gov sarebbero preziosi per la ricerca.¹⁰³ La Data Coalition ha anche esortato le agenzie a prendere in considerazione il rilascio di set di dati grezzi e non strutturati, come i registri dei call center delle agenzie, le richieste e i reclami dei consumatori, nonché le ispezioni normative e i rapporti investigativi.¹⁰⁴ Un altro esempio dei dati a cui è attualmente difficile accedere, ma che sono una questione di pubblico dominio, sono atti giudiziari elettronici ospitati in un sistema dall'ufficio amministrativo dei tribunali statunitensi.¹⁰⁵

- Dati conservati all'interno di agenzie che hanno l'autorità legale di condividere i dati e/o che hanno precedenti esperienze di condivisione dei dati. Il Census Bureau, ad esempio, dispone di maggiori collegamenti legali tra agenzie esistenti rispetto ad altre agenzie e dispone di una notevole esperienza interna di analisi dei dati . disoccupazione, retribuzione e condizioni di lavoro, prezzi e condizioni di vita) con i ricercatori.¹⁰⁷
- Dati con implicazioni limitate sulla privacy. Ad esempio, le agenzie i cui dati riguardano fenomeni naturali, piuttosto che individui, possono essere più facili da gestire dal punto di vista della privacy, ad esempio la NASA, il servizio geologico degli Stati Uniti e la National Oceanic and Atmospheric Administration. Set di dati come quelli ospitati nel Planetary Data System della NASA,¹⁰⁸ ma che non sono facilmente disponibili per i ricercatori, possono servire come valido punto di partenza per l'NRC. Aumentare la disponibilità e l'interoperabilità dei set di dati di queste agenzie farebbe avanzare la missione principale dell'NRC e potrebbe essere fatto senza mettere a repentaglio la privacy individuale.

Capitolo 4:

Progettazione organizzativa

Quale forma istituzionale dovrebbe assumere il CNR? Due considerazioni generali sono: (1) facilità di accesso ai dati; e (2) facilità di coordinamento con le risorse di calcolo.¹ Come discusso nel Capitolo 3 e approfondito nel Capitolo 5, il panorama federale della condivisione dei dati tra le agenzie è altamente frammentato, con molte agenzie riluttanti o legalmente costrette a condividere i propri dati. L'NRC dovrà coordinarsi tra le entità che forniscono l'infrastruttura di calcolo e gli stessi ricercatori. Poiché l'obiettivo dell'NRC è fornire ai ricercatori l'accesso ai dati del governo e una potenza di calcolo ad alte prestazioni, l'uno senza l'altro non riuscirà a raggiungere la missione dell'NRC.

Sulla base di un ampio lavoro a sostegno dell'Evidence Act, raccomandiamo l'uso di centri di ricerca e sviluppo finanziati a livello federale (FFRDC) e partenariati pubblico-privato (PPP) come possibili forme organizzative per l'NRC. Raccomandiamo la creazione di un FFRDC presso le agenzie governative affiliate a breve termine, poiché riteniamo che questo percorso consenta la facilitazione più semplice sia dell'infrastruttura di calcolo che dell'accesso ai dati del governo. A lungo termine, l'istituzione di un PPP potrebbe facilitare una maggiore condivisione e accesso ai dati tra il settore pubblico e quello privato. È importante sottolineare che altre opzioni includono la creazione di un'agenzia o ufficio federale completamente nuovo all'interno di un'agenzia esistente. Anche se queste opzioni potrebbero semplificare il coordinamento con le risorse di calcolo, entrambe pongono sfide per quanto riguarda l'accessibilità dei dati e la condivisione dei dati tra agenzie.

RICERCA FINANZIATA FEDERALMENTE E CENTRO DI SVILUPPO

Le FFRDC sono società senza scopo di lucro quasi governative sponsorizzate da un'agenzia federale ma gestite da appaltatori, tra cui università, altre organizzazioni senza scopo di lucro e aziende del settore privato.² Il modello FFRDC conferisce i vantaggi di uno stretto rapporto di agenzia, insieme all'amministrazione indipendente, nel facilitare l'accesso ai dati. A causa dei loro stretti rapporti di subappalto con la loro agenzia madre, tutti gli FFRDC beneficiano di un accesso ai dati che va "oltre ciò che è comune al normale rapporto contrattuale, ai dati del governo e dei fornitori, inclusi i dati sensibili e proprietari".

Un recente rapporto di Hart e Potok sul National Secure Data Service (NSDS) (si veda lo studio di caso nel capitolo 3) supporta anche il modello FFRDC come un modo ottimale per facilitare l'accesso e il collegamento dei dati amministrativi del governo . in un'agenzia esistente e lo sviluppo di un'università

CHIAVE ASPORTO

- A breve termine, l'NRC dovrebbe essere istituito come a Centro di ricerca e sviluppo finanziato a livello federale (FFRDC), che ridurrebbe i costi significativi per la protezione dei dati dalle agenzie federali.
- A lungo termine, un partenariato pubblico-privato (PPP) ben progettato, governato da funzionari di Affiliated Agenzie governative, ricercatori accademici e rappresentanti del settore tecnologico potrebbero aumentare la quantità e la qualità della R&S e ridurre i costi di manutenzione.
- Istituire l'NRC come agenzia federale autonoma o l'ufficio avrebbe dovuto affrontare numerosi sfide, in particolare nella messa in sicurezza l'accesso ai dati ospitati in altri agenzie.

guidato, servizio di condivisione dei dati, ma il rapporto alla fine ha raccomandato il modello FFRDC per diversi motivi. UN FFRDC può scalare rapidamente, perché può accedere ai dati del governo e ai talenti di alta qualità più facilmente rispetto ad altre opzioni.⁵ Un FFRDC può anche sfruttare le competenze governative esistenti. La NSF, ad esempio, sponsorizza già cinque FFRDC separati e ha una vasta esperienza nella coltivazione e nel mantenimento di reti di ricercatori.⁶

Tuttavia, il modello FFRDC ne ha alcuni limitazioni. In primo luogo, il ruolo di un FFRDC è limitato alla ricerca e allo sviluppo per la propria agenzia di sponsorizzazione che "è strettamente associata all'esecuzione di funzioni intrinsecamente governative" .

In secondo luogo, il successo di un modello FFRDC per l'NRC dipenderà dalla capacità dell'agenzia promotrice di ottenere la cooperazione attraverso il governo federale per fornire i dati necessari per la ricerca. Un modo per farlo sarebbe che più agenzie co-sponsorizzassero l'FFRDC, riducendo l'attrito contrattuale per i set di dati.⁸ Un'altra opzione sarebbe quella di creare più FFRDC ospitati in diverse agenzie, incentivando ciascuna di queste agenzie a condividere i propri dati con il rispettivo FFRDC . Un esempio analogo potrebbe includere i Laboratori Nazionali come una rete, dove ciascun Laboratorio Nazionale sarebbe un'istanza dell'NRC all'interno la propria agenzia competente.⁹

In terzo luogo, più FFRDC richiederebbero processi separati per le risorse di calcolo. A breve termine, l'NRC può alleviare questo problema stipulando contratti per crediti cloud commerciali, che è probabilmente già la soluzione a breve termine per l'NRC per fornire l'accesso al calcolo. Come discusso in precedenza, i fornitori di cloud privato hanno già una vasta esperienza nella fornitura di risorse di calcolo al governo¹⁰ e alle istituzioni accademiche.¹¹ La familiarità con questi fornitori di cloud privato può ridurre l'attrito nell'allocazione del calcolo tra i ricercatori di più FFRDC.

A lungo termine, il modello FFRDC potrebbe non essere il più efficiente. Dal punto di vista dei costi e della sostenibilità, gli FFRDC hanno tradizionalmente sofferto in modo significativo

superamenti, in quanto "operano sotto un mosaico inadeguato e incoerente di controlli federali sui costi, sulla contabilità e sulla revisione contabile, le cui carenze hanno contribuito all'uso dispendioso o inappropriato di milioni di dollari federali".¹² Un'altra preoccupazione è che, storicamente, l'infrastruttura FFRDC non ha stato regolarmente aggiornato. Un rapporto del Dipartimento dell'Energia del 2017 ha evidenziato che l'infrastruttura FFRDC era inadeguata per soddisfare la missione.¹³ L' ispettore generale della NASA ha anche evidenziato che oltre il 50% delle apparecchiature del Jet Propulsion Laboratory (un FFRDC della NASA) aveva almeno 50 anni.¹⁴ Se un FFRDC versione dell'NRC incontra queste stesse sfide, raccomandiamo che l'NRC, a lungo termine, passi a un modello di partenariato pubblico-privato.

CASO DI STUDIO: SCIENZA E TECNOLOGIA ISTITUTO POLITICO (STPI)

STPI è un FFRDC istituito dal Congresso nel 1991 per fornire obiettivi rigorosi consulenza e analisi all'Office of Science and Technology Policy e ad altre agenzie del ramo esecutivo.¹⁵ STPI è gestito dall'Institute for Defense Analysis (IDA), un'organizzazione senza scopo di lucro che gestisce anche altri due FFRDC: il Systems and Analyses Center e il Center for Comunicazioni e informatica.¹⁶ IDA non ha altre linee di attività al di fuori del quadro FFRDC.¹⁷

Il principale sponsor federale di STPI è la National Science Foundation, ma la ricerca presso STPI è anche co-sponsorizzata da altre agenzie federali, tra cui il National Institute of Health (NIH), il Department of Energy (DOE), il Department of Transportation (DOT), il Department of Defense (DOD) e Department of Health and Human Services (HHS).¹⁸ A causa della "relazione unica" tra un FFRDC e i suoi sponsor, STPI "gode di un accesso insolito a informazioni governative e aziendali altamente riservate e riservate".¹⁹

Gli stanziamenti NSF forniscono la maggior parte dei finanziamenti per STPI, inclusi 4,7 milioni di dollari nell'anno fiscale 2020,²⁰ ma un importo limitato di finanziamenti viene fornito anche da altre agenzie federali.²¹ STPI ha circa 40 dipendenti a tempo pieno e ha accesso all'esperienza di circa IDA 800 altri dipendenti.²² In qualità di FFRDC, STPI può anche stipulare contratti di consulenza, come richiesto per un particolare progetto.²³ Lo statuto che specifica i doveri di STPI le impone inoltre di consultare ampiamente i rappresentanti dell'industria privata, del mondo accademico e delle istituzioni senza scopo di lucro e di incorporare quelle opinioni nel lavoro di STPI nella misura massima praticabile.²⁴

Lo STPI è inoltre tenuto a presentare al presidente una relazione annuale sulle sue attività, in conformità con i requisiti prescritti dal presidente,²⁵ che fornisce un'ulteriore responsabilità per l'FFRDC. Secondo il rapporto 2020 di STPI, STPI ha collaborato con più agenzie federali, supportandole su 48 analisi di politiche tecnologiche separate per tutto il 2020.²⁶

PUNTI CHIAVE

- **Più co-sponsor di agenzie:** While Lo sponsor principale di STPI è il Scienza Nazionale Fondazione, un certo numero di altre le agenzie co-sponsorizzano anche STPI, riducendo le difficoltà di accesso ai dati tra le agenzie.
- **Competenza:** Mentre STPI è dotato di propri dipendenti, può anche attingere alle competenze di centinaia di dipendenti presso il Istituto per la Difesa Analyses (IDA), l'organizzazione che gestisce STPI. In qualità di FFRDC, STPI può anche contrattare competenze aggiuntive necessarie.

UN PARTENARIATO PUBBLICO-PRIVATO (PPP)

Un partenariato pubblico-privato (PPP) creerebbe una partnership tra agenzie federali e organizzazioni del settore privato per ospitare e gestire congiuntamente gli sforzi di condivisione dei dati e gestire l'infrastruttura di calcolo. Poiché diverse agenzie e membri del settore privato possono avere preferenze contrattuali, obiettivi di proprietà intellettuale e garanzie di sicurezza diverse per l'accesso ai dati, la creazione di una partnership per la condivisione dei dati all'interno di questo quadro patchwork potrebbe essere difficile nell'immediato futuro. Tuttavia, i PPP possono fornire una serie di vantaggi a lungo termine, come hanno fatto

sono stati utilizzati con successo come data clearinghouses per produrre, analizzare e condividere dati tra il settore pubblico e privato.²⁷ In effetti, riconoscendo i vantaggi del modello PPP, l'Unione Europea ha lanciato una nuova iniziativa denominata Public Private Partnerships for Big Data che offrire un ambiente sicuro per la collaborazione intersettoriale e sperimentazione utilizzando sia dati commerciali che pubblici.²⁸ In generale, i PPP per la condivisione dei dati possono aumentare la qualità e la quantità della R&S, aumentare il valore e l'efficienza della condivisione dei dati del settore pubblico e ridurre i costi a lungo termine necessari per gestire e mantenere la infrastruttura di condivisione dei dati.²⁹

CASO DI STUDIO: ALBERTA DATA PARTNERSHIPS (ADP)

Fondato nel 1997, ADP PPP è progettato per fornire una gestione a lungo termine di set di dati digitali completi per il mercato dell'Alberta. custode" dei dati governativi e Altalis è "l'operatore". ADP crea software per caricare e distribuire in modo sicuro questi set di dati spaziali provinciali agli utenti. Altalis fornisce inoltre formazione agli utenti finali ed è responsabile della pulizia, dell'aggiornamento e della standardizzazione dei set di dati.³³

Nella scelta del proprio "partner operativo" (ovvero Altalis) per la joint venture, il consiglio di amministrazione di ADP ha inizialmente emesso una "Richiesta di informazioni" che sollecitava proposte da parte di aziende del settore privato il cui core business era il miglioramento, la manutenzione, la gestione e la distribuzione di dati spaziali.³⁴ Il consiglio di amministrazione di ADP alla fine scelse Altalis, non solo perché disponeva di un'offerta superiore e delle capacità esistenti, ma anche perché Altalis era disposta ad assumersi tutti gli investimenti necessari, a proprio rischio, per costruire e gestire il sistema ADP in conformità alle specifiche ADP.³⁵

Oggi, tutti i costi di Altalis e ADP sono coperti dalle operazioni del giunto. La joint venture ottiene entrate attraverso, ad esempio, il finanziamento diretto di progetti e le tariffe per l'accesso ai dati da parte delle parti interessate, che includono comuni, agenzie di regolamentazione, organizzazioni energetiche, forestali e minerarie. /20 tra Altalis e ADP, rispettivamente, e ADP successivamente utilizza la sua quota di profitto per reinvestire in dati e miglioramenti del sistema.³⁸

L'ADP PPP afferma di aver generato efficienze per la condivisione dei dati. Ad esempio, l'ADP stima che un approccio tradizionale esclusivamente governativo al mantenimento e alla distribuzione di set di dati sarebbe stato compreso tra 65 milioni e 120 milioni di dollari cumulativamente dall'inizio di ADP, e ADP afferma di aver fornito ai propri utenti 6,8 milioni di dollari di risparmi sui costi.³⁹

PUNTI CHIAVE

■ *Utilizzo di joint venture:*

L'ADP PPP utilizza un accordo di joint venture per stabilire i termini e le condizioni per la partnership, per cui un partner è il custode per dati del governo, e l'altro è il operatore che costruisce e mantiene il software che facilita i dati condivisione.

■ *Accordi di*

compartecipazione alle entrate:

assicura un modello di compartecipazione alle entrate contribuiti da e realizzazione di valore a ciascun stakeholder. Successivamente l'ADP reinvestirà i suoi profitti per migliorare il sistema.

■ *Efficienze significative:*

Secondo il ADP, ci sono costi inferiori per creare e mantenere il file ADP che sotto a approccio convenzionale.

Un modello PPP potrebbe ridurre l'attrito del coordinamento tra dati e calcolo. Un esempio di utilizzo di un PPP per le risorse di calcolo è il COVID-19 High-Performance Computing Consortium, guidato dall'Office of Science and Technology Policy, DOE, NSF e IBM.⁴⁰ Basandosi sull'esperienza di XSEDE, il consorzio conta 43 membri dai settori pubblico e privato che offrono volontariamente risorse di calcolo gratuite ai ricercatori con proposte di ricerca relative a COVID-19.⁴¹ La natura volontaria del provisioning del calcolo, in questo caso, offre vantaggi sia ai ricercatori, che ottengono un accesso immediato al calcolo, sia membri del consorzio, che contribuiscono all'innovazione e traggono benefici dalle pubbliche relazioni.

Riconosciamo anche che le prove intorno all'efficacia dei PPP è contestata.⁴² In effetti, non esiste un modello di PPP valido per tutti; i PPP differiscono enormemente, a seconda delle responsabilità assegnate tra il settore privato e il settore pubblico, e il successo di un PPP può dipendere dalla sua struttura. Hanno organizzazioni, requisiti di accesso e strategie molto diversi per la gestione della qualità dei dati.⁴⁴ Tali punti decisionali sono cruciali. Ad esempio, alcuni studiosi sottolineano la necessità di un ambiente affidabile per i settori privato e pubblico per gestire le violazioni della privacy e dell'etica nelle industrie sensibili. Queste ulteriori considerazioni in materia di privacy, etica, sicurezza e proprietà intellettuale.

IL NRC COME GOVERNO AGENZIA

L'NRC potrebbe anche essere costruito come nuovo agenzia governativa o ufficio. I principali vantaggi di questo modello sarebbero lo sviluppo di un'istituzione distinta del settore pubblico, dedicata al calcolo e ai dati dell'IA. L'NRC potrebbe essere quello di cloud e dati ciò che gli Stati Uniti Digital Service è per la tecnologia dell'informazione del governo. Tale agenzia dovrebbe essere istituita per statuto o per mandato esecutivo. La legislazione abilitante potrebbe creare personale dedicato e professionale per costruire e sviluppare l'NRC, conferire all'NRC l'autorità di imporre dati interagenzia

condivisione e creare un piano a lungo termine basato sulla strategia nazionale per l'IA.

Ci sono, tuttavia, svantaggi significativi nella creazione di una nuova agenzia o ufficio. In primo luogo, l'NRC non potrebbe rivendicare alcun set di dati governativi e potrebbe successivamente incontrare notevoli ostacoli con la necessità di negoziare con ciascuna agenzia di origine per i dati, per non parlare dei vincoli previsti dal Privacy Act, discussi nel capitolo 5. Ciò detto, consentendo la legislazione potrebbe esentare l'agenzia dai divieti di collegamento dei dati del Privacy Act e trasferire il rischio di contenzioso per fughe di dati alla nuova agenzia. In secondo luogo, una nuova agenzia potrebbe dover affrontare maggiori sfide nel reclutare talenti di alto livello.⁴⁶ Secondo il sondaggio del 2020 sul futuro del servizio governativo, la maggioranza degli intervistati presso le agenzie federali ha convenuto che spesso perdono buoni candidati a causa del tempo necessario per assumere, e meno della metà ha concordato che le loro agenzie hanno abbastanza dipendenti per svolgere un lavoro di qualità.⁴⁷ Inoltre, molti intervistati hanno evidenziato opportunità di crescita professionale inadeguate, incapacità di competere con gli stipendi del settore privato e mancanza di una strategia di reclutamento proattiva come fattori principali che contribuiscono a una forza lavoro non adeguatamente qualificata nelle agenzie federali.⁴⁸ Gli FFRDC, al contrario, possono essere negoziati con organizzazioni esistenti, riducendo potenzialmente i costi di avvio. In terzo luogo, mentre i laboratori nazionali hanno esperienza nel stipulare contratti con enti per costruire strutture di calcolo ad alte prestazioni, non è chiaro come una nuova agenzia/ufficio federale affronterebbe tale compito. Una cosa è che un'entità come US Digital Service aiuti a sviluppare piattaforme IT per le agenzie statunitensi; un altro è costruire contemporaneamente una struttura di supercalcolo molto grande e risolvere problemi di lunga data con l'accesso ai dati. Infine, sarà importante isolare la missione di ricerca dell'NRC dall'influenza politica. Nella misura in cui una nuova agenzia potrebbe fornire meno isolamento dai cambiamenti nelle amministrazioni presidenziali e negli amministratori nominati politicamente, questa è una considerazione importante.

Sebbene questi svantaggi siano considerevoli, un'azione legislativa ambiziosa potrebbe, di fatto, rendere una nuova agenzia governativa un'opzione praticabile.

Capitolo 5: Dati

Conformità alla privacy

La visione che motiva l'NRC è sostenere la ricerca accademica nell'IA aprendo l'accesso sia alle risorse di calcolo che a quelle di dati. I dati federali possono alimentare le scoperte della ricerca di base sull'IA e riorientare gli sforzi dai domini commerciali verso quelli pubblici e sociali. Come affermato nell'invito originale dell'NRC, "i ricercatori potrebbero lavorare con le agenzie per sviluppare e testare nuovi metodi per preservare la riservatezza e la privacy dei dati, mentre i dati del governo forniranno il carburante per le scoperte dall'assistenza sanitaria all'istruzione alla sostenibilità".

Ma è possibile un NRC con dati del settore pubblico, in particolare dati amministrativi delle agenzie governative statunitensi, dati i vincoli legali? Le proposte di ricerca che si estendono ampiamente tra le agenzie per dati personali identificabili o altrimenti sensibili² giustamente susciteranno preoccupazioni sul potenziale rischio per la privacy. Il Privacy Act del 1974, la principale legge federale che disciplina i dati raccolti dalle agenzie governative, sfida fondamentalmente l'idea di un NRC come sportello unico per i dati federali.

Le sue eccezioni di ricerca lasciano qualche incertezza sugli sforzi di ricerca a tempo indeterminato che vanno oltre la ricerca statistica o la valutazione delle politiche a sostegno della missione principale di un'agenzia. Anche se le agenzie ritenessero possibile tale ricerca, i ricercatori sarebbero soggetti a vincoli di accesso e i dati stessi potrebbero potenzialmente richiedere trattamenti tecnici sulla privacy.

Forniamo le seguenti raccomandazioni in merito alla privacy dei dati e all'NRC.

In primo luogo, le agenzie potrebbero essere in grado di condividere dati amministrativi resi anonimi con l'NRC entro i limiti della legge sulla privacy ai fini della ricerca sull'IA, sulla base delle esenzioni della ricerca statistica della legge. In secondo luogo, l'NRC richiederà uno staff di professionisti della privacy che includa ruoli incaricati della conformità legale, supervisione e competenza tecnica. Questi professionisti dovrebbero costruire relazioni con colleghi tra le agenzie per facilitare l'accesso ai dati. In terzo luogo, l'NRC dovrebbe esplorare la progettazione di "stanze sicure per i dati" virtuali che consentano ai ricercatori di accedere ai microdati amministrativi grezzi in un ambiente sicuro, monitorato e basato su cloud. In quarto luogo, raccomandiamo che la task force dell'NRC coinvolga le comunità di ricerca politica e statistica e prenda in considerazione il coordinamento con le proposte per un servizio nazionale di dati sicuri, che ha affrontato ampiamente questi problemi.

Questo capitolo procede come segue. Per prima cosa esaminiamo le leggi esistenti che si applicano alle agenzie governative e le restrizioni che impongono all'accesso e alla condivisione dei dati. Descriviamo quindi le attuali pratiche dell'agenzia per la condivisione dei dati con ricercatori e agenzie ai sensi della legge sulla privacy. Infine, valutiamo le implicazioni degli attuali vincoli legali sulla condivisione dei dati NRC e la più importante proposta affine per promuovere la condivisione dei dati ai sensi dell'Evidence Act.

CHIAVE ASPORTO

- Le agenzie possono condividere informazioni amministrative anonime dati con il NRC sotto la ricerca statistica esenzione della Privacy Atto.
- La volontà e la capacità di un'agenzia di condividere i dati può dipendere dalla misura in cui un progetto di ricerca proposto si allinea con lo scopo principale di un'agenzia.
- L'NRC richiederà uno staff di professionisti della privacy per la conformità legale, la supervisione e le competenze tecniche.
- Identificabile individualmente o i dati sensibili dovranno affrontare ostacoli al rilascio e può garantire la privacy tecnica e/o l'accesso a più livelli le misure.

Notiamo all'inizio che questo capitolo prende in gran parte come un dato di fatto i vincoli statuari esistenti. A livello macro, tuttavia, le sfide nella condivisione dei dati suggeriscono anche che un intervento legislativo ambizioso potrebbe superare molti vincoli esistenti, ad esempio (a) esentando legalmente l'NRC dal divieto del Privacy Act sul collegamento dei dati; (b) concedere all'NRC il potere di assumersi le responsabilità dell'agenzia per le violazioni dei dati; (c) imporre alle agenzie di trasferire tutti i dati che sono stati condivisi in base a un accordo sull'uso dei dati o una richiesta del Freedom of Information Act (FOIA) all'NRC; e (d) richiedere che i piani di modernizzazione IT includano disposizioni per piani di condivisione dei dati con l'NRC.³

LA LEGGE SULLA PRIVACY

Le questioni relative alla privacy dei dati sono al centro dei dibattiti sulla condivisione dei dati e l'NRC non farà eccezione. La maggior parte dei dibattiti sulla privacy dei dati negli Stati Uniti oggi si concentra sul settore dei dati dei consumatori in cui le leggi sulla protezione dei dati negli Stati Uniti sono limitate a inesistenti. Al contrario, molte agenzie governative degli Stati Uniti sono soggette a una solida legge sulla privacy, il Privacy Act del 1974, che è stato approvato in risposta alle preoccupazioni sugli abusi di potere del governo.⁴ Per quasi 50 anni, questa legislazione è stata efficace nel suo obiettivo primario di impedendo al governo degli Stati Uniti di centralizzare e collegare ampiamente i dati sugli individui tra le agenzie. Tuttavia, questo approccio ha avuto un costo, vale a dire che alla maggior parte delle agenzie governative è impedito di condividere e collegare liberamente i dati oltre i confini delle agenzie, il che a sua volta ostacola gli sforzi operativi e di ricerca delle agenzie.⁵ Secondo un esperto di privacy del governo, anche quando autorizzato o incaricato di condividere i dati in modo limitato circostanze, le agenzie federali sono spesso riluttanti a farlo a causa di una miriade di fattori, in particolare la mancanza di adozione di standard di sicurezza dei dati coerenti, nonché difficoltà nel misurare e valutare i rischi per la privacy.⁶ A tal fine, molte agenzie vedono la promessa in l'adozione di misure tecniche sulla privacy, come la privacy differenziale, o la creazione di set di dati sintetici come proxy per dati effettivi, come precursore necessario per consentire la condivisione dei dati sia per scopi di ricerca che per obiettivi tra agenzie.⁷

Nei quasi 50 anni dall'approvazione della legge sulla privacy, ci sono stati sforzi periodici per affrontare l'approccio del governo alla gestione dei dati

Anche se autorizzato o

incaricato di condividere i dati in modo limitato

circostanze, le agenzie federali sono spesso riluttanti a farlo a causa di una miriade di fattori, in particolare la mancanza di adozione di standard di sicurezza dei dati coerenti, nonché difficoltà nel misurare e valutare i rischi per la privacy.

preservando la riservatezza dei dati. Gli esempi includono l'E-Government Act del 2002,⁸ il Confidential Information Protection and Statistical Efficiency Act del 2002,⁹ e, più recentemente, il Foundations for Evidence Based Policymaking Act¹⁰ e la National Data Strategy.¹¹ La maggior parte di questi sforzi è stata finalizzata alla condivisione i dati del governo per l'analisi statistica e la valutazione delle politiche e potrebbe essere necessario ampliare l'ambito delle disposizioni per sostenere la ricerca sull'IA. Riteniamo che questi sforzi siano complementari: l'NRC dovrebbe basarsi su questi sforzi, prestando maggiore attenzione alle risorse di calcolo che consentono lo sviluppo dell'IA e l'analisi avanzata dei dati.

VINCOLI DI LEGGE ON CONDIVISIONE DEI DATI

Una visione dell'NRC è che agisca come un dato magazzino per tutti i dati governativi. Ma quella visione si scontra con i vincoli fondamentali delle leggi progettate per ostacolare la condivisione ampia e illimitata dei dati tra le agenzie governative statunitensi. In mancanza di un regime generale e completo sulla privacy, simile al Regolamento generale sulla protezione dei dati (GDPR) dell'Unione europea, il panorama degli Stati Uniti è frammentato tra un mix di leggi sui consumatori specifiche del settore e alcune leggi specifiche del governo, come il Privacy Act del 1974¹² e limitato orientamenti federali di ampio respiro, come la Fair Information Practice

Principi.¹³ In particolare, il Privacy Act, che si concentra ampiamente sulla raccolta e l'utilizzo dei dati da parte delle agenzie federali e limita la condivisione tra di loro, pone delle sfide alle ambizioni dell'obiettivo dell'NRC di rendere più ampiamente disponibili set di dati governativi altrimenti limitati.

Gli sforzi esistenti, sostenuti da progetti di legge come il E-Government Act del 2002 e Fondamenti di Evidence-Based Policymaking Act, hanno tentato di aumentare l'accesso da parte dei ricercatori alle risorse di dati del governo. Tuttavia, questi approcci erano animati dagli scopi primari della valutazione delle politiche, non dalla ricerca di base sull'IA. Né considerano alcuna ambizione da parte delle stesse agenzie di perseguire la ricerca e lo sviluppo dell'IA.¹⁴

L'applicazione di queste leggi e regolamenti all'NRC, in parte, dipende da tre fattori: (1) la forma istituzionale dell'NRC, come discutiamo nel capitolo 4; (2) se gli utenti di NRC possono invocare l'attuale eccezione per la ricerca statistica del Privacy Act; e (3) se i ricercatori accedono ai dati di più agenzie federali. Qui discutiamo brevemente gli obblighi legali delle agenzie federali. Anche se l'NRC non assume la forma di una nuova agenzia federale autonoma, le agenzie che forniscono dati rimarranno soggette a questi vincoli.

LE LIMITAZIONI ED ESENZIONI DELLA NORMATIVA SULLA PRIVACY

Il Privacy Act è stato emanato in risposta alla crescita ansia per la digitalizzazione, così come lo scandalo Watergate durante la presidenza Nixon. La legge è stata motivata dalle preoccupazioni sulla capacità del governo di raccogliere ampiamente dati sui cittadini e di centralizzarli in database digitali, una pratica emergente all'epoca. È il principale regolamento limitante per la condivisione dei dati del governo e ha conseguenze per l'NRC e, più direttamente, per qualsiasi agenzia governativa che desideri condividere dati con l'NRC.

Collegamento dati

Il Privacy Act si applica ai sistemi di record, che sono definiti come "un gruppo di qualsiasi record sotto il controllo di qualsiasi agenzia da cui le informazioni sono recuperate dal nome dell'individuo o da un numero identificativo, simbolo,

o altro particolare identificativo assegnato all'individuo . in atto tra due agenzie che definiscono lo scopo, l'autorità legale e la giustificazione del programma; tali accordi possono durare 18 mesi, con possibilità di rinnovo.¹⁷ Questi limiti sono stati posti in essere per impedire l'emergere di un sistema centralizzato di registri che potesse rintracciare cittadini statunitensi o residenti permanenti in più domini governativi, nonché per limitare l'utilizzo dei dati per le finalità per le quali sono stati raccolti. In effetti, mentre il collegamento tra set di dati può essere importante per la ricerca sull'IA,¹⁸ potrebbe potenzialmente consentire l'abuso, la sorveglianza o la violazione di tali diritti come la libertà di parola consentendo la persecuzione nelle molte aree in cui un cittadino o residente statunitense interagisce con il

sistema federale.¹⁹

Poiché la restrizione sui collegamenti dati si applica ai collegamenti tra agenzie, la restrizione si applica in due scenari particolari per l'NRC. In primo luogo, se l'NRC è istituito come agenzia federale, la condivisione dei dati dell'agenzia con l'NRC andrebbe contro la limitazione dei collegamenti di dati del Privacy Act. In secondo luogo, l'accesso del personale dell'agenzia federale all'NRC potrebbe sollevare interrogativi sul collegamento dei dati tra agenzie ai sensi del Privacy Act. Tuttavia, la raccomandazione nel capitolo 3 è incentrata sulla concessione alle agenzie di un accesso semplificato alle risorse informatiche sull'NRC e ai dati delle proprie agenzie, non a dati multi-agenzia ospitati sull'NRC. Se l'NRC non è concepito come un'agenzia federale e non concede ai membri dell'agenzia l'accesso ai dati tra agenzie, le restrizioni del Privacy Act sui collegamenti di dati potrebbero non essere applicabili.

Notiamo che questo approccio alla gestione dei dati è sia insolito che non al passo con il settore privato, così come con la ricerca sull'IA in particolare. La capacità sia dell'industria²⁰ che dei ricercatori²¹ di associare più fonti di dati e punti dati a un individuo specifico (anonimizzato) è una pratica comune al di fuori del governo. In effetti, questa limitazione non è quella che molti governi²² o stati USA²³ impongono ai propri sistemi di dati. Tuttavia, le restrizioni del Privacy Act sul collegamento dei dati rimangono incontrastate, anche nei vari sforzi di riforma di cui discutiamo

sotto. Vale la pena notare che l'ampia barra del governo federale contro i collegamenti di dati comporta costi per il benessere. Ad esempio, durante la pandemia di COVID-19, l'impossibilità di condividere e collegare i dati sulla salute pubblica ha creato difficoltà nel tracciare la diffusione e la gravità del virus.²⁴ Mentre progetti come il Coronavirus Research Center²⁵ della Johns Hopkins e il COVID Tracking Project²⁶ hanno tentato di aggregare i dati disponibili, la mancanza di integrazione dei dati ha rallentato importanti risposte operative e di ricerca.²⁷ Altri paesi, ad esempio, hanno integrato i registri di immigrazione e di viaggio per valutare i casi e prevenire le epidemie ospedaliere.²⁸

Riconosciamo il potenziale per il collegamento dei dati per affrontare importanti problemi della società senza raccomandare un collegamento dati all'ingrosso e senza ostacoli. Il collegamento di dati ampio o illimitato solleva preoccupazioni legittime sia sulla privacy individuale che sulla diffusa sorveglianza del governo,²⁹ rese concrete dalle rivelazioni dell'informatore del governo Edward Snowden,³⁰ tra gli altri. Un'iniziativa per collegare i dati della Federal Aviation Administration (FAA) con i dati di altre agenzie per la risposta al COVID-19, ad esempio, incontrerebbe la resistenza del Privacy Act. La Task Force dovrebbe apprezzare queste tensioni e compromessi. In effetti, le agenzie considerano le misure tecniche per la conservazione della privacy una componente necessaria di qualsiasi strategia di dati del governo, poiché metodi come il calcolo multipartito o la crittografia omomorfa (di cui parleremo nel capitolo 8) possono consentire alcune forme di collegamenti di dati tra le agenzie, senza violare la Privacy Atto.

Nessuna divulgazione senza consenso

Un'altra restrizione fondamentale del Privacy Act è la regola "Nessuna divulgazione senza consenso", che vieta la divulgazione di documenti a qualsiasi agenzia o persona senza il previo consenso dell'individuo a cui si riferisce il record.³¹ Perché l'NRC divulgerebbe i dati dell'agenzia federale ai ricercatori (vale a dire, a "persona [e]"), questa regola, a differenza della restrizione sul record linkage, è giuridicamente rilevante e inevitabile.

La legge sulla privacy, tuttavia, contiene una serie di eccezioni a questa regola. Le più pertinenti agli sforzi di condivisione dei dati dell'NRC sono le esenzioni per: (1) "uso di routine"; (2) agenzie specificate; e (3) ricerca statistica. Sotto

la prima deroga, il Privacy Act consente alle agenzie di divulgare dati amministrativi che consentono l'identificazione personale quando tale divulgazione rientra tra gli "usi di routine" dei dati.³² L'"uso di routine" di un set di dati è definito da un'agenzia all'altra ed è semplicemente una specifica depositata presso il registro federale sul piano dell'agenzia di utilizzare e condividere i suoi dati. senza divulgazione.³⁴ Sebbene i tribunali abbiano limitato l'ampiezza con cui un'agenzia può descrivere gli "usi di routine",³⁵ un gran numero di casi d'uso può ancora essere coperto da una breve dichiarazione generale.³⁶ Dovrebbero essere condotte ulteriori ricerche sulle condizioni per quando

la condivisione dei dati per scopi di ricerca costituisce un uso ordinario.

Implicazioni per la condivisione dei dati con i ricercatori

Molto dipenderà dall'interpretazione della "statistica research", applicata alla ricerca sull'IA. Nonostante i vincoli del Privacy Act sulla condivisione dei dati, i ricercatori sono stati convenzionalmente in grado di accedere ai dati direttamente dalle agenzie, sulla base dell'eccezione di ricerca statistica al Privacy Act. Questa eccezione consente la divulgazione di record "a un destinatario che ha fornito all'agenzia un'adeguata garanzia scritta anticipata che il record sarà utilizzato esclusivamente come registrazione di ricerca statistica o registrazione e che il record deve essere trasferito in una forma che non sia identificabile individualmente."³⁷ Ciò richiede l'accesso a un set di dati di ricerca approvato o la negoziazione da parte del ricercatore di un MOU direttamente con l'agenzia, un ruolo che suggeriamo all'NRC di svolgere come intermediario, fungendo da partner negoziale per facilitare l'accesso richieste tra più ricercatori e agenzie (discusse nel Capitolo 3).

Sebbene il Privacy Act non definisca la "ricerca statistica", le leggi e le politiche successive hanno elaborato la definizione. Ad esempio, l'E-Government Act definisce "scopo statistico" per includere lo sviluppo di procedure tecniche per la descrizione, la stima o l'analisi delle caratteristiche dei gruppi, senza identificare gli individui o le organizzazioni che comprendono tali gruppi.³⁸ Nel frattempo, un "scopo non statistico" include l'uso di informazioni di identificazione personale per qualsiasi scopo amministrativo, normativo, di applicazione della legge, giudiziario o di altro tipo che influenzi i diritti, i privilegi o i vantaggi di qualsiasi individuo.³⁹ Cioè, mentre i ricercatori possono

utilizzano dati di identificazione personale per l'ampio scopo di analizzare le caratteristiche del gruppo, non possono utilizzare tali dati per scopi mirati per aiutare le agenzie con, ad esempio, specifiche funzioni giudicanti o esecutive.

Il significato preciso di "scopo statistico", tuttavia, rimane "oscuro e i criteri di valutazione possono essere difficili da individuare". La legge designa esplicitamente il Bureau of Labor Statistics, il Bureau of Economic Analysis e il Census Bureau come agenzie statistiche che hanno maggiori poteri di condivisione dei dati a fini statistici.⁴¹ Queste agenzie utilizzano regolarmente l'IA nello svolgimento delle loro attività statistiche.⁴² Mentre le definizioni di AI sono esse stesse contestate, la ricerca statistica può incapsulare almeno alcune forme di apprendimento automatico e intelligenza artificiale, se tale ricerca analizza le caratteristiche del gruppo⁴³ e non identifica gli individui.

A dire il vero, l'NRC non dovrebbe consentire ai ricercatori o agenzie per condurre una corsa finale attorno alla legge sulla privacy. A tal fine, l'NRC richiederà personale dedicato alla conformità alla privacy e alla supervisione per garantire la conformità. Le questioni chiave riguardanti l'identificabilità individuale, la sensibilità dei dati o il potenziale di collegamento e reidentificazione dovranno essere valutate da tale personale.

Implicazioni per la condivisione dei dati dell'agenzia con l'NRC

Nonostante le strade di cui sopra, le agenzie potrebbero nondimeno essere riluttanti a condividere i dati con l'NRC e i suoi ricercatori. I casi in cui le agenzie federali devono affrontare vincoli alla condivisione dei dati abbondano, anche se è del tutto legale o addirittura obbligatorio a livello federale. Ad esempio, l'Uniform Federal Crime Reporting Act del 1988 richiede alle forze dell'ordine federali di segnalare i dati sui reati all'FBI.⁴⁴ Tuttavia, nessuna agenzia federale sembra aver condiviso i propri dati con l'FBI ai sensi di questa legge.⁴⁵ Allo stesso modo, il Census Bureau è consentito dalla legislazione che lo autorizza a ottenere dati amministrativi da qualsiasi agenzia federale e gli impone di cercare di ottenere dati da altre agenzie quando possibile.⁴⁶ Tuttavia, lo statuto non richiede allo stesso modo alle agenzie di programma di fornire i propri dati al Census Bureau. Cioè, sebbene il Census Bureau sia tenuto a chiedere dati ad altre agenzie, queste agenzie non sono tenute a fornirli, e spesso non lo fanno.⁴⁷

[L] a NRC non dovrebbe consentire a ricercatori o agenzie di condurre e porre fine alla corsa la legge sulla privacy.

Mancato impegno nella condivisione dei dati, anche di fronte a un'autorizzazione legale, può derivare dall'avversione al rischio. Secondo un rapporto del GAO, le agenzie scelgono di non condividere i dati perché tendono ad essere "eccessivamente caute" nell'interpretazione dei requisiti federali sulla privacy.⁴⁸ Poiché le disposizioni legali che autorizzano o impongono la condivisione dei dati sono spesso ambigue,⁴⁹ le agenzie possono sbagliare dalla parte cautela e scegliere di non condividere i propri dati per paura del rischio negativo che l'uso dei dati da parte del destinatario possa violare la privacy o gli standard di sicurezza.⁵⁰ A peggiorare le cose, poiché le agenzie devono dedicare risorse significative per facilitare la condivisione dei dati, possono semplicemente scegliere di non dare priorità alla condivisione dei dati. La mancanza di risorse pone un problema significativo; secondo uno studio del Bipartisan Policy Center sulla condivisione dei dati delle agenzie, circa la metà delle agenzie ha citato finanziamenti inadeguati o l'impossibilità di assumere personale adeguato come ostacolo "più critico" alla condivisione dei dati.⁵¹

L'NRC può superare questi ostacoli chiarendo le disposizioni legali, assicurando che i vantaggi per le agenzie della condivisione dei dati superino i rischi e i costi e sostenendo le risorse. Ad esempio, O'Hara e Medalia descrivono come il Census Bureau sia stato in grado di ottenere buoni alimentari e dati sul welfare dalle agenzie statali. Di fronte a statuti ambigui che autorizzano il Dipartimento dell'agricoltura degli Stati Uniti (USDA) e il Dipartimento della salute e dei servizi umani (HHS) degli Stati Uniti a eseguire collegamenti di dati tra programmi sponsorizzati a livello federale, gli stati sono originariamente giunti a diverse interpretazioni statutarie. Alcuni stati hanno accettato di condividere i propri dati solo dopo che (1) l'Office of General Counsel sia dell'USDA che dell'HHS ha emesso una nota in cui chiariva che la condivisione dei dati con il Census Bureau per scopi statistici era legale e incoraggiata; e (2) gli stati erano convinti che la condivisione dei dati avrebbe consentito la costruzione di prove che avrebbero potuto aiutarli a gestire i loro programmi.⁵²

Una più ampia condivisione dei dati con l'NRC che combina più agenzie o fonti di dati esterne può essere facilitata dall'approvazione di ulteriori leggi che richiedono alle agenzie di condividere i propri dati, soggette a limitazioni specifiche su come tali dati vengono utilizzati dall'NRC. Anche allora, l'effetto di tale requisito non è certo una conclusione scontata. È necessario fare di più sia per chiarire la misura in cui è consentita la condivisione dei dati sia per fornire vantaggi che incentivino le agenzie a condividere i propri dati.

Infine, per garantire la conformità con il Privacy Act, nonché per facilitare il ruolo dell'NRC come intermediario di dati, l'NRC richiederà uno staff di professionisti della privacy che includa posizioni incaricate della conformità legale, della supervisione e dell'esperienza nei metodi tecnici. Questi professionisti dovrebbero costruire relazioni con colleghi tra le agenzie per facilitare l'accesso ai dati.

CASO STUDIO: AMMINISTRATIVO RICERCA DATI REGNO UNITO

Administrative Data Research UK (ADR UK) è un nuovo organismo, istituito nel luglio 2018, per facilitare l'accesso sicuro e ampio a set di dati amministrativi collegati di tutto il governo ai fini della ricerca pubblica.⁵³

ADR UK nasce come punto centrale di coordinamento tra quattro partenariati nazionali - ADR Inghilterra, ADR Irlanda del Nord, ADR Scozia e ADR Galles - nonché l'agenzia nazionale di statistica del Regno Unito, Office for National Statistics (ONS). ADR UK si autodefinisce un "hub strategico a livello di Regno Unito": un punto centrale che promuove l'uso di dati amministrativi per la ricerca, si impegna con i dipartimenti governativi per facilitare l'accesso sicuro ai dati e finanzia la ricerca sul bene pubblico che utilizza dati amministrativi.⁵⁴

Il finanziamento per ADR UK proveniva da un consiglio di ricerca (Economic and Social Research Council, ESRC) ed è stato inizialmente impegnato da luglio 2018 a marzo 2022. È stato fornito un totale di 59 milioni di sterline.⁵⁵

ADR UK svolge tre funzioni principali. In primo luogo, la promozione del valore e della disponibilità dei set di dati amministrativi del governo per la ricerca. ADR UK funge da sostenitore generale dell'uso di set di dati amministrativi provenienti da tutto il governo britannico. Agisce anche come motore specifico della ricerca per il bene pubblico: ha identificato aree specifiche di ricerca che sono di pressante interesse politico (ad esempio, il "mondo del lavoro"⁵⁶) e si sta concentrando sulla creazione dell'accesso a set di dati collegati per i ricercatori che affrontano questi temi prioritari.

La seconda funzione principale funge da punto di coordinamento per incoraggiare la condivisione dei dati governativi, gli standard e il collegamento dei set di dati amministrativi. Soprattutto per i suoi bandi di ricerca, ADR UK è in grado di evidenziare più set di dati, che spesso abbracciano diverse aree di competenza dei dipartimenti governativi che possono essere collegati e utilizzati nella ricerca. In tal modo, ADR UK svolge un ruolo importante nel facilitare la ricerca.

PUNTI CHIAVE

■ *Difesa proattiva per l'uso e il collegamento dei dati:*

data la gamma di agenzie e fonti di dati nel governo, avere un'unica voce coordinata di difesa per l'uso dei dati e il collegamento di set di dati pubblici per il bene pubblico è una funzione importante.

■ *Portare talenti esterni nell'uso dei dati del governo:*

ADR UK ha due schemi:

Borse di studio per la ricerca e sovvenzioni per lo sviluppo di metodi, che si rivolgono a talenti esterni eccezionali con l'intenzione di creare consapevolezza e utilizzare set di dati pubblici nella ricerca all'avanguardia.

■ *Piccoli finanziamenti a fondo perduto per accelerare i metodi di ricerca*

che utilizzano set di dati di grandi dimensioni:

Pubblicando inviti alla ricerca che rispondono in modo ampio temi, ADR UK è in grado di raggruppare una serie di set di dati per rispondere a domande di ricerca ed evitare un singolo focus disciplinare.



SFORZI COMPLEMENTARI PER MIGLIORARE L'APPROCCIO FEDERALE ALLA GESTIONE DEI DATI

Le barriere alla condivisione dei dati create dal Privacy Act hanno da tempo rappresentato una sfida per i ricercatori interessati a utilizzare i dati del governo per valutare o informare la politica.⁶⁸ Le comunità di ricerca politica e statistica, sia all'interno che all'esterno del governo federale, si sono impegnate in ammirevoli sforzi di riforma facilitare la condivisione dei dati per la valutazione delle politiche.⁶⁹

Le basi per il processo decisionale basato sull'evidenza Act (EBPA) del 2018, che ha attuato riforme per migliorare l'accesso ai dati per il processo decisionale basato su prove, è un risultato chiave di questi sforzi fino ad oggi. Tuttavia, molte delle disposizioni della legge che hanno contribuito ad affrontare alcuni degli ostacoli al collegamento e alla condivisione dei dati non sono state approvate dal Congresso. Queste disposizioni, note collettivamente come National Secure Data Service (NSDS), rimangono una priorità assoluta per facilitare ulteriori progressi nella condivisione dei dati a fini di ricerca. Secondo la Data Foundation senza scopo di lucro, uno dei maggiori sostenitori dell'NSDS, il suo passaggio "creerà il ponte tra le capacità di dati decentralizzate del governo con una nuova entità che massimizza congiuntamente le responsabilità di accesso ai dati con le protezioni della riservatezza".⁷⁰

L'NSDS è concepito come un'entità legale indipendente all'interno del governo federale che avrebbe l'autorità legale per acquisire e utilizzare i dati. Tuttavia, questa autorità è attualmente concepita come emanante dall'EBPA, che si concentra sull'utilizzo di dati statistici a fini di costruzione di prove. Potrebbe essere necessaria una fonte più ampia di autorità per scopi di ricerca sull'IA ai sensi dell'NRC, che può essere distinta dagli obblighi dell'agenzia. Una chiara area di sovrapposizione è l'invito della proposta all'NSDS di facilitare le proprie risorse di calcolo, che potrebbero essere armonizzate con le esigenze di calcolo dell'NRC. Analogamente alla discussione del capitolo 4 sulle opzioni organizzative, i sostenitori dell'NSDS identificano un'esigenza fondamentale sia per una fonte di finanziamento affidabile sia per un posizionamento ponderato dell'NSDS all'interno di un'agenzia esistente o come agenzia indipendente o FFRDC. Le aree di terreno comune tra l'NRC

e NSDS, così come l'esperienza e lo slancio alla base della proposta, suggeriscono fortemente che l'NRC si impegni e si coordini con questi sforzi.

Un'altra iniziativa complementare è la Federal Data Strategy (FDS), lanciata nel 2018 dal ramo esecutivo e guidata dall'OMB. FDS è uno sforzo a livello di governo per riformare il modo in cui l'intero governo federale gestisce i suoi dati. Il piano richiama la necessità di un "collegamento sicuro dei dati" attraverso tecniche di privacy tecnica,⁷¹ e incorpora una direttiva dell'ordine esecutivo del 2019 sul mantenimento della leadership americana nell'intelligenza artificiale per "[e] migliorare l'accesso a dati federali di alta qualità e completamente tracciabili, modelli e risorse informatiche per aumentare il valore di tali risorse per la ricerca e lo sviluppo dell'IA, pur mantenendo la sicurezza, la protezione, la privacy e le protezioni della riservatezza, coerenti con le leggi e le politiche applicabili. limitazioni" e per "[p]rovidere formati di schemi tecnici sugli inventari", con particolare attenzione alle fonti di dati aperti (ovvero, dati non sensibili o identificativi individuali).⁷³ I set di dati identificati da questo processo potrebbero essere candidati chiave per popolare l'NRC.

Sebbene sia l'NSDS che l'FDS possano promuovere la condivisione dei dati, questi sforzi sono attualmente concentrati principalmente sulla promozione di scopi di valutazione delle politiche. Fortunatamente, c'è molta sovrapposizione e complementarità tra queste iniziative e l'NRC, che illustra l'ampia importanza di meccanismi più efficaci per condividere i dati federali in modo sicuro e nel rispetto della privacy.

Capitolo 6: Tecnica Privacy e Virtuale Stanze sicure per i dati

Discutiamo ora il ruolo dei metodi di riservatezza tecnica per l'NRC. Negli ultimi decenni, i ricercatori hanno ideato una varietà di metodi computazionali che consentono l'analisi dei dati preservando la privacy. Questi metodi sono molto promettenti per consentire la condivisione di dati governativi a fini di ricerca. Notiamo all'inizio che i metodi tecnici sono solo un meccanismo per rafforzare le protezioni della privacy. Sebbene efficaci, tali metodi potrebbero non essere né sufficienti né universalmente appropriati. L'applicazione di un metodo particolare non esime dalla necessità di verificare se i dati stessi aderiscono a standard di privacy articolati. I metodi discussi qui non sono "sostituti" delle raccomandazioni discusse in precedenza e non giustificano mai di per sé la raccolta di dati altrimenti problematici.

L'uso dei dati dell'NRC introduce due minacce alla privacy individuale. Il primo tipo prevede la divulgazione accidentale da parte delle agenzie (divulgazione dell'agenzia): un'agenzia carica un set di dati nell'NRC che non dispone di una protezione della privacy sufficiente e contiene informazioni identificative su un individuo. Un ricercatore, analizzando questo set di dati da solo o insieme ad altri set di dati NRC, scopre queste informazioni e identifica nuovamente l'individuo.¹ Il secondo tipo prevede la divulgazione accidentale da parte dei ricercatori (divulgazione del ricercatore). In questo caso, un ricercatore rilascia prodotti calcolati su dati NRC limitati (ad es. modelli di machine learning addestrati, pubblicazioni). Tuttavia, i prodotti rilasciati mancano di una protezione della privacy sufficiente e un consumatore esterno del prodotto di ricerca apprende informazioni sensibili su una o più persone nel set di dati originale utilizzato dal ricercatore.²

Raccomandiamo che, a causa dell'infanzia e dell'incertezza che circonda gli usi delle tecnologie che migliorano la privacy, la privacy dovrebbe essere affrontata principalmente attraverso politiche di accesso ai dati. Mentre ci saranno circostanze che suggeriscono, o addirittura impongono, trattamenti tecnici, le policy di accesso, discusse nel Capitolo 3, sono la linea di difesa principale: garantiscono che i set di dati sensibili siano protetti controllando chi può accedere ai dati. Consigliamo un criterio di accesso a più livelli, con set di dati più sensibili inseriti in livelli più limitati. Ad esempio, i dati di accesso altamente ristretti possono corrispondere ai dati sanitari individuali del VA, mentre i dati di accesso minimamente ristretti possono corrispondere alle misurazioni oceaniche del NOAA. Le proposte che richiedono l'accesso a dati altamente riservati dovrebbero affrontare standard di revisione più elevati e i ricercatori potrebbero essere limitati all'accesso a un solo set di dati ad accesso limitato alla volta. Questo approccio rispecchia i regimi attuali in cui i ricercatori seguono una formazione specifica per lavorare con determinati tipi di dati.³

CHIAVE ASPORTO

- Le misure tecniche di riservatezza sono utili, ma non sostituiscono la protezione della riservatezza dei dati tramite criteri di accesso.
- In alcuni casi, l'NRC o le agenzie potrebbero voler fare accesso ai dati subordinato a l'uso della privacy tecnica le misure.
- Agenzie contribuenti e l'NRC dovrebbe collaborare a determinare la privacy tecnica misure basate su set di dati sensibilità, utilità del set di dati e implicazioni sull'equità.
- L'NRC deve avere tecnico personale della privacy per amministrare i trattamenti tecnici sulla privacy, nonché per supportare la ricerca sulla privacy in contraddittorio.
- L'NRC dovrebbe esplorare l'adozione di "dati virtuali". camere sicure" che consentono ricercatori a cui accedere dati amministrativi grezzi o microdati in un ambiente sicuro, monitorato e basato su cloud ambiente.

I trattamenti tecnici sono una linea di difesa diversa:

Riducono significativamente le possibilità di deanonimizzazione di un set di dati. Esistono una serie di metodi tecnici che possono consentire l'analisi garantendo la privacy:

- Tecniche come k-anonymity e γ -diversity tentano di offrire l'anonimizzazione basata sul gruppo riducendo la granularità dei singoli record nei dati tabulari.⁴ Benché efficaci in contesti semplici e facili da implementare, entrambi i metodi sono suscettibili agli attacchi da parte di avversari che possiedono ulteriori informazioni sulle persone nel set di dati.

- Una delle tecniche più diffuse è la privacy differenziale,⁵ che fornisce garanzie dimostrabili sulla privacy, anche quando un avversario possiede informazioni aggiuntive sui record nel set di dati.

Tuttavia, la privacy differenziale richiede l'aggiunta di quantità casuali di "rumore" statistico ai dati e a volte può compromettere l'accuratezza delle analisi dei dati. Sebbene la privacy differenziale sia diventata un punto controverso rispetto al nuovo sistema di prevenzione della divulgazione del Census Bureau,⁶ la tecnica rimane una potente difesa contro i malintenzionati che cercano di sfruttare i dati pubblici ai fini della reidentificazione.

- I ricercatori hanno anche identificato altri metodi promettenti. Recenti lavori lo hanno dimostrato l'apprendimento automatico può essere utilizzato per generare set di dati "sintetici", che rispecchiano i set di dati del mondo reale in modi importanti, ma sono costituiti da esempi interamente sintetici.⁷ Altri lavori si sono concentrati sull'incorporazione di metodi dalla crittografia, tra cui il calcolo multipartitico sicuro⁸ e la crittografia omomorfa.⁹

I metodi che oscurano i dati introducono fondamentali tensioni con il modo in cui i ricercatori di apprendimento automatico sviluppano modelli. Ad esempio, quando si considerano questioni di equità algoritmica, in alcuni casi le protezioni della privacy possono indebolire il potere di valutare se un metodo tecnico come la privacy differenziale si traduca in disparità demografiche, in particolare per piccoli sottogruppi.

campioni su cui un modello di apprendimento automatico funziona male, è fondamentale per il modo in cui i ricercatori migliorano i modelli. Richiede la comprensione degli attributi e delle caratteristiche dei dati al fine di comprendere meglio le carenze di un algoritmo. Pertanto, metodi come la privacy differenziale, che rendono i dati grezzi più opachi, ostacoleranno inevitabilmente il processo di analisi degli errori. I dati sintetici in genere catturano le relazioni tra le variabili solo se tali relazioni sono state intenzionalmente incluse nel modello statistico che ha generato i dati¹¹ e, pertanto, potrebbero non essere adatte a determinati modelli di intelligenza artificiale che scoprono relazioni impreviste tra i dati.

Sebbene la crittografia omomorfa potrebbe non richiedere presupposti simili sulla struttura dei dati, i metodi esistenti sono computazionalmente costosi.

Pur promettendo, comprendendo e applicando questi metodi è un processo scientifico in evoluzione. L'NRC è pronto a contribuire alla loro evoluzione sostenendo direttamente la ricerca sulla loro applicazione.

CRITERI E PROCESSO PER ADOZIONE

L'NRC conterrà una ricca gamma di set di dati, ognuno dei quali presenta implicazioni uniche sulla privacy rispetto a diversi tipi di formati di dati (ad esempio, record tabulari individuali, testo non strutturato, immagini). L'inclusione di un set di dati nell'NRC solleva una questione di scelta: quale trattamento tecnico della privacy dovrebbe essere applicato (ad esempio, k-anonimato vs. privacy differenziale) e come dovrebbe essere applicato? Questa domanda richiede spesso determinazioni tecniche su diverse impostazioni algoritmiche, ma tali scelte tecniche possono anche avere importanti conseguenze sostanziali.¹²

Innanzitutto, raccomandiamo che queste determinazioni siano effettuate rispetto ai seguenti fattori:

- **Sensibilità del set di dati:** set di dati diversi porranno rischi per la privacy che variano per tipo e grandezza. Le cartelle cliniche, ad esempio, sono più sensibili dei modelli meteorologici. Il metodo di privacy scelto dovrebbe riflettere questa sensibilità. Come discutiamo nel capitolo 3, questi metodi di privacy dovrebbero corrispondere ed essere stratificati.

alla classificazione FedRAMP appropriata per il set di dati.

- **Utilità del set di dati:** come discusso in precedenza, l'applicazione di un metodo di privacy può distorcere i dati originali, diminuendo l'accuratezza e l'utilità dell'analisi. Perché metodi diversi influenzano livelli diversi di distorsione, la scelta del metodo dovrebbe essere informata dall'utilità percepita dei dati. I set di dati ad alta utilità, in cui le analisi accurate sono estremamente importanti (ad esempio, strumenti diagnostici medici), possono richiedere metodi che producono meno distorsioni.
- **Equità:** alcune misure sulla privacy possono avere un impatto sproporzionato sui sottogruppi sottorappresentati nei dati.¹³ Nel determinare quale metodo applicare, è necessario valutare la presenza di sottogruppi sensibili e la loro relazione con gli obiettivi del set di dati.

Per ogni dato set di dati, raccomandiamo che le agenzie che forniscono i dati collaborino con il personale NRC per identificare e raccomandare eventuali trattamenti per la privacy. Le agenzie di origine e il personale dell'NRC possederanno competenze di dominio e di ricerca per effettuare valutazioni sull'equilibrio tra privacy, utilità ed equità, ma le agenzie dovrebbero consultare il personale e i ricercatori dell'NRC sui trattamenti più appropriati. Dato il costo della revisione, tali trattamenti della privacy dovrebbero essere considerati molto meno ampiamente per i set di dati a basso rischio.

CAMERE DI SICUREZZA DEI DATI VIRTUALI

Per le singole proposte di ricerca che sarebbe molto ostacolato da misure tecniche sulla privacy, l'NRC dovrebbe esplorare l'uso di "stanze sicure per i dati" virtuali che consentano ai ricercatori di accedere ai microdati amministrativi grezzi in un ambiente sicuro e monitorato. Attualmente, il Census Bureau implementa queste stanze sicure in luoghi fisici e modera l'accesso ai dati grezzi interagenzia attraverso la sua rete di centri dati di ricerca statistica federale (FSRDC). Tuttavia, l'NRC non dovrebbe adottare il modello FSRDC all'ingrosso. In effetti, le barriere all'utilizzo degli FSRDC sono elevate e "solo i ricercatori più persistenti hanno successo".¹⁴ Ad esempio, richiedere l'accesso e ottenere

l'approvazione per l'utilizzo di un FSRDC richiede almeno sei mesi, richiede l'ottenimento dello "Special Sworn Status", che comporta un nulla osta di sicurezza di livello due ed è limitato ai richiedenti che sono cittadini statunitensi o residenti negli Stati Uniti per tre anni.¹⁵ A complicare ulteriormente le cose, le agenzie hanno diversi processi di revisione e approvazione per i progetti di ricerca che desiderano accedere ai dati delle agenzie utilizzando un FSRDC.¹⁶ Infine, anche dopo che l'approvazione è stata concessa, i ricercatori possono accedere ai dati solo di persona andando a luoghi, come lo stesso FSRDC.¹⁷

Per essere chiari, alcune di queste restrizioni sono uniche per l'ufficio censimento. La legge degli Stati Uniti prevede che qualsiasi set di dati del censimento che non protegga completamente la riservatezza possa essere utilizzato solo dal personale del censimento.¹⁸ I ricercatori che tentano di accedere a tali dati devono quindi passare attraverso il rigoroso processo per diventare un appaltatore giurato del censimento. La misura in cui queste restrizioni si applicheranno all'NRC dipenderà dal fatto che l'NRC abbia istituzionalmente sede nel Census Bureau, cosa che alla fine sconsigliamo.¹⁹ Altri problemi, tuttavia, come la mancanza di uniformità interagenzia nella concessione dell'accesso ai set di dati è non un problema esclusivo di Census, ma un problema comune a tutto il governo federale (vedi capitolo 3).

Un altro problema comune, non necessariamente legato agli FSRDC o al Census Bureau, è l'uso di una data room fisica per accedere ai microdati grezzi. L'NRC dovrebbe esplorare un modello di stanza sicura virtuale, in base al quale i ricercatori possono accedere da remoto a tali microdati. Ad esempio, nel settore privato, l'organizzazione di ricerca imparziale e obiettiva, NORC, con sede presso l'Università di Chicago, è un ambiente riservato e protetto in cui i ricercatori autorizzati possono archiviare, accedere e analizzare in modo sicuro i microdati sensibili da remoto.²⁰ Alcune agenzie del governo federale hanno anche implementato le proprie stanze sicure per i dati virtuali. Dati di ricerca virtuale del Centro per i servizi Medicare e Medicaid Center (VRDC), ad esempio, garantisce ai ricercatori l'accesso diretto ai file di dati approvati attraverso una rete privata virtuale . un'enclave

dovrebbe essere fisico o virtuale.²² Come articolato nelle risposte alla RFI dalla Società Americana di Biochimica e Biologia Molecolare e dalla Federazione per la

Societies for Experimental Biology, un'enclave virtuale faciliterebbe notevolmente l'accesso dei ricercatori ai dati e può essere progettata e amministrata in modo da preservare la privacy e la sicurezza.²³

Un National Research Cloud non può funzionare in modo efficace se l'accesso a determinati set di dati è in definitiva legato a una National Research Room.

ARGOMENTO DI STUDIO:

LABORATORIO DI POLITICA DELLA CALIFORNIA

Il California Policy Lab (CPL) è un'università dell'istituto di ricerca della California che fornisce supporto alla ricerca e ai dati per aiutare i governi statali e locali della California a elaborare politiche pubbliche basate su prove.²⁴ CPL offre una varietà di servizi ai governi, inclusi servizi di analisi dei dati e infrastrutture sicure per l'hosting e il collegamento di grandi quantità di dati raccolti da enti governativi.²⁵ Questi servizi aiutano a colmare il divario tra il mondo accademico e il governo, aiutando i responsabili politici ad accedere ai ricercatori e fornendo ai ricercatori un modo sicuro per accedere ai dati amministrativi. CPL mira a costruire partenariati basati sulla fiducia con enti governativi e consentire loro di prendere decisioni politiche supportate empiricamente.

CPL stipula accordi sull'utilizzo dei dati con vari enti governativi in tutta la California, tra cui, ad esempio, il California Department of Public Health e la Los Angeles Homeless Services Authority.²⁶ Questi accordi consentono a CPL di archiviare dati amministrativi in un formato collegabile, promuovendo ampie analisi longitudinali tra vari domini del settore pubblico.

Per aiutare a gestire le esigenze dei vari accordi sull'utilizzo dei dati e semplificare la conformità, CPL applica i requisiti più severi per tutti i dati individuali in tutti i dati che memorizza.²⁷ Ogni serie di dati amministrativi sono quindi soggetti a rigide restrizioni tecniche e controlli approfonditi.²⁸ CPL gestisce i dati in un data hub locale presso l'UCLA. Questo hub di dati utilizza "enclavi virtuali" modellate su camere bianche con air gap tipicamente utilizzate per dati governativi sensibili.²⁹ Le enclavi virtuali sono macchine virtuali che vietare qualsiasi connessione in uscita.

CPL crea un file nuova enclave virtuale per ciascun progetto di ricerca e consente l'accesso solo a specifici ricercatori a set di dati specifici per ogni progetto.³⁰ I ricercatori possono lavorare solo con i dati in l'enclave e can utilizzare solo gli strumenti forniti nell'ambiente.

I processi di accesso ai dati variano, in base ai requisiti degli enti governativi e alla maggior parte dei dati di CPL gli accordi d'uso hanno uno scopo limitato e quindi richiedono l'approvazione dell'ente governativo competente prima di essere utilizzati in un progetto.³¹

In generale, CPL aiuta i ricercatori a capire come ottenere l'accesso a vari tipi di dati amministrativi. Per alcuni set di dati, CPL ha formalizzato le domande sul proprio sito Web.³² CPL preseleziona le proposte di progetto e invia progetti promettenti ai suoi partner governativi per l'approvazione finale. I ricercatori conducono quindi questi progetti approvati sull'infrastruttura sicura di CPL. Per altri set di dati senza processi di accesso formalizzati, CPL indirizza i ricercatori verso gli individui all'interno gli enti governativi.³³ CPL può quindi assumere la gestione dei progetti approvati che mirano a utilizzare i dati archiviati sul proprio hub in base ai propri contratti di utilizzo dei dati. In alternativa, gli enti governativi e i ricercatori stessi possono elaborare nuovi accordi sull'utilizzo dei dati per progetti specifici.

PUNTI CHIAVE

■ *Enclavi virtuali:*

Il CPL utilizza pesantemente secure enclavi virtuali per i ricercatori a accedere, lavorare ed eseguire collegamenti di dati tra set di dati sensibili.

■ *Fare da*

intermediario: CPL facilita e semplifica l'accesso ad amministrativo dati fungendo da intermediario tra i ricercatori e relativo stato agenzie.

IMPLICAZIONI PER IL NRC

PERSONALE DEDICATO

Come discusso in precedenza, sarà fondamentale per l'NRC mantenere uno staff professionale dedicato specializzato in tecnologie per la privacy. In primo luogo, non tutte le agenzie o i dipartimenti che cercano di inserire dati nell'NRC avranno l'esperienza sia per determinare il metodo di privacy che soddisfa le aspettative di utilità dei dati e le richieste di privacy dei dati, sia per applicarlo al set di dati di interesse. Il personale specializzato dell'NRC sarà essenziale per assistere tali agenzie e dipartimenti.

In secondo luogo, anche laddove agenzie e dipartimenti possiedono le competenze necessarie, il personale NRC porterà una prospettiva unica dalle loro collaborazioni in tutto il governo. Laddove il personale di un dipartimento specifico può prevedere solo rischi specifici per l'insieme di dati, il personale dell'NRC sarà in grado di prevedere i casi in cui la presenza di altri dati nell'NRC potrebbe sollevare altre preoccupazioni. Infatti, lavorando con le agenzie governative affiliate e i rappresentanti delle agenzie, il personale dell'NRC può anche aiutare queste agenzie a interiorizzare tali vantaggi, aiutandole a comprendere l'intera gamma di rischi per la privacy in relazione ai loro dati.³⁴ Tale governance collaborativa sarà necessaria per garantire che le valutazioni sulla privacy considerano tutte le implicazioni delle tecnologie di accesso e privacy. Infine, non va trascurato il fatto che mentre la gestione dei dati in generale richiede competenze tecniche, queste varie tecnologie che migliorano la privacy richiedono anche competenze molto specifiche e altamente qualificate. Utilizzando set di dati sintetici come esempio, al personale dell'NRC potrebbe essere chiesto di creare dati sintetici per conto di un'agenzia o di convalidare il lavoro svolto presso un'agenzia per garantire che sia svolto correttamente e bene. Qualunque sia il compito, ci sono effetti a cascata a valle attraverso l'ecosistema della ricerca se non gestiti ed eseguiti con attenzione.

UN FOCUS SULLA VALUTAZIONE E LA RICERCA DELLA PRIVACY MIGLIORARE LE TECNOLOGIE

Sarà necessario valutare continuamente lo stato delle protezioni della privacy sull'NRC, sia da parte dei membri del personale dell'NRC che supportando i ricercatori sulla privacy e la sicurezza presso le istituzioni accademiche. La ricerca sulla privacy e la sicurezza tecnica è per sua natura contraddittoria: i ricercatori adottano la posizione degli avversari per sondare il

Un National Research *Cloud* non può
funzionare in modo efficace se
l'accesso a determinati set di dati è in
definitiva legato a una *sala di ricerca nazionale*.

punti deboli di un sistema/set di dati. Nel contesto dell'NRC, ciò richiederà la simulazione di attacchi mentre i ricercatori cercano di identificare nuovamente gli individui all'interno di specifici set di dati dell'NRC. Questo tipo di ricerca è necessario per far avanzare il campo e l'NRC può essere posizionato in modo speciale per supportare un centro di ricerca dedicato alla ricerca di tecnologie che migliorano la privacy. Ciò consentirebbe alla comunità di ricerca di creare metodi di privacy più solidi per garantire l'anonimato, identificare i difetti e autoregolamentare un ecosistema di dati in evoluzione.

Capitolo 7: Salvaguardie per la ricerca etica

Il ritmo dei progressi nell'IA ha suscitato un ampio dibattito sui principi che dovrebbero governare il suo sviluppo e la sua attuazione. Nonostante la promessa della tecnologia per la crescita economica e i benefici sociali, l'intelligenza artificiale pone anche seri rischi etici e sociali.

Ad esempio, gli studi hanno dimostrato che i sistemi di intelligenza artificiale possono diffondere la disinformazione¹, danneggiare il lavoro e l'occupazione², dimostrare pregiudizi algoritmici in base a età, genere, razza e disabilità³ e perpetuare le disuguaglianze sistemiche.⁴

Questo capitolo considera come l'NRC dovrebbe garantire che le sue risorse siano impiegate in modo responsabile ed etico. Un numero crescente di ricerche sull'equità, la responsabilità e la trasparenza dell'IA ha sollevato domande serie e legittime sui valori impliciti nella ricerca sull'IA e sul suo impatto sulla società. È destinato ad aiutare ad affrontare queste preoccupazioni consentendo maggiori opportunità per la ricerca accademica. Allo stesso tempo, l'ampliamento dell'accesso alle risorse non è sufficiente per garantire che la ricerca accademica sull'IA non esacerbi le disuguaglianze esistenti o perpetui pregiudizi sistematici. Inoltre, l'NRC deve anche essere preparato a gestire e ad agire in caso di denunce di pratiche di ricerca non etiche da parte dei ricercatori.

Mentre vi è abbondanza di quadri etici proposti per l'IA (vedi Appendice C per quelli pubblicati dalle agenzie federali), non esiste un insieme di principi accettati sanciti dalla legge, come la Regola comune per la ricerca sui soggetti umani, che stabilisca chiaramente i confini per la ricerca etica con l'IA.⁶ In mancanza di tale orientamento, una questione fondamentale per l'NRC è come istituzionalizzare la considerazione delle preoccupazioni etiche. Questo capitolo inizia discutendo due potenziali approcci per le proposte di ricerca: revisione ex ante nella fase della proposta per l'accesso alle risorse NRC (ad esempio, calcolo, set di dati) e revisione ex post dopo la conclusione della ricerca. Separatamente, discutiamo le linee guida per l'NRC su questioni relative alle pratiche di ricerca. Uno dei vantaggi di iniziare con l'accesso tramite lo stato di Principal Investigator (PI) (Capitolo 2) è che i ricercatori (a) avranno spesso seguito una formazione di base da parte delle loro istituzioni di origine in materia di conformità alla ricerca, privacy, sicurezza dei dati e pratiche per la ricerca che utilizza l'uomo soggetti; e (b) essere soggetto a standard di ricerca e revisione tra pari (ad esempio, attraverso la revisione IRB ove applicabile). Questi meccanismi non sono sufficienti per coprire molti progetti di ricerca sull'IA, come quando la revisione dei soggetti umani è ritenuta inapplicabile.

Pertanto, adattiamo le nostre raccomandazioni al disegno istituzionale dell'NRC.

In primo luogo, raccomandiamo che l'NRC richieda l'inclusione di una dichiarazione sull'impatto etico per i PI richiedere l'accesso oltre il calcolo di base o per la ricerca utilizzando set di dati limitati.

Ciò fornisce un livello di revisione etica per i progetti con le risorse più elevate che sono già tenuti a sottoporsi a un processo di applicazione personalizzato. In secondo luogo, per altre categorie di ricerca (ad esempio, ricerche condotte con accesso al computer di livello base, dove non è presente alcuna revisione personalizzata

CHIAVE ASPORTO

- Ricercatori richiedenti l'accesso al calcolo al di là dell'allocazione predefinita e/o dei dati limitati (ovvero, quelli sottoposti a un processo di richiesta personalizzato) dovrebbero essere tenuti a fornire una dichiarazione sull'impatto etico come parte della loro richiesta.
- L'NRC dovrebbe stabilire un processo per gestire i reclami sulla ricerca non etica pratiche o uscite.
- Idoneità basata su Investigatore principale lo stato assicurerà qualche recensione sotto il Anche la Regola Comune come attraverso la revisione tra pari, ma consigliamo le università prendono in considerazione modelli più completi per valutare le implicazioni etiche dell'IA ricerca.

contemplato), raccomandiamo che l'NRC stabilisca un processo per la gestione dei reclami che possono derivare da pratiche e risultati di ricerca non etici. In terzo luogo, dati i limiti dei meccanismi precedenti, raccomandiamo l'esplorazione di una serie di misure per affrontare le preoccupazioni etiche nel calcolo dell'IA, come l'approccio adottato dal National Institutes of Health per incentivare l'integrazione della bioetica nella ricerca in corso.

MECCANISMI DI REVISIONE ETICA

DA PRIMA

La revisione ex ante valuta la ricerca ancora da eseguire.⁷ Le agenzie di finanziamento e i consigli di ricerca di tutto il mondo si affidano alle revisioni tra pari ex ante per valutare il merito intellettuale e il potenziale impatto sociale delle proposte di ricerca, sulla base di criteri prestabiliti.⁸ I comitati di revisione istituzionale (IRB) valutano comunemente ricerca accademica che coinvolge soggetti umani prima del suo inizio.⁹ Tuttavia, gran parte della ricerca relativa all'IA potrebbe non rientrare sotto la supervisione dell'IRB, in quanto la ricerca potrebbe non utilizzare soggetti umani o fare affidamento su dati esistenti (non raccolti dai proponenti) su persone che sono pubblicamente disponibili,¹⁰ utilizzato con il permesso della parte che ha raccolto i dati, o è reso anonimo. Potenziali questioni etiche possono, quindi, sfuggire alla revisione dell'IRB.¹¹

Creazione di una revisione etica ex ante trasversale processo, tuttavia, sarebbe impegnativo. In primo luogo, come discutiamo nel capitolo due, raccomandiamo di non esaminare caso per caso tutte le richieste PI per l'accesso al calcolo e ai dati NRC, in quanto tale processo richiederebbe un notevole sovraccarico amministrativo. Nella fase in cui i ricercatori richiedono semplicemente l'accesso al computer, la ricerca può essere così varia e in una fase iniziale che non c'è molto di concreto da rivedere. E nella misura in cui ogni PI richiederebbe una revisione specifica del progetto, tale processo lo sarebbe oneroso.

In secondo luogo, è improbabile che la revisione ex ante faccia i conti con il molte implicazioni etiche delle decisioni di progettazione che hanno luogo dopo l'inizio della ricerca.¹² La progettazione della ricerca può cambiare sostanzialmente rispetto alle proposte iniziali man mano che i progetti procedono. La revisione ex ante potrebbe identificare alcune preoccupazioni, ma è improbabile che tutte.¹³ La natura dell'apprendimento automatico lo è

intrinsecamente incerto - e le previsioni possono essere difficili da spiegare - così come altamente dipendenti dai dati utilizzati per costruire e addestrare i modelli. dei dati utilizzati in un progetto specifico per valutare con qualsiasi attendibilità.¹⁵

In terzo luogo, ci sono preoccupazioni uniche nel discorso accademico sulla valutazione della ricerca da parte del governo. Autorizzare il governo a condurre una revisione etica (separata dalla revisione IRB ai sensi della Common Rule, che è tipicamente delegata alle istituzioni accademiche) con standard vaghi può implicare preoccupazioni accademiche sul discorso, così come sottoporre proposte a valutazioni guidate politicamente che possono passare dall'amministrazione al amministrazione.

Se l'NRC dovesse creare un processo per la revisione ex ante delle proposte di ricerca per questioni etiche, tale consiglio dovrebbe probabilmente essere composto da esperti scientifici ed etici, in modo simile a come l'NSF conduce il proprio processo, anche se forse con l'aggiunta di membri dalle organizzazioni della società civile che si concentrano sulla lotta ai danni causati dall'IA. L'NSF riunisce gruppi di esperti del mondo accademico, industriale, delle aziende private e delle agenzie governative come revisori tra pari, guidati da funzionari di programma e direttori di divisione dell'NSF. di natura interdisciplinare, rendendo impegnative le valutazioni etiche.

EX POST

Le valutazioni ex post forniscono una valutazione dopo che la ricerca si è conclusa.¹⁷ Nel mondo accademico, i ricercatori sottopongono i risultati della ricerca a riviste o conferenze per la revisione tra pari ex post; è in questa fase precedente alla pubblicazione che i revisori o gli editori possono far emergere questioni etiche non identificate dai processi ex ante. Nel settore pubblico, ad esempio, il Privacy and Civil Liberties Oversight Board (PCLOB) conduce revisioni ex post sulle pratiche antiterrorismo da parte dei dipartimenti e delle agenzie del ramo esecutivo per garantire che siano coerenti con le leggi, i regolamenti e le politiche vigenti in materia di privacy e libertà civili. ¹⁸ PCLOB ha anche recentemente iniziato a valutare il

l'uso di nuove tecnologie nella raccolta e nell'analisi di intelligence straniera¹⁹ e per identificare proposte legislative che rafforzino la sua supervisione dell'IA per l'antiterrorismo.²⁰

RACCOMANDAZIONI

Sebbene non raccomandiamo es

Prima della revisione delle proposte di ricerca, raccomandiamo che l'NRC istituisca un processo per gestire i reclami sulle pratiche e sui risultati della ricerca etica. Su questo punto, si raccomanda al CNR di collaborare con l'Ufficio di

Research Integrity (ORI) presso il Department of Health and Human Services per modellare i loro processi e procedure per la gestione dei problemi di cattiva condotta della ricerca.²¹ L'ORI ha una notevole esperienza nella supervisione delle preoccupazioni sulle pratiche di ricerca etiche. Le parti potrebbero presentare una petizione all'NRC per revocare l'accesso quando si dimostra che la ricerca viola manifestamente gli standard o le pratiche di ricerca etica generale applicabili al dominio disciplinare di un ricercatore. Notiamo che l'NRC potrebbe voler adottare uno standard elevato per tale violazione, date le considerazioni sul discorso accademico. Ad esempio, le agenzie federali o le parti esterne che desiderano revocare l'accesso al calcolo o ai dati da parte dei PI dovrebbero presentare un reclamo scritto con prove a sostegno. Le decisioni di revocare l'accesso dovrebbero richiedere il contributo della leadership esecutiva e del consulente legale dell'NRC.

Per i PI che richiedono l'accesso oltre il calcolo di livello base o per set di dati limitati, raccomandiamo di richiedere il completamento delle dichiarazioni sull'impatto etico da presentare con le proposte di ricerca. Una recente proposta per affrontare la mancanza di "processi di revisione etica e sociale professionale ampiamente applicati" nell'informatica ha pilotato tale requisito in un processo di sovvenzione, richiedendo una descrizione dei potenziali impatti sociali ed etici e degli sforzi di mitigazione da parte dei ricercatori.²² Limitiamo questo approccio alle proposte per l'accesso al calcolo oltre l'allocazione predefinita o alle richieste di accesso a set di dati limitati, poiché i problemi di amministrabilità sono più deboli per i ricercatori che stanno già richiedendo l'accesso al calcolo o ai dati oltre i livelli predefiniti. Per tali applicazioni, un processo di revisione di una proposta specifica avverrà già da parte di un gruppo di esperti di revisione esterno (Capitolo 2) e, proprio come la NSF richiede dichiarazioni di "Impatti più ampi"²³ dichiarazioni sulle considerazioni etiche del lavoro

[E] gli approcci etici incorporati

possono . . . identificare[] e affrontare[] i

problemi man mano che la ricerca procede,

a differenza della revisione ex ante, in

cui potrebbe essere troppo presto per

individuare un problema, e della

revisione ex post, che potrebbe essere troppo tardi.

potrebbe facilmente essere incluso. È importante notare che le dichiarazioni sull'impatto etico sarebbero solo una componente delle applicazioni NRC e dovrebbero essere valutate insieme ad altri materiali applicativi. Oltre a richiedere ai ricercatori di riflettere attentamente e documentare i potenziali impatti del proprio lavoro, le dichiarazioni possono anche servire come documentazione utile di potenziali impatti negativi ed essere utili al personale dell'NRC nel determinare se fornire l'accesso a specifici tipi di dati. Tali valutazioni possono anche essere utili per riviste, conferenze o università che affrontano le preoccupazioni ex post sugli impatti etici.

Successivamente, raccomandiamo che l'NRC assuma uno staff professionale dedicato alla supervisione etica, simile a quanto proponiamo per quanto riguarda la privacy dei dati nei capitoli 5 e 6. Oltre al personale dedicato alla gestione delle questioni di conformità legale, l'NRC necessita di personale con formazione specializzata in IA etica (così come competenza in altri sottodomini) per fornire consulenza interna di esperti ai candidati NRC, nonché per aiutare nella valutazione delle dichiarazioni di impatto etico. Allo stesso modo, gli esperti di privacy dei dati possono identificare problemi di privacy etica specificamente correlati ai dati, ad esempio se il consenso è stato correttamente ottenuto e documentato. Per garantire che le decisioni siano basate sul merito, il personale dell'NRC che sovrintende a tali questioni deve operare indipendentemente da altre agenzie federali ed essere isolato da interferenze politiche.

Riconosciamo che questi meccanismi di revisione etica potrebbe non identificare tutti i casi in cui i ricercatori utilizzano l'NRC in un modo per condurre ricerche che sollevano questioni etiche. Pochi meccanismi di revisione potrebbero, in particolare alla luce della notevole ambiguità presente negli standard etici (vedi Appendice C). Tuttavia, questi meccanismi possono aumentare i punti di controllo accademici chiave (revisione IRB e revisione tra pari) in modo amministrabile che non solleva serie preoccupazioni sul discorso accademico.

Infine, raccomandiamo che le parti non NRC esplorino una serie di misure per affrontare le preoccupazioni etiche nel calcolo dell'IA. Questi possono includere un processo di revisione etica o approcci ampiamente utilizzati in bioetica dai National Institutes of Health, in particolare per incentivare l'inserimento di esperti di etica nei progetti di ricerca.²⁴ Tali approcci etici integrati possono avere il particolare vantaggio di identificare e affrontare i problemi man mano che la ricerca procede, a differenza della revisione ex ante, in cui potrebbe essere troppo presto per individuare un problema, e della revisione ex post, che potrebbe essere troppo tardi. Ci aspettiamo che questa sia un'area di indagine attiva man mano che vengono convalidati nuovi approcci. L'NRC, potenzialmente in collaborazione con l'NSF, dovrebbe prendere in considerazione l'idea di offrire finanziamenti per progetti che incorporano esperti di dominio etico in team, al fine di sostenere questa proposta.

Capitolo 8: Gestione

Rischi per la sicurezza informatica

Mentre l'NRC

ha il potenziale
per livellare il

campo di gioco per la

ricerca sull'intelligenza

artificiale, creerà anche

un bersaglio allettante

per una vasta gamma di cattivi attori.

Sebbene l'NRC abbia il potenziale per livellare il campo di gioco per la ricerca sull'IA, creerà anche un bersaglio allettante per una vasta gamma di cattivi attori.

La sicurezza informatica, lo sforzo per proteggere i sistemi da incidenti che potrebbero compromettere le operazioni o causare danni alle risorse e alle parti interessate, sarà un obiettivo fondamentale dell'NRC.

Richiederà un quadro di sicurezza informatica che gestisca i potenziali incidenti durante tutto il loro ciclo di vita, comprendendo: (1) preparazione; (2) rilevamento e analisi; (3) contenimento, sradicamento e recupero; e (4) attività post-incidente, che comprende collettivamente il monitoraggio, il rilevamento, il ripristino e la segnalazione degli incidenti.¹ Pratiche di sicurezza informatica efficaci integrano la valutazione del rischio basata su impatto, immediatezza

e probabilità e contribuiranno a guadagnare la fiducia degli utenti e contrastare la sovversione e interferenze di attori stranieri o altre parti avversarie. Un'attenta progettazione amministrativa dell'NRC con la sicurezza informatica in prima linea stabilirà uno standard elevato man mano che i sistemi di informazione diventeranno più centrali per la nostra infrastruttura nazionale.

In questo capitolo affrontiamo questi problemi di sicurezza informatica. Per prima cosa forniamo una panoramica dei tipi comuni di vulnerabilità e attacchi e valutiamo la loro rilevanza per l'NRC. Successivamente, forniamo una panoramica del panorama normativo del governo federale, per quanto riguarda la sicurezza informatica, con particolare attenzione ai framework FISMA e FedRAMP. Infine, chiudiamo con una discussione sulle misure di sicurezza e progettazione del sistema più adatte a garantire che l'integrità dell'NRC non sia compromessa.

MOTIVAZIONI PER POTENZIALI ATTACCHI

Eventuali attacchi contro l'NRC potrebbero richiedere una serie di approcci, ciascuno dei quali comporterebbe conseguenze sostanziali per l'NRC.² In primo luogo, gli avversari potrebbero lanciare un attacco contro l'NRC con l'intenzione di interrompere le sue operazioni o la sua capacità di aiutare la ricerca. Ad esempio, gli avversari potrebbero attaccare direttamente l'infrastruttura dell'NRC disabilitando o interferendo con i server dell'NRC. Di conseguenza, i ricercatori non sarebbero in grado di accedere ai server NRC o di utilizzarli efficacemente. Lanciando tali attacchi, gli avversari possono soffocare l'NRC, aumentando così i costi per il governo federale.³ In alternativa, gli avversari potrebbero tentare di attaccare specifici progetti di ricerca sull'NRC,

CHIAVE ASPORTO

- Dissuadere gli attori malintenzionati dall'attaccare l'NRC richiederà più che aderire all'attuale FISMA e standard FedRAMP.
- L'NRC dovrebbe centralizzare responsabilità di sicurezza per i set di dati con il personale del programma piuttosto che rimandare alle agenzie originarie.
- Misure tecniche i L'NRC dovrebbe indagare includendo informazioni riservate cloud, apprendimento federato e crittografia misure basate come crittografia omomorfa e calcolo multipartito sicuro.

rallentando così il ritmo di tale ricerca o compromettendo la qualità dei risultati della ricerca. Possono anche avviare attacchi di "avvelenamento dei dati" sui set di dati NRC, compromettendo così la qualità dei risultati della ricerca.

In secondo luogo, i malintenzionati potrebbero anche avviare operazioni informatiche contro l'NRC, con l'intenzione di rubare risorse computazionali. In questo caso, lo scopo non sarebbe quello di interrompere l'NRC, ma di riutilizzare il potere computazionale per scopi illeciti (ad esempio, l'estrazione di criptovalute).⁴ Ad esempio, gli individui potrebbero fingere di essere ricercatori, sostenendo di utilizzare i crediti cloud per scopi di ricerca legittimi effettivamente usandoli per fini alternativi. Gli individui potrebbero anche infiltrarsi nella rete dell'NRC, sottraendo risorse computazionali da altri progetti e riducendo la funzionalità per gli utenti legittimi.

In terzo luogo, gli avversari potrebbero rappresentare una minaccia per l'uscita dell'NRC di un desiderio di rubare o utilizzare i dati e la ricerca prodotti alloggiati all'interno del sistema. L'NRC promette di essere un obiettivo attraente perché ospiterà i dati di una serie di agenzie diverse. Se un avversario volesse rubare dati equivalenti dalle agenzie stesse, dovrebbe entrare in ogni agenzia in modo indipendente.

Tuttavia, la potenziale combinazione di set di dati sull'NRC, compresi i set di dati di proprietà dei ricercatori, può aumentare i potenziali guadagni derivanti dall'accesso a queste informazioni.

Inoltre, gli avversari possono tentare di irrompere nell'NRC per rubare i prodotti generati dai ricercatori dell'NRC. Ciò potrebbe includere modelli di apprendimento automatico addestrati o risultati di ricerche specifiche.

Allo stesso modo, i cattivi attori potrebbero determinare quell'esecuzione le intrusioni nell'NRC è un modo efficace per prendere di mira le agenzie governative affiliate. Poiché un incentivo alla partecipazione per le agenzie è il supporto informatico che l'NRC offrirà, uno dei maggiori rischi informatici è rappresentato dagli attori malintenzionati che tentano di utilizzare l'NRC per hackerare i loro sistemi. Per questo motivo, il rischio per la sicurezza informatica per il governo può essere notevole. D'altra parte, come abbiamo discusso nel capitolo 3, l'NRC offre anche un'opportunità per migliorare e armonizzare la conformità agli standard di sicurezza, man mano che le agenzie si spostano nel cloud.

Può esistere una serie di altre motivazioni. Le operazioni riuscite contro l'NRC, in quanto entità federale, avrebbero un valore simbolico e attirerebbero l'attenzione. Gli attacchi ransomware potrebbero portare a profitti significativi. L'NRC potrebbe anche essere un obiettivo di spionaggio, sia da parte di attori statali nazionali che cercano di acquisire set di dati sensibili (ad esempio, infrastrutture della rete energetica) sia da parte di entità del settore privato che cercano di rubare proprietà intellettuale o monitorare gli ultimi progressi tecnologici.

In caso di successo, qualsiasi attacco potrebbe minare l'NRC. Ad esempio, i ricercatori verrebbero dissuasi dall'utilizzare l'NRC e potrebbero investire i propri sforzi in cloud privati alternativi. Ciò potrebbe accadere perché i ricercatori credono

l'NRC sarebbe inefficace da utilizzare (ad esempio, a causa di frequenti interruzioni del server) o perché ritengono che i loro prodotti di ricerca sarebbero protetti in modo inadeguato.

Le agenzie e i dipartimenti federali potrebbero essere dissuasi dall'affidare all'NRC set di dati sensibili. Le entità federali potrebbero rischiare l'imbarazzo e incontrare ostacoli nell'esecuzione dei loro obiettivi politici se i set di dati venissero trapelati accidentalmente.

Se l'NRC non è sufficientemente sicuro, tali entità possono scegliere di evitare del tutto la condivisione dei dati.

FISMA, FEDRAMP ED ESISTENTI STANDARD FEDERALI

In quanto entità federale, l'NRC sarà soggetto agli standard e ai regolamenti federali. In questa sezione, forniamo una panoramica di alto livello delle due normative più rilevanti: il Federal Information Systems Management Act (FISMA) e il Federal Risk and Authorization Management Program (FedRAMP).⁵ FISMA si applica tradizionalmente ai sistemi non cloud che supportano un'agenzia singola, mentre l'autorizzazione FedRAMP è richiesta per i sistemi cloud.⁶ Concludiamo discutendo le critiche a queste normative.

Fisma

La legge federale sulla gestione dei sistemi informativi (FISMA) è stata approvata per la prima volta nel 2002, con lo scopo di fornire un quadro completo per garantire l'efficacia dei controlli di sicurezza per i sistemi informativi federali.⁷ La legge è stata successivamente modificata nel 2014 e ha

da allora è stata aumentata attraverso altre singole azioni legislative ed esecutive, e la nostra discussione si concentra sull'impatto collettivo delle normative di conformità FISMA.⁸

FISMA si applica a tutte le agenzie federali, appaltatori o altre fonti che forniscono sicurezza delle informazioni per i sistemi informativi che supportano le operazioni e le risorse dell'agenzia.⁹ Investe la responsabilità in diverse entità. In primo luogo, il National Institute of Standards and Technology (NIST) ha il compito di sviluppare standard e linee guida uniformi per l'implementazione dei controlli di sicurezza, valutando la rischiosità dei diversi sistemi informativi e altre metodologie.¹⁰ In secondo luogo, l'Office of Management and Budget (OMB) ha il compito di con la supervisione della conformità dell'agenzia alla FISMA e riferire al Congresso sullo stato della conformità alla FISMA. un programma di sicurezza delle informazioni basato sul rischio in conformità con gli standard NIST e le politiche OMB.¹³ Le agenzie sono inoltre tenute a condurre valutazioni periodiche per garantire la continuità dell'efficienza e dell'efficacia in termini di costi.¹⁴

Vale la pena menzionare qui diversi requisiti NIST. Ai sensi del NIST SP 800-18, le agenzie sono tenute a identificare i sistemi informativi pertinenti che rientrano nell'ambito di competenza di FISMA. Le agenzie devono inoltre classificare ciascuno di questi sistemi in base a un livello di rischio, seguendo le linee guida stabilite in FIPS 199 e NIST 800-60.¹⁵ Schemi NIST 800-53

sia i controlli di sicurezza che le agenzie dovrebbero seguire sia il modo in cui le agenzie dovrebbero condurre le valutazioni del rischio.¹⁶ Le agenzie devono riassumere ulteriormente sia i requisiti di sicurezza che i controlli implementati in "piani di sicurezza", come delineato nel NIST 800-18.¹⁷ Infine, i funzionari dell'organizzazione sono tenuti condurre revisioni annuali della sicurezza in conformità con NIST 800-37.

FEDRAM

Alla fine degli anni 2000, le agenzie federali iniziarono a esprimere problemi di sicurezza come ostacolo all'adozione del cloud computing.¹⁸ In risposta, il Congresso ha approvato il Programma federale di gestione dei rischi e delle autorizzazioni del 2011

(FedRAMP) per fornire un approccio economico e basato sul rischio per l'adozione e l'utilizzo dei servizi cloud da parte del governo federale.¹⁹ L'approvazione FedRAMP è esentata se: (i) il cloud è privato dell'agenzia; (ii) il cloud si trova fisicamente all'interno di una struttura federale; e (iii) l'agenzia non fornisce servizi cloud dal sistema informativo basato su cloud a entità esterne.²⁰ Come FISMA, i requisiti di sicurezza FedRAMP sono regolati dagli standard NIST, tra cui NIST SP 800-53, FIPS 199, NIST 800-37, e altri.²¹ Tuttavia, a differenza di FISMA, i due percorsi di FedRAMP per ricevere un'autorizzazione ad operare implicano che i fornitori che lavorano con più agenzie non devono necessariamente sottoporsi al processo di approvazione completo con ciascuna agenzia. Ciò significa che sia i fornitori di servizi cloud che le agenzie sono in grado di risparmiare molto tempo e denaro.

CRITICHE A FISMA E FEDRAMP

Questi regolamenti non sono senza colpa. In particolare, i critici sottolineano il fatto che, nonostante la loro esistenza, le intrusioni informatiche nelle infrastrutture governative sono comuni e in aumento. mandato dalla FISMA.²³ Gli errori individuati includevano un mancato

proteggere le informazioni di identificazione personale, la documentazione IT inadeguata, la scarsa correzione dei bug, il mancato aggiornamento dei sistemi legacy e l'autorità inadeguata conferita ai responsabili delle informazioni dell'agenzia.²⁴ I rapporti del Government Accountability Office (GAO) sono giunti a conclusioni simili.²⁵ A loro volta, alcuni hanno criticato l'approccio del governo alla sicurezza informatica all'ingrosso, sostenendo che pone troppa enfasi sul semplice rilevamento delle intrusioni. di quegli intrusi per navigare nella rete.²⁷

FedRAMP affronta le proprie critiche. Uno studio recente ha rilevato che ottenere l'autorizzazione può richiedere molto tempo e denaro, impiegando fino a due anni e in alcuni casi costando milioni di dollari.²⁸ Anche se parti di FedRAMP sono progettate per essere riutilizzabili tra le agenzie, le agenzie spesso ritardano il processo imponendo separato,

requisiti aggiuntivi. Sono state rilevate una serie di ragioni per queste carenze, tra cui un comitato di autorizzazione congiunto a corto di personale, una mancanza di fiducia tra le agenzie per quanto riguarda l'autorizzazione a operare (ATO) e un processo di autorizzazione eccessivamente complesso che porta a errori da parte delle agenzie e dei fornitori di servizi cloud.²⁹ Le raccomandazioni proposte per affrontare queste carenze includono maggiori finanziamenti per il Joint Authorization Board di FedRAMP, incentivi per incoraggiare il riutilizzo delle ATO e meccanismi per migliorare l'efficienza del processo di autorizzazione.³⁰

Il 12 maggio 2021 l'amministrazione Biden ha rilasciato un ordine esecutivo (EO) sul miglioramento della sicurezza informatica della nazione,³¹ e l'OMB hanno pubblicato una bozza di strategia federale per un commento pubblico il 7 settembre 2021.³² Firmato all'indomani della violazione del fornitore di software SolarWinds e dell'attacco ransomware a Colonial Pipeline, il EO presenta diverse nuove iniziative. In primo luogo, invita il governo federale ad adottare "l'architettura zero-trust" e migliorare i processi di indagine post-attacco. In secondo luogo, cerca di migliorare la collaborazione tra il settore pubblico e quello privato migliorando i requisiti di divulgazione e istituendo un Comitato di revisione della sicurezza informatica pubblico-privato (sul modello del Consiglio nazionale per la sicurezza dei trasporti). Infine, cerca un approccio alla sicurezza informatica più coerente a livello di governo, chiedendo la creazione di un playbook per standardizzare la risposta informatica tra le agenzie federali, insieme a un sistema di rilevamento e risposta agli attacchi a livello di governo.

Anche se è troppo presto per determinare se l'OE e la strategia proposta saranno efficaci, sembra affrontare le carenze individuate nel panorama esistente. Cerca di migliorare la documentazione e la reattività agli attacchi e suggerisce un cambiamento nel pensiero sulla sicurezza informatica. Non è chiaro, tuttavia, se affronterà i problemi di approvvigionamento sottostanti e la mancanza di fiducia tra le agenzie che i critici ritengono abbiano ostacolato l'efficacia di FedRAMP. Ma data la possibilità di archiviare dati altamente sensibili nell'NRC, l'adozione di un'architettura zero-trust all'inizio è una considerazione cruciale per assicurarne l'integrità.

STANDARD DI SICUREZZA NRC E MISURE DI PROGETTAZIONE DEL SISTEMA

Qui, presentiamo raccomandazioni sulla sicurezza informatica politica per l'NRC informata dal panorama delle normative federali esistenti e considerazioni uniche che un cloud di ricerca nazionale porrà.

PROCESSO PER LA DETERMINAZIONE DEI RISCHI E DELLA SICUREZZA

Nell'attuale panorama normativo, le agenzie sono responsabili della determinazione delle categorie di rischio e dei controlli di sicurezza appropriati per i set di dati che si trovano sui loro server. Tuttavia, ciò solleva una potenziale sfida poiché le agenzie iniziano a condividere i propri dati con l'NRC, rendendo poco chiaro chi manterrà l'autorità per classificare il rischio di questi set di dati e determinare controlli di sicurezza appropriati.

Da un lato, le agenzie stesse potrebbero continuare a mantenere la discrezionalità sulla classificazione di sicurezza e sui controlli per i set di dati che inseriscono nell'NRC. In questo approccio decentralizzato, gran parte delle responsabilità di sicurezza assegnate dalla FISMA rimarrebbero alle agenzie, indipendentemente dal fatto che i dati esistessero sui server NRC. D'altra parte, l'NRC potrebbe assumersi la responsabilità di tutte le decisioni in materia di sicurezza. I set di dati aggiunti all'NRC verrebbero quindi classificati in base alla valutazione del rischio dell'NRC e protetti con i controlli che il personale dell'NRC ritiene appropriati. Questo approccio "centralizza" le responsabilità di sicurezza conferendole all'NRC dopo la negoziazione una tantum per ciascun set di dati.

Sebbene entrambi gli approcci abbiano i loro meriti, raccomandiamo l'approccio centralizzato per diversi motivi. In primo luogo, l'approccio centralizzato garantisce l'uniformità interna. Il paradosso della regolamentazione federale della sicurezza informatica è che sebbene il NIST abbia articolato una serie di standard relativi al rischio e ai controlli, le agenzie interpretano questi standard in modo diverso, portando a discrepanze nell'implementazione e nella classificazione in tutto il governo federale. Seguire le classificazioni di sicurezza di ciascuna agenzia per i dati sull'NRC produrrebbe classificazioni inutilmente complesse e incoerenti per un singolo sistema. Ciò minaccia di diminuire l'usabilità dell'NRC e dell'aggiunto

la complessità potrebbe probabilmente indebolire la sicurezza aumentando la probabilità di errori. Consentire all'NRC di imporre le proprie classificazioni consente l'uniformità all'interno dell'NRC e l'allineamento con i livelli di accesso suggeriti nel capitolo 3 del presente Libro bianco. Questo approccio può anche semplificare la gestione delle pratiche di sicurezza in un potenziale mix di provider di cloud computing.

In secondo luogo, l'NRC rappresenta una preziosa opportunità armonizzare gli standard federali di sicurezza informatica tra le diverse agenzie. Le valutazioni e le implementazioni adottate dall'NRC devono essere generalizzate alla piena diversità dei set di dati federali. Pertanto, le pratiche dell'NRC possono servire da modello per le linee guida del NIST, che qualsiasi agenzia è libera di adottare.

In terzo luogo, l'approccio centralizzato eliminerà gli ostacoli alla condivisione dei dati. I problemi di sicurezza spesso ostacolano la condivisione dei dati delle agenzie. In uno schema in cui le agenzie mantengono il controllo su tutte le determinazioni di sicurezza, le agenzie potrebbero richiedere classificazioni di sicurezza eccessivamente elevate o poco pratiche da implementare. L'approccio centralizzato imporrebbe alle agenzie l'onere di articolare con specificità il motivo per cui le politiche di sicurezza o le linee guida di classificazione dell'NRC sono inadeguate per un particolare set di dati.

Infine, i ricercatori dovrebbero anche avere voce in capitolo nel determinare i controlli di sicurezza appropriati, poiché una risorsa pubblica di questa portata che non può attrarre utenti è destinata a fallire. Poiché i controlli di sicurezza implicano l'usabilità, l'NRC non dovrebbe optare per controlli che sostanzialmente inibiscono o disincentivano i ricercatori dallo sfruttare le proprie risorse. L'NRC deve trovare il giusto equilibrio tra usabilità e sicurezza.

CONSIDERAZIONI TECNICHE

Il governo federale dispone già di una serie di opzioni tecniche e contromisure per gli attacchi informatici. Le minacce e le difese della sicurezza informatica sono, ovviamente, in continua evoluzione, quindi ne discutiamo solo come punto di partenza: una sicurezza informatica solida ea lungo termine passa attraverso una vigilanza continua e una definizione delle priorità che riconosca la natura mutevole del campo.

ARCHIVIO DATI

I meccanismi di archiviazione dei dati dovrebbero garantire un'adeguata protezione dall'accesso esterno. La crittografia può essere utilizzata per proteggere i dati sensibili a riposo, per essere successivamente decrittografati quando necessario. L'isolamento fisico attraverso ambienti con air-gap è un'altra caratteristica di progettazione che può eliminare la possibilità che le interfacce di rete wireless vengano utilizzate per connettere i dati a minacce esterne dannose. Tuttavia, anche il gap d'aria non è una soluzione infallibile: ci sono modi per "saltare" i gap d'aria, ad esempio nascondendosi nelle chiavette USB (che è presumibilmente il modo in cui il malware Stuxnet ha notoriamente compromesso le centrifughe nucleari iraniane).³³ Attacchi più recenti aggirano la necessità per la trasmissione elettronica del tutto sfruttando altri segnali che trapelano dati, come frequenze FM, audio, calore, luce e campi magnetici. Questi tipi di minacce portano a casa la necessità di un approccio completo e in evoluzione alla sicurezza informatica.

PROTOCOLLI DI RETE

I pacchetti di dati inviati attraverso le reti vengono trasmessi secondo una serie di protocolli Internet standardizzati a livello internazionale. Seguendo il modello Open Systems Interconnection (OSI), i livelli concettuali coinvolti nella rete di computer possono essere classificati in sette dimensioni: livello fisico, collegamento dati, rete, trasporto, sessione, presentazione e applicazione.³⁴

SICUREZZA DURANTE L'ESECUZIONE

Quando si considerano le tecnologie di sicurezza in fase di esecuzione, tre caratteristiche di progettazione rilevanti per gli ambienti cloud sono l'uso di cloud riservati, apprendimento federato e misure basate sulla crittografia come la crittografia omomorfa e il calcolo multipartito sicuro. Un numero crescente di fornitori offre opzioni di "cloud riservato" come soluzione tecnica emergente per un calcolo cloud completamente sicuro dal punto di vista informatico che è protetto durante l'esecuzione. ambienti. Per

ad esempio, la virtualizzazione consente a un sistema operativo di eseguire un altro sistema operativo al suo interno come ambiente virtuale con firewall aggiuntivo o altra rete

barriere, simulando efficacemente un altro dispositivo all'interno del computer host.

CALCOLO DISTRIBUITO E APPRENDIMENTO FEDERATO

Un altro paradigma informatico, noto come distribuito informatica o apprendimento federato, considera le situazioni in cui più parti hanno singoli frammenti di dati che sono interessati a sfruttare in forma aggregata, senza condividerli apertamente. L'apprendimento federato affronta questa situazione, ad esempio, dimostrando come i telefoni cellulari degli utenti possono inviare informazioni, possibilmente private in modo differenziato, a server centrali senza esporre i dettagli precisi delle informazioni di un singolo individuo. Un secondo scenario più rilevante per la natura decentralizzata su larga scala dell'NRC è il calcolo distribuito, in cui molte istituzioni condividono collettivamente il calcolo, simile per certi aspetti al calcolo di crowdsourcing. Questi approcci consentono a più parti di sfruttare l'infrastruttura computazionale esistente, pur mantenendo alcune garanzie sulla privacy.

MISURE BASATE SULLA CRITTOGRAFIA

Infine, ci sono due tipi di misure basate sulla crittografia degne di nota.

I ricercatori di crittografia hanno sviluppato metodi per calcolare operazioni matematiche su dati crittografati, noti come crittografia omomorfica. Questa impresa impressionante ha implicazioni preziose perché elimina la necessità di decrittografia, che può potenzialmente esporre i valori intermedi del calcolo e concedere l'accesso alle chiavi di crittografia pubbliche e segrete durante il calcolo. Inizialmente, erano possibili solo schemi di crittografia parzialmente omomorfici che supportavano operazioni aritmetiche limitate come l'addizione e la moltiplicazione. Ma recentemente sono stati sviluppati schemi di crittografia completamente omomorfici che consentono ciò che è noto come calcolo "arbitrario" per casi d'uso promettenti nella medicina predittiva e nell'apprendimento automatico. Detto questo, la standardizzazione è ancora in corso per un'adozione più ampia e la crittografia omomorfica (per progettazione) è malleabile, una proprietà della crittografia che di solito è indesiderabile in quanto consente agli aggressori di modificare i testi cifrati crittografati senza bisogno di conoscere i propri

valore decifrato. Queste e altre limitazioni di qualsiasi approccio tecnico meritano di essere prese in considerazione quando si considera quali tecnologie adottare e per cosa scopo.

A complemento del modello informatico distribuito e decentralizzato discusso in questo White Paper c'è il sottocampo noto come calcolo multipartitico sicuro (noto anche come calcolo che preserva la privacy), che presenta metodi per consentire a più parti di calcolare congiuntamente una funzione su tutti i rispettivi input, mantenendo quelli ingressi privati da altri soggetti. Questi metodi sono maturati nelle loro origini da una curiosità teorica a tecniche con applicazione pratica in studi su documenti fiscali e scolastici, gestione di chiavi crittografiche per il cloud e altro ancora.³⁶ Ciò rende i metodi di calcolo multipartitico sicuri un potenziale candidato per applicazioni relative a calcolo distribuito.

In definitiva, sarà fondamentale per l'NRC conoscere continuamente gli standard di sicurezza più efficaci (comprese altre strategie creative come il red teaming o i bug bounties³⁷ per identificare le vulnerabilità) in questo spazio in rapida evoluzione.

Capitolo 9:

Proprietà intellettuale

Chi dovrebbe detenere i diritti di proprietà intellettuale sui risultati sviluppati utilizzando le risorse NRC?1 Quando la ricerca privata è finanziata, sovvenzionata o influenzata dal governo federale, le leggi e le regole si sono evolute, in modo che sia il ricercatore che il governo abbiano determinati diritti sulla proprietà intellettuale sviluppata nell'ambito della ricerca. Mentre la protezione della proprietà intellettuale è teoricamente progettata per incentivare la ricerca e l'innovazione, alcuni segnali indicano che i ricercatori di intelligenza artificiale in particolare sono già disponibili a condividere i frutti della loro ricerca. In effetti, oltre 2.000 ricercatori hanno firmato una petizione del 2018 per boicottare un nuovo giornale di intelligenza artificiale avviato da Nature, perché prometteva di mettere i suoi articoli dietro un paywall. software e algoritmi disponibili al pubblico.3 Inoltre, come discuteremo in seguito, il progresso di tecniche come il transfer learning dipende dalla capacità dei ricercatori di distribuire liberamente i frutti delle loro ricerche.

Questo capitolo esamina gli accordi di condivisione della proprietà intellettuale esistenti tra ricercatori e il governo, ed esamina se e fino a che punto il governo dovrebbe mantenere i diritti di proprietà intellettuale sui risultati dei ricercatori, come condizione per utilizzare l'NRC.4 Mentre le prove sui diritti di proprietà intellettuale ottimali variano, raccomandiamo che: (1) ricercatori accademici e le università dovrebbero mantenere gli stessi diritti di proprietà intellettuale previsti dal Bayh-Dole Act per i brevetti sviluppati nell'ambito della ricerca finanziata a livello federale; (2) Il governo dovrebbe mantenere i suoi diritti d'autore e i diritti sui dati ai sensi della Guida uniforme, ma stipulare accordi su tali diritti ove applicabile per incentivare l'utilizzo di NRC e l'innovazione dell'IA; e (3)

Il governo dovrebbe prendere in considerazione le condizioni per richiedere ai ricercatori di condividere i loro risultati di ricerca con una licenza ad accesso aperto.

[A] i ricercatori accademici e le università dovrebbero mantenere gli stessi diritti di proprietà intellettuale previsti dal Bayh-Dole Act per i brevetti sviluppati nell'ambito della ricerca finanziata a livello federale.

DIRITTI BREVETTI

Una domanda fondamentale è se gli utenti di NRC debbano conservare i diritti di brevetto sulle invenzioni sostenute dal CNR. Il Bayh-Dole Act regola i diritti di brevetto per le invenzioni sviluppate nell'ambito di accordi di finanziamento federali e la sua applicabilità dipende dal

CHIAVE ASPORTO

- Per armonizzarsi con il processo di sovvenzione federale, l'NRC dovrebbe adottare lo stesso approccio all'assegnazione dei diritti di brevetto, diritti d'autore e diritti sui dati agli utenti NRC come si applica agli accordi di finanziamento federali.

- L'NRC dovrebbe contrarsi sui diritti di proprietà intellettuale del governo, ove applicabile, per incentivare l'utilizzo di NRC e

Innovazione dell'intelligenza artificiale.

- L'NRC dovrebbe prendere in considerazione condizioni per richiedere ricercatori da condividere uscite sotto un open licenza di origine.

natura dell'accesso all'NRC; ad esempio, se i crediti cloud vengono ripartiti utilizzando sovvenzioni federali, come descritto nel Capitolo 2, possono essere considerati accordi di finanziamento federali.⁵ In tali casi, il Bayh-Dole Act consente ai ricercatori di detenere il titolo del brevetto e di concedere in licenza i diritti di brevetto.⁶ Tuttavia, questi diritti di brevetto sono soggetti a determinate restrizioni: ad esempio, l'agenzia finanziatrice ha una licenza gratuita e non esclusiva per utilizzare l'invenzione "per o per conto degli Stati Uniti" e l'agenzia può utilizzare "[m] arch- in rights" per concedere ulteriori licenze.⁷

La questione politica più ampia riguarda quella del governo esercizio dei suoi diritti di brevetto e se e come i brevetti stimolino l'innovazione nell'IA. Alcuni commentatori hanno sostenuto che gli Stati Uniti soffrono di un eccesso di brevetti nel software,⁸ e l'intelligenza artificiale non fa eccezione.⁹ Il numero totale di domande di brevetto AI ricevute ogni anno dall'Ufficio brevetti e marchi degli Stati Uniti è più che raddoppiato, passando da 30.000 nel 2002 a oltre 60.000 nel 2018,¹⁰ e alcuni sostengono che questa proliferazione di ampi brevetti di intelligenza artificiale, in particolare quelli depositati da società commerciali, stia ostacolando l'innovazione futura . esclusività per giustificare l'assunzione dei costi di commercializzazione, come, ad esempio, nel contesto farmaceutico.¹² Per la parte sostanziale dei brevetti universitari, inclusa l'IA, questa logica potrebbe non avere molto peso.¹³

Alcune ricerche mostrano che i brevetti in realtà potrebbero non avere alcun effetto netto sulla quantità o sulla qualità della ricerca sull'IA condotta nel contesto universitario. In uno studio empirico sui docenti delle migliori università di informatica e ingegneria elettrica negli Stati Uniti, la ricerca ha rilevato che la prospettiva di ottenere diritti di brevetto sui frutti delle loro ricerche non motiva i ricercatori a condurre ricerche più o di qualità superiore.¹⁴ L'85% dei professori ha riferito che i diritti di brevetto non erano tra i primi quattro fattori che motivano le loro attività di ricerca, e il 57% dei professori ha riferito di non sapere se o come la loro università condivide i ricavi delle licenze con gli inventori.¹⁵ Lo schema dei brevetti adottato da l'NRC, quindi, potrebbe non avere una forte influenza sull'adozione da parte dei ricercatori.

Il governo dovrebbe. . . prendere in considerazione le condizioni per richiedere ai ricercatori NRC di divulgare o condividere i propri risultati di ricerca con una licenza ad accesso aperto.

Detto questo, in pratica, c'è un vantaggio nel trattare le innovazioni derivanti dall'uso di NRC in modo coerente con Bayh-Dole. In particolare se i crediti cloud vengono assegnati attraverso l'espansione di programmi come NSF CloudBank, sarebbe fonte di confusione avere diritti di brevetto distinti al di fuori della ricerca e della sovvenzione cloud. Inoltre, molti uffici universitari di trasferimento tecnologico sembrano avere forti preferenze per i diritti di brevetto.¹⁶ Nella misura in cui le università considerano il mantenimento dei diritti di brevetto una condizione per l'utilizzo dell'NRC, l'allineamento dei diritti di brevetto dell'NRC con Bayh-Dole può essere preferito, ma l'evidenza alla base di questa raccomandazione non è forte.

COPYRIGHT, DIRITTI SUI DATI E LA GUIDA UNIFORME

La Guida Uniforme (2 CFR § 200) razionalizza e consolida i requisiti del governo per ricevere e utilizzare le sovvenzioni federali per ridurre gli oneri amministrativi.¹⁷ Grants.gov lo descrive come un "quadro governativo per la gestione delle sovvenzioni", una base di regole per le agenzie federali nell'amministrazione dei finanziamenti federali.¹⁸ La Guida uniforme include disposizioni su, ad esempio, principi di costo, requisiti di audit e requisiti per il contenuto dei premi federali.¹⁹

La Guida uniforme è applicabile a "federal awards",²⁰ ma le disposizioni in materia di proprietà intellettuale non richiedono al governo di far valere i propri diritti sui risultati dei ricercatori.²¹ Se e come il governo assegna i propri diritti di proprietà intellettuale ai sensi della Guida uniforme è quindi una questione importante.

Questa sezione copre innanzitutto il copyright del governo e i diritti sui dati relativi alla proprietà intellettuale ai sensi della Guida uniforme e

discute in che modo la condivisione dei diritti d'autore e dei dati potrebbe influire sul panorama dell'innovazione dell'IA. Esaminiamo quindi la misura in cui il governo dovrebbe mantenere i propri diritti sulla ricerca generata utilizzando l'NRC. Sebbene le prove siano contrastanti, alla fine raccomandiamo che il governo mantenga i propri diritti d'autore e i diritti sui dati ai sensi della Guida uniforme, ma applichi contratti su tali diritti ove applicabile, per incentivare l'utilizzo dell'NRC e l'ulteriore innovazione dell'IA.

DIRITTO D'AUTORE

In base alla legge sul copyright degli Stati Uniti, i ricercatori NRC possono ottenere i diritti d'autore su vari aspetti del loro lavoro. Ad esempio, i ricercatori dell'NRC potrebbero desiderare di proteggere i diritti d'autore del software che hanno utilizzato per costruire il modello, poiché il software è considerato un'opera letteraria ai sensi del Copyright Act.²² I ricercatori possono persino ottenere i diritti d'autore su vari aspetti del modello, comprese le scelte dei parametri di addestramento, architetture modello ed etichette di formazione, se possono dimostrare che tali scelte richiedevano creatività.

Ai sensi della Uniform Guidance,²⁵ il beneficiario di fondi federali può tutelare i diritti d'autore su qualsiasi opera sviluppata o acquisita in virtù di un premio federale. Tuttavia, anche se i ricercatori sono autorizzati a mantenere i diritti d'autore, l'ente aggiudicatore federale si riserva un "diritto esente da royalty, non esclusivo e irrevocabile di riprodurre, pubblicare o utilizzare in altro modo l'opera per scopi federali e di autorizzare altri a farlo".²⁶ In particolare, questo diritto è limitato a "scopi federali", nel senso che i terzi che acquisiscono licenze per le opere protette da copyright dei ricercatori non possono utilizzarle per scopi esclusivamente commerciali.²⁷

Non è chiaro fino a che punto i diritti d'autore su NRC i risultati dovrebbero essere pienamente conferiti al ricercatore per stimolare la ricerca di base sull'IA. Una classe di ricerca sull'IA e il risultato dello sviluppo che ha ricevuto una significativa attenzione accademica è stato se le opere creative generate dall'intelligenza artificiale, come la musica del Jukebox di OpenAI,²⁸ possano o debbano ricevere la protezione del copyright.²⁹ Tuttavia, la comunità della tecnologia e del copyright ha appena raggiunto

un consenso sul fatto che l'interesse pubblico nella ricerca sull'IA richieda la concessione del copyright in questi scenari. Da un lato, in un sondaggio condotto tra scienziati di intelligenza artificiale, esperti di politiche tecnologiche e studiosi di copyright, circa il 54% degli intervistati ha convenuto che la protezione del copyright è un importante incentivo per gli autori a rendere il proprio lavoro disponibile in commercio, e il 63% ha convenuto che un aumento della numero di opere prodotte dall'IA disponibili in commercio stimolerebbe un'ulteriore crescita e ricerca sull'IA. intervento di un autore umano.³¹

Nonostante l'importante dibattito sul diritto d'autore sulle opere creative generate dai modelli di intelligenza artificiale, tali opere sono solo un sottoinsieme della possibile protezione del diritto d'autore nel contesto dell'IA. Come discusso in precedenza, i ricercatori potrebbero teoricamente cercare una protezione aggiuntiva del copyright, tra le altre cose, sul loro codice, architettura o modello. Qui, l'innovazione dell'IA può dipendere dalla condivisione di questi elementi protetti da copyright. Ad esempio, il transfer learning utilizza i modelli ML esistenti e li "mette a punto" per un'attività target correlata,³² e sono emersi vari approcci di fine tuning per eseguire il transfer learning su diversi classi di compiti.³³

DIRITTI DEI DATI

Ai sensi della Guida uniforme, quando i "dati" sono "prodotti" nell'ambito di un premio federale, il governo si riserva il diritto di: (1) ottenere, riprodurre, pubblicare o utilizzare in altro modo tali dati; e (2) autorizzare altri a ricevere, riprodurre, pubblicare o altrimenti utilizzare tali dati.³⁴

In particolare, ciò non limita l'uso di tali dati per scopi di governo federale. In altre parole, tali dati possono essere divulgati per qualsiasi uso. La questione in sospeso, quindi, è se questi "dati", che non sono definiti esplicitamente nella Guida uniforme, riguardino i dati generati per scopi di intelligenza artificiale e apprendimento automatico. Di seguito, esaminiamo due classi di dati generati per scopi di intelligenza artificiale (dati sintetici ed etichette di dati) e in che modo la condivisione di questi dati potrebbe influire sull'innovazione dell'IA.

Una classe di dati generati per scopi di intelligenza artificiale è dati sintetici. I ricercatori si sono rivolti a modelli generativi profondi come Variational Autoencoders³⁵ e Generative Adversarial Networks³⁶ per generare dati sintetici per addestrare i loro modelli di machine learning. Come notato dall'Organizzazione mondiale della proprietà intellettuale, i dati sintetici sono una classe di dati completamente nuova che non rientra perfettamente nell'attuale legge sulla proprietà intellettuale. Le disposizioni sul diritto d'autore della Guida Uniforme (descritte sopra), l'ampia classe di dati sintetici, siano essi "creativi" o meno, possono anche implicare la disposizione dei diritti sui dati. Da un lato, i dati di addestramento sono spesso custoditi con cura,³⁸ quindi i requisiti per la condivisione di dati sintetici, che vengono spesso utilizzati per addestrare modelli di intelligenza artificiale, potrebbero non essere un punto di partenza per gli utenti NRC. D'altra parte, molti studiosi hanno scritto sulla promessa dei dati sintetici di consentire effettivamente un'ulteriore condivisione dei dati preservando la privacy e i segreti commerciali dei ricercatori.³⁹ In effetti, la condivisione di set di dati sintetici stimolerebbe ulteriore ricerca e innovazione in campi come l'assistenza sanitaria, la condivisione dei dati è stata limitata.⁴⁰

Un'altra classe di dati generati per l'intelligenza artificiale sono i dati etichettati, vale a dire i dati che sono stati etichettati e classificati per fornire verità di base per i modelli di machine learning supervisionati.⁴¹ Sebbene siano state sviluppate tecniche per ridurre i costi associati all'etichettatura dei dati,⁴² rimane comunque una risorsa e compito dispendioso in termini di tempo. Ad esempio, Cognilytics Research riporta che il 25% del tempo totale impiegato per costruire modelli di apprendimento automatico è dedicato all'etichettatura dei dati.⁴³ I ricercatori che utilizzano l'NRC potrebbero quindi cercare di proteggere il proprio investimento nell'etichettatura dei dati scegliendo di non condividere le proprie etichette con altri, se i dati sottostanti sono proprietari.⁴⁴ Tuttavia, riconoscendo la difficoltà dell'etichettatura dei dati, alcuni ricercatori hanno creato piattaforme online per condividere le etichette dei dati.⁴⁵ Nel caso di ImageTagger, una piattaforma di etichettatura e condivisione dei dati per RoboCup Soccer, gli sviluppatori volevano risolvere il problema che nessuna singola squadra, agendo da sola, potrebbe facilmente costruire i propri set di formazione di alta qualità. marea crescente che solleva tutte le barche, migliorando la qualità non solo dei dati governativi

come set di dati di addestramento, ma anche tutte le ricerche successive che utilizzano tali dati. Inoltre, la condivisione delle etichette dei dati potrebbe essere strumentale per condurre pregiudizi e equità dei risultati della ricerca NRC ove necessario, come discusso nel capitolo 7.

48

MANTENERE I DIRITTI DI PI NELLA GUIDA UNIFORME

Come suggerisce la discussione precedente, la condivisione dei risultati della ricerca sull'IA coperti da diritti d'autore e diritti sui dati potrebbe essere vantaggiosa per l'innovazione dell'IA. Noi quindi raccomandare che l'NRC mantenga almeno gli stessi diritti sui diritti d'autore e sui dati previsti dalla Guida uniforme, ottenendo numerosi vantaggi aggiuntivi. In primo luogo, analogamente alla nostra raccomandazione nel capitolo 3 secondo cui le agenzie federali dovrebbero essere autorizzate a utilizzare le risorse di calcolo dell'NRC, il mantenimento dello stesso schema di allocazione IP della Guida uniforme potrebbe produrre benefici per il benessere migliorando il processo decisionale del governo utilizzando l'IA. Ad esempio, le agenzie federali possono ridurre il costo delle funzioni di governance di base e aumentare l'efficienza e l'efficacia dell'agenzia utilizzando etichette di dati condivise dai ricercatori NRC o perfezionando i modelli generati dai ricercatori NRC. In secondo luogo, ne risulterebbe il mantenimento dello schema di assegnazione della proprietà intellettuale uniforme

in maggiore coerenza nel panorama dei premi federali.

In effetti, come accennato in precedenza nel contesto dei brevetti, potrebbe creare confusione discostarsi dalla Guida uniforme, soprattutto se la sovvenzione del credito cloud è ripartita attraverso programmi come CloudBank ma la sovvenzione per la ricerca è amministrata come premio federale.

Insomma, raccomandiamo almeno al governo conservare i suoi diritti d'autore e i diritti sui dati ai sensi della Guida uniforme. Tuttavia, ribadiamo anche che la Guida uniforme serve semplicemente come quadro utile, non come regola immutabile. Dove l'allocazione PI Guida Uniforme dissuaderebbe i ricercatori dall'utilizzare l'NRC o ostacolerebbe l'innovazione dell'IA in scenari specifici, il governo può e dovrebbe modificare esplicitamente i propri diritti e stipulare contratti separati con i ricercatori su quali diritti il governo conserva, se del caso.

CONSIDERAZIONI PER OPEN REPERIMENTO

Il governo dovrebbe andare oltre i suoi diritti e mandato che i ricercatori condividano la loro ricerca NRC uscite con altri sotto una licenza open-source? Come prima cosa, notiamo che le agenzie possono modificare gli schemi di assegnazione della proprietà intellettuale ai sensi della Guida uniforme, ma non ai sensi del Bayh-Dole Act. Alcune agenzie federali integrano e/o sostituiscono i diritti di proprietà intellettuale stabiliti nella Guida uniforme con restrizioni che sono più specifiche per la proprietà intellettuale sviluppata per quella particolare agenzia o nell'ambito di una concessione specifica.⁴⁹ Ad esempio, il Dipartimento del lavoro richiede che la proprietà intellettuale sviluppato nell'ambito di un premio federale non deve solo rispettare i termini specificati nella Guida uniforme, ma anche essere disponibile per la licenza aperta al pubblico.⁵⁰ I beneficiari NSF sono inoltre tenuti a condividere i propri dati con altri.⁵¹ Tuttavia, il governo non può modificare assegnazione della proprietà dei brevetti ai sensi del Bayh-Dole Act, a meno che la legge stessa non venga modificata o a meno che l'NRC non sia amministrato come federale aggiudicazione, rendendo la legge inapplicabile.

Richiedere ai ricercatori di rendere open source la loro ricerca i risultati possono essere possibili, ma le considerazioni al riguardo sono complesse. Da un lato, un requisito open source potrebbe influire negativamente sulla commercializzazione a valle, data l'ampia gamma di potenziali ricerche sull'IA. Le linee guida o il Bayh-Dole Act potrebbero creare confusione per i ricercatori nel navigare tra i premi federali e comprendere le interazioni delle licenze open source in molteplici situazioni.⁵⁴ Inoltre, richiedere ai ricercatori di condividere i risultati della ricerca comporta una serie di problemi di privacy e sicurezza informatica.⁵⁵ Se i ricercatori consentito di utilizzare l'NRC per condurre ricerche classificate,⁵⁶ ad esempio, mantenere i risultati della ricerca proprietari servirebbe l'interesse nazionale.⁵⁷ In questo caso, tuttavia, l'NRC dovrebbe considerare di limitare qualsiasi requisito di open source alla ricerca che ha meno privacy e sicurezza implicazioni.

D'altra parte, come discusso, la condivisione dei risultati della ricerca con altri ricercatori NRC potrebbe essere vantaggiosa,

e molti studiosi sostengono che i ricercatori di intelligenza artificiale dovrebbero rendere open source il loro software per stimolare l'innovazione. che cercano di mantenere i propri diritti di brevetto, ma le sole divulgazioni dei brevetti software sono spesso limitate ed eccessivamente ampie e non riescono a migliorare il benessere sociale. innovazione. La crescita dei solidi movimenti open source e open science suggerisce anche che un requisito di open source per l'NRC non costituirebbe un ostacolo completo all'utilizzo dell'NRC.⁶¹

Un forte argomento a favore dell'obbligo dell'open-sourcing deriva anche dalla crescente dipendenza del settore privato dai segreti commerciali per la protezione della proprietà intellettuale nell'IA. Al impedendo la divulgazione, fornendo protezione per una durata potenzialmente illimitata e collegandosi immediatamente e ampiamente a qualsiasi output con un valore economico percepibile. cruciale nell'IA e si traduce in un significativo consolidamento del settore dell'IA e in livelli subottimali di innovazione dell'IA. Giocatori. Poiché l'NRC dovrebbe esplicitamente evitare di replicare queste sfide del settore privato, ciò fornisce ulteriore supporto a una raccomandazione secondo cui l'NRC dovrebbe contemplare la richiesta ai ricercatori di condividere i loro risultati di ricerca.

In sintesi, mentre l'intelligenza artificiale solleva una serie di nuovi problemi di proprietà intellettuale (ad esempio,

se l'output dell'intelligenza artificiale è di per sé idoneo per la protezione della proprietà intellettuale), riteniamo che il governo possa evitare molte di queste complicazioni monitorando Bayh-Dole e l'Uniform Guidance. Il governo dovrebbe anche prendere in considerazione le condizioni per richiedere ai ricercatori NRC di divulgare o condividere i loro risultati di ricerca con una licenza ad accesso aperto.

Conclusione

Come abbiamo articolato in questo Libro bianco, l'ambizioso invito a presentare un NRC ha un potenziale di trasformazione per il panorama della ricerca sull'IA.

La sua più grande promessa è garantire un accesso più equo agli ingredienti fondamentali per la ricerca sull'IA: calcolo e dati. Il livellamento di questo campo di gioco potrebbe spostare l'attuale ecosistema da uno che si concentra su problemi commerciali ristretti a uno che promuove la ricerca IA di base e non commerciale per garantire la competitività nazionale a lungo termine, per risolvere alcuni dei problemi più urgenti e per interrogare rigorosamente i modelli IA.

Come abbiamo spiegato in questo Libro bianco, l'NRC solleva una serie di questioni politiche, legali e normative. Come tali risorse di elaborazione possono essere fornite in modo rapido e di facile utilizzo, ma senza precludere i potenziali risparmi sui costi derivanti da una risorsa di proprietà pubblica? In che modo l'NRC può essere progettato per aderire al Privacy Act del 1974, animato da preoccupazioni per un sistema nazionale di registri che sorveglia i suoi cittadini? Come possiamo garantire che NRC mitighi, piuttosto che aumentare, le preoccupazioni sull'uso non etico dell'IA? E come si può evitare che l'NRC diventi il principale obiettivo degli attacchi informatici?

Queste sono domande difficili e speriamo di aver abbozzato il nostro tentativo iniziale di risposte sopra. Siamo fiduciosi, se progettato bene, l'NRC potrebbe aiutare a riallineare lo spazio di innovazione dell'IA da uno che è fissato con il profitto privato a breve termine a uno che è intriso di valori pubblici a lungo termine.

Glossario degli acronimi

ADP	Partnership dati dell'Alberta	FedRAMP	Gestione federale dei rischi e delle autorizzazioni Programma
ADR Regno Unito	Ricerca sui dati amministrativi nel Regno Unito	FFRDC	Ricerca e sviluppo finanziati dal governo federale Centro
ADRF	Strumento di ricerca sui dati amministrativi	FIPS	Standard federali per l'elaborazione delle informazioni
AI	intelligenza artificiale	Fisma	Modernizzazione federale della sicurezza delle informazioni Atto
API	Interfaccia di programmazione applicazioni	FSRDC	Centro dati federale per la ricerca statistica
ARPA	Agenzia per i progetti di ricerca avanzata	GAO	Ufficio per la responsabilità del governo degli Stati Uniti
ARPANET	Agenzia per i progetti di ricerca avanzata Rete	PCP	Piattaforma cloud di Google
QUELLI	autorità ad operare	GDPR	Regolamento generale sulla protezione dei dati
QUELLI	Autorizzazione ad operare	GPS	Sistema di posizionamento globale
AWS	Servizi Web Amazon	GPU	unità di elaborazione grafica
AaaS	Calcolo come servizio	GSA	Amministrazione dei servizi generali degli Stati Uniti
CIPSEA	Legge sulla protezione delle informazioni riservate e sull'efficienza statistica	DUE	Stanford Institute for Human-Centered Intelligenza artificiale
CSM	Centri per i servizi Medicare e Medicaid	ETTORE	Risorsa terascale di elaborazione di fascia alta
CPL	Laboratorio di politica della California	HHS	Dipartimento della salute e dei servizi umani degli Stati Uniti
processore	Unità centrale di elaborazione	HIPAA	Portabilità dell'assicurazione sanitaria e Legge sulla responsabilità
DFARS	Regolamento sull'acquisizione federale della difesa Supplemento	HPC	calcolo ad alte prestazioni
DHS	Dipartimento per la sicurezza interna degli Stati Uniti	http	Protocollo di trasferimento ipertestuale
VENIRE	Dipartimento della Difesa degli Stati Uniti	HTTPS	Protocollo di trasferimento ipertestuale sicuro
DAI	Dipartimento dell'Energia degli Stati Uniti	—	Comunità di intelligence degli Stati Uniti
PUNTO	Dipartimento dei trasporti degli Stati Uniti	IDEA	Istituto per le analisi della difesa
DUE	Accordo sull'utilizzo dei dati	IPTO	Ufficio Tecniche informatiche
EBP	Fondamenti per Evidence Based Legge sulla definizione delle politiche o legge sulle prove	IRB	Comitato istituzionale di revisione
EO	Ordine esecutivo	JV	joint venture
ESRC	Consiglio per la Ricerca Economica e Sociale	LIDAR	Rilevamento e misurazione della luce
EULA	Contratto di licenza per l'utente finale	ML	apprendimento automatico
FAA	Amministrazione federale dell'aviazione	MOU	protocollo d'intesa
LONTANO	Regolamento sull'acquisizione federale		
FDS	Strategia federale dei dati		

NAIR	Ricerca nazionale sull'intelligenza artificiale Legge sulla task force sulle risorse	RIST	Organizzazione di ricerca per l'informazione Scienze e tecnologia
NASA	Nazionale Aeronautica e Spazio Amministrazione	SDSC	Centro dei supercomputer di San Diego
NDA	Legge sull'autorizzazione alla difesa nazionale	SRCC	Centro di calcolo di ricerca di Stanford
NIH	Istituto Nazionale della Salute	SSL	Livello di socket sicuri
NISE	Direzione per Computer e Scienze dell'informazione e ingegneria	STPI	Istituto di politica scientifica e tecnologica
NIST	Istituto nazionale di standard e Tecnologia	TLS	Sicurezza del livello di trasporto
NIST SP	Pubblicazioni speciali del NIST	UC Berkeley	Università della California, Berkeley
NOAA	Nazionale oceanica e atmosferica Amministrazione	UC San Diego	Università della California, San Diego
CHI	Centro nazionale di ricerca sull'opinione pubblica	UCLA	Università della California, Los Angeles
NRC	Nube di ricerca nazionale	USDA	Dipartimento dell'Agricoltura degli Stati Uniti
NSCAI	Commissione per la sicurezza nazionale sull'artificiale Intelligenza	VRDC	Centro dati di ricerca virtuale di CMS
NSDS	Servizio nazionale di dati sicuri	OMPI	Organizzazione mondiale della proprietà intellettuale
NSF	Fondazione Nazionale della Scienza		
ODNI	Ufficio del direttore dell'intelligence nazionale		
OCSE	Organizzazione per la cooperazione economica e Sviluppo		
OMB	Ufficio di gestione e bilancio degli Stati Uniti		
NOI	Ufficio per le statistiche nazionali		
ORNL	Laboratorio nazionale di Oak Ridge		
OLCF	Oak Ridge Leadership Computing Facility		
ANCHE	Interconnessione di sistemi aperti		
OT	Altra transazione		
PCLOB	Consiglio di vigilanza sulla privacy e le libertà civili		
PHI	informazioni sanitarie protette		
PS	Centro per la popolazione di Stanford Scienze della salute		
PI	investigatore principale		
IPI	informazioni di identificazione personale		
PPP	partenariato pubblico-privato		
Ricerca e sviluppo	ricerca e sviluppo		
RFI	Richiesta di informazioni		
RFP	Richiesta di proposta		

Appendice

A. CONFRONTI DEI COSTI DELL'INFRASTRUTTURA INFORMATICA

Questa appendice fornisce un esempio di confronto della stima dei costi tra un servizio cloud commerciale, AWS, e un sistema HPC governativo dedicato, Summit. In sintesi, le nostre stime mostrano che le istanze AWS P3 con hardware paragonabile a Summit sarebbero 7,5 volte più costose dei costi stimati in condizioni di utilizzo costante e 2,8 volte rispetto ai costi stimati di Summit in condizioni di domanda fluttuante.

La tabella 3 elenca i tre modelli di infrastruttura utilizzati in questo confronto. Summit è stato utilizzato come sistema HPC governativo di riferimento perché è uno dei sistemi più recenti del DOE e dispone di hardware adatto per la ricerca sull'IA.¹ L'altro modello di infrastruttura utilizzato è AWS EC2

P3.2 Entrambi sono comunemente usati nella ricerca sull'IA e nelle applicazioni HPC generali. Anche altre piattaforme cloud commerciali, come GCP o Azure, potrebbero fornire in modo fattibile l'infrastruttura per l'INRC. Qui è stato utilizzato AWS EC2 P3 perché AWS dispone di un solido calcolatore dei costi che lo consente carichi di lavoro variabili.

Il numero di istanze AWS è stato impostato in modo tale che quelli i modelli avrebbero esattamente lo stesso numero di GPU di Vertice. Le GPU erano la variabile fissa perché le GPU lo sono l'hardware più importante per le applicazioni di ricerca AI, in particolare il deep learning. Entrambe le istanze Summit e AWS P3 utilizzano GPU NVIDIA V100.

Conduciamo il nostro confronto dei costi per i due modelli di infrastruttura nell'arco di cinque anni, poiché i documenti RFP iniziali di Summit includono un contratto di manutenzione di cinque anni. AWS, tuttavia, fornisce solo piani tariffari annuali o triennali, quindi abbiamo estrapolato il costo quinquennale in base al suo piano triennale.

Per il preventivo dei costi di Summit, abbiamo basato il nostro calcolo sui dettagli del budget nella richiesta originale di proposta (RFP) del Dipartimento dell'Energia (DOE) del gennaio 2014.³ La RFP include un importo di 155 milioni di dollari

budget massimo per la costruzione del Summit, un massimo previsto di \$ 15 milioni per i costi di ingegneria non ricorrenti⁴ e circa \$ 15 milioni per la manutenzione quinquennale⁵ più interessi basati sui titoli del Tesoro USA a scadenza costante quinquennale come specificato nel prezzo programma.⁶ In base ai calcoli, abbiamo stimato che Summit costa circa 192 milioni di dollari in totale, il che è coerente con la rendicontazione pubblica del costo di Summit.⁷

Per la stima dei costi di AWS, abbiamo utilizzato i prezzi AWS calcolatrice, scegliendo Stati Uniti orientali (Virginia settentrionale) come data center e tariffe pubblicamente disponibili nell'ambito del piano tariffario più economico possibile (EC2 Instance Savings Plans). Per approssimare uno sconto negoziato, abbiamo applicato uno sconto del 10% basato sulla tariffa negoziata di una grande università.

Poiché i costi della piattaforma cloud commerciale si adattano al numero di istanze effettivamente in uso, sono stati calcolati due costi per ciascun modello AWS che rappresenta gli estremi di utilizzo: (1) con l'infrastruttura in uso costante; (2) con l'infrastruttura soggetta a fluttuazioni drammatiche di utilizzo ogni giorno. Per il calcolo del picco di traffico giornaliero, impostiamo il modello in modo che venga eseguito cinque giorni alla settimana con 8,4 ore al giorno al massimo delle prestazioni. Il numero massimo di istanze utilizzate è lo stesso che si userebbe per constant utilizzare mentre il numero minimo è zero. Questo carico di lavoro l'impostazione si basa sul presupposto che le GPU utilizzate per l'addestramento dei modelli di intelligenza artificiale rimangano inattive per il 30% del tempo.⁸ Queste stime dovrebbero fornire limiti superiori e inferiori rigidi sui costi per l'utilizzo di ciascun tipo di istanza.

La Figura 1 riporta i costi sull'asse y per un periodo di cinque anni sull'asse x. La linea turchese indica il costo di un sistema simile a Summit e le linee viola e blu indicano il costo delle stesse istanze AWS sotto

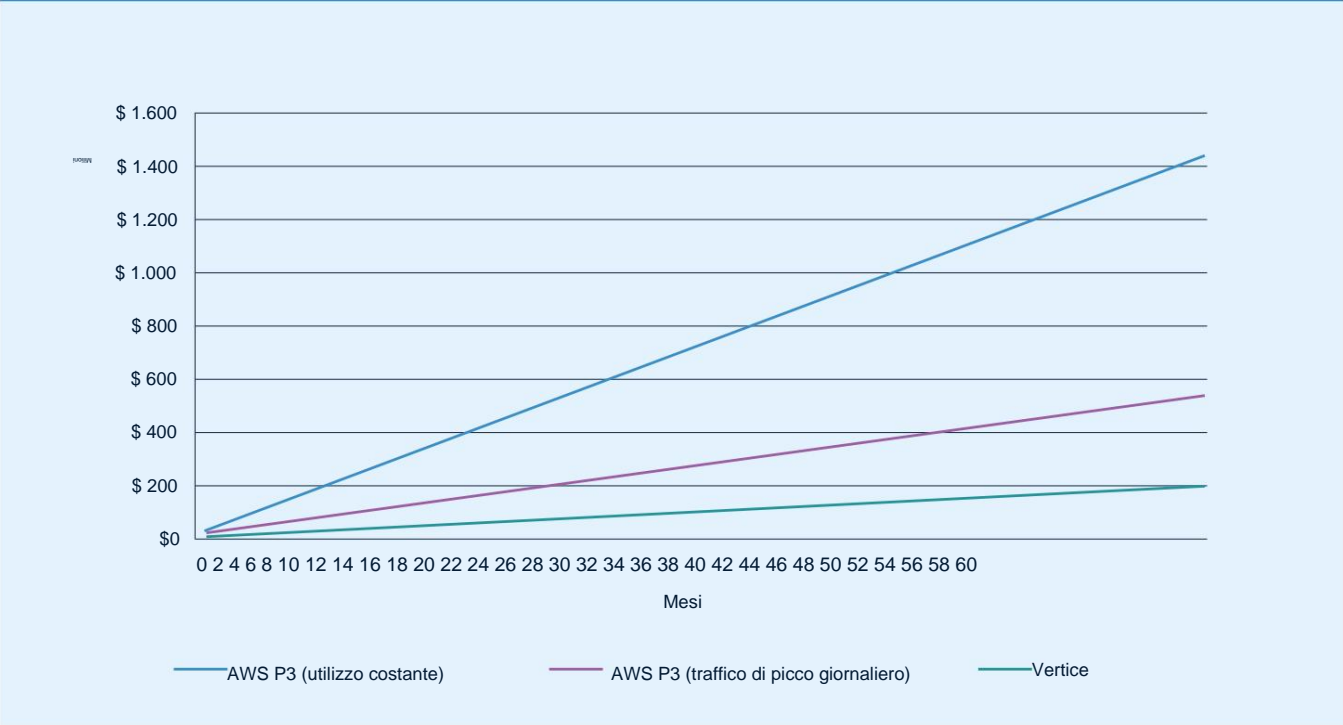
utilizzo variabile e costante. Nel complesso, questa semplice analisi corrobora l'analisi condotta da Compute Canada, che ha rilevato che il cloud commerciale "va da 4 a 10 volte superiore al costo di proprietà e gestione dei nostri cluster".⁹ In cinque anni e con un utilizzo costante, le istanze AWS P3 con un hardware paragonabile a Summit sarebbe 7,5 volte più costoso dei costi stimati. In condizioni di domanda fluttuante, le istanze AWS P3 costerebbero 2,8 volte i costi stimati di Summit.

Notiamo che questa semplice analisi omette molti fattori potenziali (si veda la discussione nel Capitolo 2), ma fornisce un punto di partenza puntano a comprendere le notevoli implicazioni in termini di costi per la decisione make-or-buy.

TABELLA 3: CONFRONTO SUMMIT E AWS

	GPU	RAM	Larghezza di banda della rete
Vertice IBM AC922	27.648 (NVIDIA Volta V100)	2,8 milioni di dollari	200 GB/sec
AWS p3dn.24xlarge (3456 nodi)	27.648 (NVIDIA Volta V100)	2,6 PB	100 GB/sec

FIGURA 1: COSTO STIMATO DELLE ISTANZE AWS RISPETTO AL SUMMIT IN 3 ANNI



B. SET DI DATI PRIVATI FACILITATI CONDIVISIONE

Sfide IP uniche sorgono se i ricercatori sono autorizzati condividere i propri set di dati privati con l'NRC. In effetti, i ricercatori che "caricano" dati proprietari potrebbero essere preoccupati per come altri utenti NRC utilizzano tali dati.¹⁰

Attraverso le interviste condotte per questo libro bianco, le parti interessate aziendali che rappresentano l'industria dell'intrattenimento, così come altre industrie creative, hanno ulteriormente espresso il timore che i ricercatori possano caricare e condividere dati su cui non detengono diritti. Tuttavia, se l'NRC decide di facilitare la condivisione privata dei dati, dovrebbe prendere in considerazione l'adozione di due requisiti per affrontare queste preoccupazioni: (1)

L'NRC dovrebbe richiedere a tutti gli utenti di affermare di avere i diritti di proprietà intellettuale originali sui dati o che i dati sono già di dominio pubblico; e (2) L'NRC dovrebbe disporre di uno schema che consenta ai suoi utenti di concedere in licenza i propri dati.

(a) Gli utenti NRC devono possedere i diritti di proprietà intellettuale sui dati che stanno caricando

I ricercatori che caricano i dati devono concordare che essi possiede i diritti di proprietà intellettuale sui dati prima del caricamento o che i dati sono già di dominio pubblico. Questo dovrebbe essere il caso se i ricercatori condividono i dati ampiamente con altri ricercatori o semplicemente utilizzare i propri dati per uso privato.

Ovviamente, nonostante l'obbligo per gli uploader garantire la proprietà legittima o lo stato di dominio pubblico dell'IP caricato, gli uploader possono comunque caricare dati di cui non possiedono i diritti IP. Ciò può accadere perché gli ingegneri informatici e i ricercatori non sono informati sulla legge sulla proprietà intellettuale, prevedono che il fair use scuserà il loro comportamento o semplicemente sperano di non essere scoperti.¹¹ Le parti interessate del settore erano anche preoccupate che i ricercatori di intelligenza artificiale

estrarrebbe "fatti" da un'opera protetta da copyright (ad esempio, determinate melodie nel ritornello di una canzone) o applicherebbe determinati algoritmi all'opera e rivendicherebbe "erroneamente" un copyright sull'opera trasformata. In ogni caso, questo assemblaggio di dati di input protetti rappresenta la "più chiara responsabilità del copyright nel processo di apprendimento automatico" perché l'assemblaggio di dati protetti viola il diritto alla riproduzione e qualsiasi pre-elaborazione dei dati potrebbe

violare il diritto alle opere derivate.¹²

Nelle interviste, le parti interessate aziendali hanno espresso il desiderio di ostacolare il caricamento di opere protette da copyright chiedendo allo stesso NRC di valutare se i dati caricati sono già protetti da copyright. I dati di diligenza possono essere completati manualmente o utilizzando sistemi automatizzati come Content ID, utilizzato anche da aziende come YouTube.¹³ La prima opzione sarebbe molto laboriosa,¹⁴ mentre la seconda potrebbe essere proibitivamente costosa,¹⁵ quindi il valore di affrontare queste preoccupazioni deve essere soppesato rispetto a questi onerosi costi.

Infine, non è chiaro fino a che punto il caricamento e la condivisione di dati protetti da copyright per l'apprendimento automatico equivale a fair use.¹⁶ Il caso più analogo è *Author's Guild v. Google Books*.¹⁷ In quel caso, Google ha scansionato oltre 20 milioni di libri, molti dei quali erano protetti da copyright, e ha assemblato un corpus di - testi leggibili per potenziare il suo servizio Google Libri.¹⁸ Il 2° Circuito ha ritenuto che le riproduzioni non autorizzate di opere protette da copyright da parte di Google Libri costituissero un fair use trasformativo, soprattutto perché Google Libri forniva informazioni sui libri attraverso piccoli frammenti, senza minacciare il nucleo tutelabile dei titolari dei diritti espressione nei libri.¹⁹ Mentre alcuni hanno affermato che la holding della Gilda degli autori protegge categoricamente l'utilizzo di materiale protetto da copyright nei set di dati per scopi di apprendimento automatico,²⁰ molti giuristi non sono così sicuri di una holding così ampia, soprattutto perché il fair use è così intensivo di fatti.²¹ In effetti, mentre Google Libri utilizzava opere protette da copyright per uno scopo non espressivo, Sobel osserva che i modelli di apprendimento automatico possono essere sempre più in grado di ricavare valore dagli aspetti espressivi di un'opera.²² Pertanto, fino a quando i tribunali e i legislatori non forniranno maggiore chiarezza sull'applicabilità del fair use nel contesto dell'apprendimento automatico, l'NRC dovrebbe comunque richiedere a chi carica i dati di attestare di possedere i diritti sui dati.

(b) Gli utenti devono essere in grado di concedere in licenza i propri dati ad altri utenti.

Se l'NRC abilitasse la condivisione dei dati privati, gli utenti dovrebbero chiarire quali diritti hanno gli altri utenti dell'NRC sui dati condivisi degli autori del caricamento. L'NRC avrebbe due opzioni di base per la creazione di schemi di licenza IP: (1) L'NRC

potrebbe consentire ai ricercatori di utilizzare qualsiasi licenza IP desiderino quando condividono i propri dati privati; o (2) L'NRC potrebbe imporre una licenza uniforme su tutta la linea per tutti i dati che viene caricato.

(1) Scelta della licenza da parte del ricercatore

Consentire ai ricercatori di creare le proprie licenze IP accordi quando la condivisione di dati privati con altri ricercatori sarebbe la soluzione più agevole da

la prospettiva dell'autore del caricamento; consentirebbe loro di condividere esattamente ciò che vogliono e limitare l'uso solo a determinati contesti. Questa scelta di licenza sembra essere importante per chi condivide i dati.²³ In effetti, molti data scientist e ingegneri hanno scritto guide che consigliano i membri della comunità open source su come dovrebbero scegliere licenze specifiche per il loro lavoro.²⁴ GitHub, una piattaforma di condivisione di codice open source, consente i suoi utenti possono scegliere tra dozzine di licenze,²⁵ e FigShare, una piattaforma di condivisione dei dati per i ricercatori, supporta allo stesso modo una serie di diverse licenze Creative Commons.²⁶ Alcuni set di dati

hanno anche i propri accordi di licenza IP personalizzati. Il set di dati accademici di Twitter, ad esempio, è concesso in licenza in base all'accordo con gli sviluppatori di Twitter e alle politiche di utilizzo non commerciale, non a una licenza open source esistente.²⁷

Tuttavia, ci sono degli svantaggi a tale flessibilità. Solo perché potrebbero essere consentite licenze diverse non significa che queste licenze saranno pienamente comprese da tutti gli utenti. L'adozione di più licenze può comportare un aumento delle violazioni accidentali. In effetti, uno studio condotto dall'Institute of Electrical and Electronics Engineers ha rilevato che "sebbene gli sviluppatori [di software] comprendessero chiaramente i casi che coinvolgono una licenza, hanno avuto difficoltà quando erano coinvolte più licenze"²⁸ e, in particolare, sono stati trovati a "mancare di conoscenza e comprensione per separare le interazioni di licenza in più situazioni." ²⁹

In particolare, i ricercatori che non hanno familiarità con le concessioni fornite da diverse licenze dati, in contesti in cui è implementata più di una licenza, possono portare alla violazione di determinate licenze. Ad esempio, quando i ricercatori sono stati intervistati in merito alla loro comprensione degli accordi di trasferimento del copyright nell'IP

processo di commercializzazione, hanno dimostrato solo un punteggio medio del 33% in un sondaggio di verifica delle conoscenze.³⁰

(2) Contratto di licenza uniforme

La seconda opzione a disposizione dell'NRC sarebbe quella di imporre che tutti i dati privati siano concessi in licenza con un'unica licenza uniforme. Per la stessa amministrazione NRC, questa potrebbe essere l'opzione più semplice, poiché gli utenti potrebbero essere avvisati al momento dell'accesso sull'uso appropriato dei dati. Lo svantaggio di questa strategia è che potrebbe scoraggiare aspiranti ricercatori che condividerebbero i dati sotto una licenza più ristretta.³¹ Data la volontà di consentire ai ricercatori per innovare liberamente, potrebbero sorgere dubbi sull'adozione di un accordo di licenza restrittivo. Tuttavia, diverse opzioni di accordi di licenza sarebbero ancora disponibili per l'adozione e questo percorso richiederebbe la scelta di un accordo uniforme tra queste opzioni, con la possibilità di consentire l'opt-out di questa licenza predefinita.

Se l'NRC dovesse implementare una licenza uniforme, potrebbe considerare gli accordi di licenza sfruttati dai cloud di ricerca istituzionali, come l'Harvard Dataverse, come un'analogia nella determinazione delle migliori pratiche per i propri accordi di licenza. Il modello adottato dal Dataverse è un utilizzo predefinito del CC0 Public Domain

Dedica "a causa del riconoscimento del suo nome nella comunità scientifica" e del suo "uso da parte di archivi e riviste scientifiche che richiedono il deposito di dati aperti". dati che disciplina per essere utilizzati in qualsiasi contesto, anche commerciale, e consentirebbe anche la riproduzione e la creazione di derivati dai dati.

In alternativa, l'NRC potrebbe avere una licenza aperta predefinita, consentendo anche ai ricercatori di scegliere tra una manciata di licenze più restrittive, se lo desiderano. Ad esempio, l'Harvard Dataverse consente in particolare agli uploader di rinunciare al CC0 se necessario e specificare termini di utilizzo personalizzati. L'Australian Research Data Commons e Anche la piattaforma di condivisione dei dati FigShare³³ utilizza una licenza CC0 predefinita, ma consente comunque ai ricercatori di utilizzare una licenza Creative Commons condizionata. Questi condizionati le licenze possono, ad esempio, richiedere l'attribuzione al

proprietario originale, impedirne la riproduzione esatta o consentirne l'uso solo per contesti non commerciali. Ciò può anche aiutare a soddisfare i ricercatori che cercano di caricare set di dati che incorporano dati di terze parti che detengono una licenza più restrittiva, poiché un "set di dati combinato adotterà le condizioni più restrittive delle sue parti componenti".

Se l'NRC segue questa strada per dare agli utenti il scelta di una licenza più ristretta, trasferirebbe anche alcune responsabilità agli utenti - o allo stesso NRC - facendo affidamento sugli utenti per rispettare la licenza. Gli approcci all'esecuzione varierebbero, a seconda dell'entità della responsabilità nell'esecuzione e, per estensione, la responsabilità che l'NRC cerca di assumersi. Ad esempio, nell'Harvard Dataverse, se un uploader decide di rinunciare a una licenza aperta predefinita e perseguire il proprio accordo di licenza personalizzato sui dati caricati, le Condizioni generali d'uso di Dataverse assolvono questo particolare cloud dalle responsabilità di applicazione delle risorse affermando che "non ha alcun obbligo di aiutare o sostenere nessuna delle parti dell'Accordo nell'esecuzione o nell'applicazione dell'Uso dei dati

Termini dell'accordo."³⁵

C. STATO ATTUALE DEI QUADRI ETICI AI

I quadri etici dell'IA (o principi, linee guida) tentano di affrontare le preoccupazioni etiche relative allo sviluppo, alla distribuzione e all'uso dell'IA all'interno delle potenziali organizzazioni. Discutiamo brevemente l'attuale panorama dei quadri etici dell'IA, pur osservando che questo è ancora un argomento emergente senza un ampio consenso.

Tra il 2015 e il 2020, governi, tecnologia aziende, organizzazioni internazionali, organizzazioni professionali e ricercatori di tutto il mondo hanno pubblicato circa 117 documenti relativi all'etica dell'IA. preminenza di concetti essenzialmente contestati nell'etica dell'IA - vale a dire parole come equità, equità, privacy che hanno significati diversi per pubblici diversi³⁸ - così come la mancanza di storia professionale vincolante e meccanismi di responsabilità, questi quadri sono spesso di alto livello e

autoregolamentazione, ponendo poca minaccia a potenziali violazioni della condotta etica.³⁹

Quadri federali

Negli Stati Uniti non esiste una guida centrale quadro sullo sviluppo responsabile e l'applicazione dell'IA in tutta la Confederazione. Alcune agenzie governative hanno adottato o sono in procinto di adottare il proprio framework di intelligenza artificiale, mentre altre non hanno pubblicato tali linee guida. Di seguito sono pubblicati quadri etici federali sull'IA a partire da agosto 2021:

- Dopo 15 mesi di deliberazione con dirigente Esperti di IA, nel febbraio 2020 il Dipartimento della Difesa (DOD) ha adottato una serie di principi etici per l'uso dell'IA che si allineano con l'esistente Missione DOD e parti interessate.⁴⁰
- La General Services Administration (GSA), incaricata dall'Office of Management and Budget (OMB) nel Piano d'azione Federal Data Strategy 2020, ha sviluppato un Data Ethics Framework nel febbraio 2020 per aiutare il personale federale a prendere decisioni etiche mentre acquisiscono, gestiscono, e utilizzare i dati.⁴¹
- Il Government Accountability Office (GAO) ha sviluppato un quadro di responsabilità dell'IA nel giugno 2020 per le agenzie federali e altre entità coinvolte nella progettazione, sviluppo, implementazione e monitoraggio continuo dei sistemi di IA per contribuire a garantire la responsabilità e l'uso responsabile dell'IA.⁴²
- L'Ufficio del direttore dell'intelligence nazionale (ODNI) ha pubblicato i Principi di etica dell'IA per la comunità dell'intelligence nel luglio 2020 per guidare lo sviluppo etico e l'uso dell'IA da parte della comunità dell'intelligence (IC) per risolvere i problemi di intelligence.⁴³
- La Commissione per la Sicurezza Nazionale sull'Artificiale Intelligence (NSCAI) ha pubblicato una serie di best practice nel luglio 2020 (successivamente riviste e

integrato nella relazione finale 2021 della Commissione) per le agenzie fondamentali per la sicurezza nazionale da attuare come paradigma per lo sviluppo responsabile e la messa in campo dei sistemi di IA.⁴⁴

Mentre questi quadri possono aiutare a guidare l'approccio dell'NRC all'etica, ci asteniamo dal raccomandare un quadro specifico per diversi motivi. In primo luogo, nonostante le crescenti richieste di etica applicata nella comunità dell'IA, lo sviluppo di un quadro etico per l'IA è ancora un'area emergente. La mancanza di uno standard governativo unificato pone sfide all'istituzione della revisione etica dell'NRC processi.

In secondo luogo, ci sono, infatti, differenze significative tra i quadri etici pubblicati dalle varie agenzie federali. Ad esempio, NSCAI ha delineato le differenze tra le sue pratiche raccomandate e quelle di DOD e IC.⁴⁵ Inoltre, tra i cinque quadri di cui sopra, il quadro GSA si è concentrato solo sulla condotta etica dei dipendenti federali quando si tratta di dati, mentre altri si sono concentrati sullo sviluppo etico e l'applicazione di sistemi di intelligenza artificiale in particolare.

In terzo luogo, il quadro etico per l'adozione della tecnologia IA può essere diverso da un quadro per la valutazione della ricerca. La maggior parte delle agenzie federali sviluppa framework per guidare l'uso di soluzioni basate sull'intelligenza artificiale per attività specifiche dell'agenzia. Ad esempio, i principi etici del DOD si applicano solo ai sistemi IA specifici per la difesa da combattimento o non da combattimento. ricerca prevista per il NRC. Il lavoro sui quadri può tuttavia fornire un utile punto di partenza per il processo etico di NRC.

D. PERSONALE E COMPETENZE

Come osservato in tutto il Libro bianco, il successo di l'NRC dipenderà dalle risorse umane, sia all'interno dell'NRC che attraverso il governo, per risolvere le numerose sfide che l'NRC promette di affrontare. Mentre ci asteniamo dal fornire un organigramma, elenchiamo le dimensioni in cui il personale e le competenze saranno fondamentali per il successo dell'NRC. Questo elenco non vuole essere esaustivo, ma per evidenziare l'importanza vitale dell'essere umano risorse.

Aree Risorse Umane

- Informatica °
 - Amministratori di sistema °
 - Ingegneri di data center °
 - Ingegneri del software di ricerca °
 - Sviluppatori di applicazioni di ricerca
- Dati
 - ° Responsabili dei dati
 - ° Rapporti con le agenzie ° Data architects
 - ° Scienziati dei dati
- concedere amministratori
- Funzionari appaltanti •
 - Personale di supporto e formazione
- Personale privacy (tecnico e legale) •
 - Personale etico
- Personale di sicurezza informatica

Note di chiusura

Sintesi

1 Klaus Schwab, La quarta rivoluzione industriale (2016).

2 Tae Yano & Moonyoung Kang, Approfondimento di Wikipedia nell'elaborazione del linguaggio naturale, Carnegie Mellon U. (2008), <https://www.cs.cmu.edu/~taey/pub/wiki.pdf>.

3 Cfr., ad esempio, Anthony Alford, Google Trains Two Billion Parameter AI Vision Model, InfoQ (22 giugno 2021), <https://www.infoq.com/news/2021/06/google-vision-transformer/>; Anthony Alford, OpenAI annuncia il modello di linguaggio AI GPT-3 con 175 miliardi di parametri, InfoQ (2 giugno 2020), <https://www.infoq.com/news/2020/06/openai-gpt3-language-model/>.

4 AlphaGo, DeepMind (2021), <https://deepmind.com/research/case-studies/alphago-the-story-so-far/>.

5 Benjamin F. Jones & Lawrence H. Summers, A Calculation of the Social Returns to Innovation (Nat'l Bureau of Econ. Research, Working Paper No.

27863, 2020); JG Tewksbury, MS Crandall & WE Crane, Misurare i vantaggi sociali dell'innovazione, 209 Sci. Mag. 658-62 (1980); si veda anche Accademie Nazionali di Scienze, Ingegneria e Medicina, Rendimenti degli investimenti federali nel sistema dell'innovazione (2017)

6 Stuart Zweben e Betsy Bizot, Sondaggio Taulbee 2019: il totale delle iscrizioni CS agli studenti universitari aumenta di nuovo, ma con un numero inferiore di nuovi studenti; La produzione del dottorato si riprende dal calo dell'anno scorso (2019).

7 Jathan Sadowski, When Data is Capital: Datafication, Accumulation, and Extraction, 2019 Big Data & Soc'y 1 (2019).

8 Amy O'Hara & Carla Medalia, Condivisione dei dati nel sistema statistico federale: impedimenti e possibilità, 675 Annali Am. Acad. pol. & Soc. Sci. 138, 140-41 (2018).

9 National Security Comm'n on Artificial Intelligence, Final Report 186 (2021).

10 Stan. U.Inst. per l'Intelligenza Artificiale incentrata sull'uomo, 2021 Rapporto sull'Indice di Intelligenza Artificiale 118 (2021).

11 id.

12 id.

13 Neil C. Thompson, Shuning Ge & Yash M. Sherry, Building the Algorithm Commons: Who Discovered the Algorithms that Underpin Computing in the Modern Enterprise?, 11 Global Strategy J. 17-33 (2020).

14 Si veda, ad esempio, US Gov't Accountability Office, Federal Agencies Need to Address Aging Legacy Systems (2016); Ufficio per la responsabilità del governo degli Stati Uniti, cloud computing: le agenzie hanno aumentato l'utilizzo e realizzato vantaggi, ma i dati sui costi e sui risparmi devono essere monitorati meglio (2019).

15 David Freeman Engstrom, Daniel E. Ho, Catherine M. Sharkey e Mariano-Florentino Cuéllar, Government by Algorithm: Artificial Intelligence in Federal Administrative Agencies 6, 71-72 (2020).

16 William M. (Mac) Thornberry National Defense Authorization Act per l'anno fiscale 2021, Pub. L. n. 116-283, § 5106.

17 L'amministrazione Biden lancia la National Artificial Intelligence Research Resource Task Force, The White House (10 giugno 2021), <https://www.whitehouse.gov/ostp/news-updates/2021/06/10/the-biden-administration-launches-the-national-artificial-intelligence-research-resource-task-force/>.

18 William M. (Mac) Thornberry National Defense Authorization Act per l'anno fiscale 2021, Pub. L. n. 116-283, § 5107 (g).

19 National Security Comm'n on Artificial Intelligence, supra nota 9, at 191.

20 Cfr., ad esempio, Cloudbank, <https://www.cloudbank.org>; Scheda informativa: National Secure Data Service Act promuove la condivisione responsabile dei dati nel governo, Data Coalition (13 maggio 2021), <https://www.datacoalition.org/fact-sheet-national-secure-data-service-act-advances-responsible-sharing-of-data-in-government/>.

21 Steve Lohr, Universities and Tech Giants Back National Cloud Computing Project, NY Times (30 giugno 2020), <https://www.nytimes.com/2020/06/30/technology/national-cloud-computing-project.html>; John Etchemendy e Fei-Fei Li, National Research Cloud: garantire la continuazione dell'innovazione americana, Stan. U.Inst. for Human-Centered Artificial Intelligence, (28 marzo 2020), <https://hai.stanford.edu/news/national-research-cloud-secured-continuation-american-innovation>.

22 Jennifer Villa e Dave Troiano, Scegliere la propria infrastruttura di deep learning: The Cloud vs. On-Prem Debate, Determinated AI (30 luglio 2020), <https://determined.ai/blog/cloud-v-onprem/>; Is HPC Going to Cost Me a Fortune?, InsideHPC (ultima visita 23 luglio 2021), <https://insidehpc.com/hpc-basic-training/is-hpc-going-to-cost-me-a-fortune/>.

23 Si veda, ad esempio, US Plans \$1.8 Billion Spend on DOE Exascale Supercomputing, HPCwire (11 aprile 2018), <https://www.hpcwire.com/2018/04/11/us-plans-1-8-billion-spend-on-doe-exascale-supercomputing/>; Federal Government, Advanced HPC (ultima visita 23 luglio 2021), <https://www.advancedhpc.com/pages/federal-government>; Gli Stati Uniti continuano a guidare il mondo nel supercalcolo, US Dep't. Energy (18 novembre 2019), <https://www.energy.gov/articles/united-states-continues-lead-world-supercomputing>.

24 Cfr. NSF finanzia cinque nuovi sistemi assegnati da XSEDE, Nat'l Sci. Trovato. (10 agosto 2020), <https://www.xsede.org/-/nsf-funds-five-new-xsede-allocated-systems>.

25 Cloudbank, supra nota 20.

26 Cfr., ad esempio, National Data Service, <http://www.nationaldataservice.org>; L'Open Science Data Cloud, <https://www.opensciencedatacloud.org>; Harvard Dataverse, <https://dataverse.harvard.edu>; FigShare, <https://figshare.com>.

27 FedRAMP, <https://www.fedramp.gov>.

28 Cfr. scheda informativa: National Secure Data Service Act promuove la condivisione responsabile dei dati nel governo, Data Coalition (13 maggio 2021), <https://www.datacoalition.org/fact-sheet-national-secure-data-service-act-advances-responsible-data-sharing-in-government/>.

29 Cfr. Administrative Data Research Facility, Coleridge Initiative, <https://coleridgeinitiative.org/adrf/> (ultima visita 26 luglio 2021).

30 Cfr. Landsat Data Access, US Geological Survey, <https://www.usgs.gov/core-science-systems/nli/landsat/landsat-data-access> (ultima visita 23 luglio 2021); Alimentato. Geographic Data Comm., La proposta di valore per le applicazioni Landsat (2014); Crista L. Straub, Stephen R. Koontz e John B. Loomis, Valutazione economica delle immagini Landsat (2019).

- 31 Cfr. Bipartisan Pol'y Ctr., Barriers to Using Government Data: Extended Analysis of the US Commission on Evidence-Based Policymaking's Survey of Federal Agencies and Offices 18-20 (2018); vedi anche US Dep't of Health & Human Services, Lo stato della condivisione dei dati negli Stati Uniti Department of Health and Human Services 4 (2018) (che descrive come i dati presso l'agenzia sono "in gran parte conservati in silos con una mancanza di consapevolezza organizzativa di quali dati vengono raccolti in tutto il Dipartimento e come richiedere l'accesso").
- 32 Legge sulla privacy, 5 USC § 552a (1974).
- 33 Michael S. Bernstein et al., ESR: Ethics and Society Review of Artificial Intelligence Research, Cornell. U. (9 luglio 2021), <https://arxiv.org/pdf/2106.11521.pdf>.
- 34 Courtenay R. Bruce et al., An Embedded Model for Ethics Consultation: Characteristics, Outcomes, and Challenges, 5 AJOB Empirical Bioethics 8 (2014).

introduzione

- 1 Invito all'azione del cloud di ricerca nazionale, Stan. U.Inst. Per l'intelligenza artificiale centrata sull'uomo, <https://hai.stanford.edu/national-research-cloud-joint-letter>.
- 2 Cfr. id.; John Etchemendy e Fei-Fei Li, National Research Cloud: garantire la continuazione dell'innovazione americana, Stan. U.Inst. Per l'intelligenza artificiale centrata sull'uomo (28 marzo 2020), <https://hai.stanford.edu/news/national-research-cloud-ensuring-continuation-american-innovation>.
- 3 William M. (Mac) Thornberry National Defense Authorization Act per l'anno fiscale 2021, Pub. L. n. 116-283, § 5106.
- 4 L'amministrazione Biden lancia la National Artificial Intelligence Research Resource Task Force, The White House (10 giugno 2021), <https://www.whitehouse.gov/ostp/news-updates/2021/06/10/the-biden-administration-launches-the-national-artificial-intelligence-research-resource-task-force/>.
- 5 Legge sulla privacy del 1974, 5 USC § 552a (2012).
- 6 Foundations for Evidence-Based Policymaking Act del 2017, Pub. L. n. 115-435, 132 Stat. 5529 (2019).
- 7 Cfr. scheda informativa: National Secure Data Service Act promuove la condivisione responsabile dei dati nel governo, Data Coalition (13 maggio 2021), <https://www.data-coalition.org/fact-sheet-national-secure-data-service-act-anticipi-responsabile-condizione-dei-dati-nel-governo/>.
- 8 Cfr., ad esempio, riconoscimento facciale e legge sulla moratoria della tecnologia biometrica, S. 4084, 116th Cong. (2020); Bhaskar Chakravorti, La "rivoluzione antitrust" di Biden trascura l'IA — a rischio degli americani, Wired (27 luglio 2021), <https://www.wired.com/story/opinion-bidens-antitrust-revolution-overlooks-ai-at-american-pericolo/>.

Capitolo 1

- 1 Cfr. Stephen Breyer, Regulation and Its Reform (1982); Clifford Winston, Fallimento del governo contro fallimento del mercato (2006).
- 2 società più grandi per capitalizzazione di mercato, capitalizzazione di mercato delle società (2021), <https://companiesmarketcap.com>.
- 3 Stan. U.Inst. per l'Intelligenza Artificiale incentrata sull'uomo, 2021 Rapporto sull'Indice di Intelligenza Artificiale 93 (2021).
- 4 Si veda, ad esempio, Mary L. Gray e Siddarth Suri, Ghost Work: How to Stop Silicon Valley from Building a New Global Underclass (2019); Craig Webster e Stanislav Ivanov, robotica, intelligenza artificiale e natura in evoluzione del lavoro 132-35 (2019); Weiyu Wang & Keng Siau, Intelligenza artificiale, apprendimento automatico, automazione, robotica, futuro del lavoro e futuro dell'umanità: un'agenda di revisione e ricerca, 30 J. Database Mgmt. 61 (2019).
- 5 AlphaFold: una soluzione a una grande sfida di biologia di 50 anni, DeepMind (30 novembre 2020), <https://deepmind.com/blog/article/alphafold-a-solution-to-a-50-year-vecchia-grande-sfida-in-biologia>.
- 6 Tanha Talaviya et al., Attuazione dell'intelligenza artificiale in agricoltura per l'ottimizzazione dell'irrigazione e l'applicazione di pesticidi ed erbicidi, 4 Intelligenza artificiale in agricoltura 58 (2020).
- 7 Greg Allen e Taniel Chan, Intelligenza artificiale e sicurezza nazionale, Harv. Kennedy Sch. Belfer Ctr. (luglio 2017), <https://www.belfercenter.org/publication/artificial-intelligence-and-national-security>.
- 8 Stan. U.Inst. per l'intelligenza artificiale centrata sull'uomo, supra nota 3.
- 9 Jeffrey Ding, Decifrare il sogno dell'IA in Cina (2018).
- 10 Fugaku viene ampiamente utilizzato per iniziative di ricerca sull'IA. Vedere Atsushi Nukariya et al., HPC and AI Initiatives for Supercomputer Fugaku and Future Prospects, Fujitsu (11 novembre 2020), <https://www.fujitsu.com/global/about/resources/publications/technicalreview/2020-03/articolo09.html>.
- 11 Eng'g & Physical Sciences Research Council, The Impact of HECToR (2014).
- 12 Joshua New, Perché gli Stati Uniti hanno bisogno di una strategia nazionale di intelligenza artificiale e come dovrebbe essere (2018).
- 13 Maggie Miller, La Casa Bianca istituisce il National Artificial Intelligence Office, The Hill (12 gennaio 2021), <https://thehill.com/policy/cybersecurity/533922-white-house-establishes-national-artificial-intelligence-office>.
- 14 Cfr. comunicazione azione rapida. sull'informatica strategica, aggiornamento dell'iniziativa nazionale sull'informatica strategica: pioniere del futuro dell'informatica (2019), <https://www.nitrd.gov/pubs/National-Strategic-Computing-Initiative-Update-2019.pdf>.
- 15 Commissione nazionale per la sicurezza sull'intelligenza artificiale, rapporto finale (2021).
- 16 The COVID-19 High Performance Computing Consortium, COVID-19 HPC Consortium, <https://covid19-hpc-consortium.org>.
- 17 Cfr. Aaron L. Friedberg, Science, the Cold War, and the American State, 20 Diplomatic Hist. 107, 112 (1996); Sean Pool e Jennifer Erickson, L'alto ritorno sull'investimento per la ricerca finanziata con fondi pubblici, Ctr. per l'Am. Progress (10 dicembre 2012), <https://www.americanprogress.org/issues/economy/reports/2012/12/10/47481/the-high-return-on-investment-for-publicly-funded-research/>.
- 18 Peter L. Singer, Innovazioni sostenute a livello federale: 22 esempi di importanti progressi tecnologici derivanti dal sostegno federale alla ricerca 14-15 (2014).
- 19 Consiglio Nazionale delle Ricerche, Sostegno del governo alla ricerca informatica 136-55 (1999).
- 20 National Security Comm'n on Artificial Intelligence, supra nota 15, a 185.
- 21 Philippe Aghion, Benjamin F. Jones e Charles I. Jones, Intelligenza artificiale e crescita economica, in The Economics of Artificial Intelligence: An Agenda 237 (2019).

- 22 Ian Moll, Il mito della quarta rivoluzione industriale, 68 Theoria 1 (2021); si veda anche Tim Unwin, 5 Problems with 4th Industrial Revolution, ICT Works (23 marzo 2019), <https://www.ictworks.org/problems-fourth-industrial-revolution/>.
- 23 Si veda, ad esempio, Geoffrey A. Manne e Joshua D. Wright, Google and the Limits of Antitrust: The Case Against the Antitrust Case Against Google, 34 Harv. J.L. & Pub. Pol'y 1 (2011); Lina M. Khan, Il paradosso antitrust di Amazon, 126 Yale L.J. 710 (2016).
- 24 David Patterson et al., Carbon Emissions and Large Neural Network Training, Cornell U. (23 aprile 2021), <https://arxiv.org/pdf/2104.10350.pdf>. Per essere chiari, tuttavia, lo studio ha rilevato che l'addestramento di altri modelli di PNL sofisticati ma più piccoli come Meena e T5 richiedeva rispettivamente circa 96 e 48 tonnellate di anidride carbonica. Id. Un altro studio ha rilevato che i modelli NLP all'avanguardia di addestramento hanno prodotto circa 626.000 libbre (313 tonnellate) di anidride carbonica, cinque volte le emissioni di un'auto media negli Stati Uniti. Emma Strubell, Ananya Ganesh e Andrew McCallum, Considerazioni sull'energia e le politiche per l'apprendimento profondo nella PNL, Cornell U. (2019), <https://arxiv.org/pdf/1906.02243.pdf>.
- 25 Calcola la tua impronta di carbonio, The Nature Conservancy, <https://www.nature.org/en-us/get-involved/how-to-help/carbon-footprint-calculator/>.
- 26 Studi economici in altri campi mostrano anche che l'aumento dell'accesso, dell'offerta o della qualità di determinati beni senza adeguati meccanismi di determinazione dei prezzi o interventi normativi può portare a un uso eccessivo e uno spreco. Vedi, ad esempio, Chengri Ding & Shunfeng Song, Traffic Paradoxes and Economic Solutions, 1 J. Urban Mgmt. 63 (2012) (strade e congestione del traffico); Ari Mwachofi e Assaf F. Al-Assaf, Deviazioni del mercato sanitario dal mercato ideale, 11 Sultan Qaboos Univ. Med. J. 328 (2011) (medici e qualità delle cure).
- 27 Si veda Emily M. Bender et al., On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?, 2021 Proceedings ACM Conf. su Equità, Responsabilità e Trasparenza 610 (2021).
- 28 Cfr. Joy Buolamwini e Timnit Gebru, Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification, 81 Proceeding Machine Learning Res. 1 (2018); Inioluwa Deborah Raji e Joy Buolamwini, Auditing attuabile: indagare sull'impatto della denominazione pubblica dei risultati delle prestazioni distorte dei prodotti commerciali di IA, 2019 Proceedings AAAI/ACM Conf. su AI, etica e società 429 (2019).
- 29 Cfr. Virginia Eubanks, Automating Inequality (2018); Cathy O'Neil, Armi di distruzione matematica (2016).
- 30 Cfr. Christopher Whyte, Deepfake News: AI-Enabled Disinformation as a Multi-Level Public Policy Challenge, 5 J. Cyber Pol'y 199 (2020); Jeffrey Dastin, Amazon Scrapes Secret AI Recruiting Tool che ha mostrato pregiudizi nei confronti delle donne, Reuters (10 ottobre 2018), <https://www.ceiving.com/article/us-amazon-com-jobs-automation-insight/amazon-scraps-strumento-di-reclutamento-segreto-ai-che-ha-mostrato-pregiudizi-contro-le-donne-idUSKCN1MK08G>; James Vincent, Google ha "corretto" il suo algoritmo razzista rimuovendo i gorilla dalla sua tecnologia di etichettatura delle immagini, The Verge (12 gennaio 2018), <https://www.theverge.com/2018/1/12/16882408/google-racist-algoritmo-di-riconoscimento-foto-gorilla-ai>.
- 31 Kate Crawford, Atlas of AI: Power, Politics, and the Planetary Costs of Artificial Intelligence 211 (2021) ("I sistemi di intelligenza artificiale sono costruiti per vedere e intervenire nel mondo in modi che avvantaggiano principalmente gli stati, le istituzioni e le serve. In questo senso, i sistemi di intelligenza artificiale sono espressioni di potere che emergono da forze economiche e politiche più ampie, create per aumentare i profitti e centralizzare il controllo per coloro che li esercitano.).
- 32 Elizabeth Gibney, L'intelligenza artificiale autodidatta è la migliore di sempre al gioco di strategia Go, Nature (18 ottobre 2017), <https://www.nature.com/news/self-taught-ai-is-best-yet-at-game-go-1.22858>.
- 33 Bill Schackner, Prestigious Computer Science School di Carnegie Mellon ha un nuovo leader, Pittsburgh Post-Gazette (8 agosto 2019), <https://www.postgazette.com/news/education/2019/08/08/Carnegie-Mellon-University-computer-science-Martial-Hebert-dean-artificial-intelligence-google-robotics/stories/201908080096>.
- 34 Bipartisan Pol'y Ctr., Cementing American Artificial Intelligence Leadership: AI Research & Development (2020).
- 35 Nur Ahmed & Muntasir Wahed, La de-democratizzazione dell'IA: Deep Learning e il divario informatico nella ricerca sull'intelligenza artificiale, Cornell U. (22 ottobre 2020), <https://arxiv.org/pdf/2010.15581.pdf>.
- 36 Id.
- 37 Fei-Fei Li, America's Global Leadership in Human-Centered AI Can't Come From Industry Alone, The Hill (6 luglio 2021), <https://thehill.com/opinion/technology/561638-americas-global-leadership-in-human-centered-ai-cant-come-from-industry?rl=1>.
- 38 Cade Metz, I ricercatori di intelligenza artificiale guadagnano più di 1 milione di dollari, anche in un'organizzazione no profit, NY Times (19 aprile 2018), <https://www.nytimes.com/2018/04/19/technology/artificial-intelligence-stipendi-openai.html>.
- 39 Stan. U. Inst. per l'intelligenza artificiale centrata sull'uomo, supra nota 3, a 118.
- 40 Michael Gofman e Zhao Jin, Artificial Intelligence, Education, and Entrepreneurship, SSRN (17 settembre 2019), https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3449440.
- 41 Jathan Sadowski, When Data is Capital: Datafication, Accumulation, and Extraction, 2019 Big Data & Soc'y 1 (2019).
- 42 Ad esempio, i ricercatori hanno chiesto a gran voce a Facebook di condividere alcuni dei suoi dati proprietari in modo da poter comprendere meglio l'effetto dei social media sulla politica e sul discorso sociale. Simon Hegelich, Facebook Needs to Share More with Researchers, Nature (24 marzo 2020), <https://www.nature.com/articles/d41586-020-00828-5>.
- 43 Ashlee Vance, This Tech Bubble Is Different, Bloomberg (14 aprile 2011), <https://www.bloomberg.com/news/articles/2011-04-14/this-tech-bubble-is-different>.
- 44 Amy O'Hara & Carla Medalia, Condivisione dei dati nel sistema statistico federale: impedimenti e possibilità, 675 Annali Am. Acad. pol. & Soc. Sci. 138, 140-41 (2018).
- 45 Dario Amodei & Danny Hernandez, AI and Compute, Open AI (16 maggio 2018), <https://openai.com/blog/ai-and-compute/>.
- 46 Cfr., ad esempio, Ahmed & Wahed, supra nota 35; Ian Sample, "Non possiamo competere": perché le università stanno perdendo i loro migliori scienziati di intelligenza artificiale, The Guardian (1 novembre 2017), <https://www.theguardian.com/science/2017/nov/01/cant-compete-università-perdere-i-migliori-scienziati>.
- 47 Neil C. Thompson, Shuning Ge & Yash M. Sherry, Building the Algorithm Commons: Who Discovered the Algorithms that Underpin Computing in the Modern Enterprise?, 11 Global Strategy J. 17-33 (2020).
- 48 Minkyung Baek, RoseTTAFold: previsione accurata della struttura proteica accessibile a tutti, U. Wash. Inst. per Protein Design (15 luglio 2021), <https://www.ipd.uw.edu/2021/07/rosettafold-accurate-protein-structure-prediction-accessible-to-all/>; Minkyung Baek et al., Previsione accurata delle strutture proteiche e delle interazioni utilizzando una rete neurale a tre tracce, Sci. Mag. (15 luglio 2021), <https://science.sciencemag.org/content/sci/early/2021/07/19/science.abj8754.full.pdf>.
- 49 Come la diplomazia ha contribuito a porre fine alla corsa al sequenziamento del genoma umano, Nature (24 giugno 2020), <https://www.nature.com/articles/d41586-020-01849-w>.
- 50 Joel Klinger et al., A Narrowing of AI Research?, Cornell U. (17 novembre 2020), <https://arxiv.org/pdf/2009.10385.pdf>.
- 51 Id.
- 52 Alex Tamkin et al., Understanding the Capabilities, Limitations, and Societal Impact of Large Language Models, Cornell U. (4 febbraio 2021), <https://arxiv.org/pdf/2102.02503.pdf>.

53 Quei 5 erano Boston, San Francisco, San Jose, Seattle e San Diego. Vedi Robert D. Atkinson, Mark Muro e Jacob Whiton, The Case for Growth Centers: How to Spread Tech Innovation Across America, Brookings (9 dicembre 2019), <https://www.brookings.edu/research/growth-centers-come-diffondere-l'innovazione-tecnologica-in-tutta-l'America/>.

54 Intervista con il Professor Erik Brynjolfsson, Direttore, Stanford Digital Economy Lab (2021).

55 Solon Barocas e Andrew D. Selbst, L'impatto disparato dei Big Data, 104 Cal. L. Rev. 671 (2016).

capitolo 2

1 Lo status di "Principal Investigator" può variare da università a università, ma in genere rappresenta la facoltà principale che è idonea a supervisionare progetti di ricerca presso le proprie istituzioni di origine.

2 Vedi Beth Jensen, AI Index Diversity Report: An Unmoving Needle, Stan. U. Inst. per l'intelligenza artificiale centrata sull'uomo (3 maggio 2021), <https://hai.stanford.edu/news/ai-index-diversity-report-unmoving-needle>.

3 Per una prospettiva, ad esempio, sull'importanza della modellazione e della simulazione in fisica, vedere Karen E. Wilcox, Omar Ghattas & Patrick Heimbach, The Imperative of Physics-Based Modeling and Inverse Theory in Computational Science, 1 Nature Comp. Sci. 166 (2021).

4 15 USC § 9415 (sottolineatura aggiunta).

5 id. (enfasi aggiunta).

6 I resoconti contemporanei confermano questo obiettivo centrale. La National Security Commission on AI, ad esempio, descrive la proposta come "fornire a ricercatori e studenti verificati un accesso sovvenzionato a risorse di calcolo scalabili" con un riferimento specifico al "divario di calcolo" che ha lasciato "università di livello medio e inferiore" [mancando] delle risorse necessarie per una ricerca IA all'avanguardia". Nat'l Security Comm'n on AI, Final Report 191, 197 (2021) (corsivo aggiunto). Dopo l'annuncio della legislazione NRC, Jeff Dean, SVP di Google Research e Google Health, ha osservato: "Una risorsa nazionale di ricerca sull'intelligenza artificiale contribuirà ad accelerare i progressi degli Stati Uniti nell'intelligenza artificiale e nelle tecnologie avanzate fornendo ai ricercatori accademici l'accesso alle risorse di cloud computing necessarie per esperimenti su larga scala". Brandi Vincent, Il Congresso è sempre più vicino alla creazione di un cloud nazionale per la ricerca sull'IA, NextGov (2 luglio 2020), <https://www.nextgov.com/emerging-tech/2020/07/congress-inches-closer-creating-national-cloud-ai-research/166624/> (corsivo aggiunto). Altri hanno suggerito che i "ricercatori" sotto NRC potrebbero includere individui di piccole imprese, start-up, organizzazioni non profit e alcune aziende tecnologiche.

Un co-sponsor della legislazione, ad esempio, ha suggerito che le risorse NRC dovrebbero essere fornite a "sviluppatori" e "imprenditori". Portman, Heinrich Introduce Bipartisan Legislation to Develop National Cloud Computer for AI Research, Rob Portman, senatore degli Stati Uniti per l'Ohio (4 giugno 2020), <https://www.portman.senate.gov/newsroom/press-releases/portman-heinrich-introduce-bipartisan-legislation-develop-national-cloud>.

7 domande frequenti sulle piccole imprese, US Small Bus. Amministratore Ufficio dell'Avv. (ottobre 2020), <https://cdn.advocacy.sba.gov/wp-content/uploads/2020/11/05122043/Small-Business-FAQ-2020.pdf>.

8 Louise Balle, Information on Small Business Startups, Houston Chron., <https://smallbusiness.chron.com/information-small-business-startups-2491.html>.

9 Tali entità potrebbero potenzialmente collaborare con partner accademici e l'NRC dovrebbe ovviamente anche stabilire regole sull'idoneità dei collaboratori.

10 Lo status di PI fornisce un livello di standardizzazione tra i docenti rispetto ad altri parametri, come il percorso di ruolo o la designazione come docente di ricerca. Ad esempio, l'Università del Michigan nomina persone focalizzate sulla ricerca a tempo pieno come "facoltà di ricerca", che non è una posizione di ruolo. Al contrario, i docenti di ricerca della Purdue possono beneficiare del ruolo di ruolo. Distinto dalla categorizzazione utilizzata da entrambe le università, il MIT designa i ricercatori a tempo pieno come "personale accademico" piuttosto che come docenti. Tutti e tre i tipi di ricercatori, tuttavia, si qualificano per lo status di ricercatore principale presso le rispettive università. Alcune università vanno oltre fornendo lo status di PI temporaneo a persone con status di PI non affiliate all'università per un singolo progetto (comprese tutte e tre le università menzionate in precedenza).

11 Richieste di risorse per la comunità e l'istruzione, CloudBank, <https://www.cloudbank.org/training/cloudbank-community#toc-eligibilit-36nfpqrS>.

12 Richiedi un account, Compute Canada, <https://www.computeCanada.ca/research-portal/account-management/apply-for-an-account/>.

13 Nat'l Sci. Bd., Indicatori scientifici e ingegneristici 2016, Ricerca e sviluppo accademici 72 (2016).

14 id.

15 Iscrizione al college negli Stati Uniti dal 1965 al 2019 e proiezioni fino al 2029 per i college pubblici e privati, Statista (gennaio 2021), <https://www.statista.com/statistics/183995/us-college-enrollment-and-projections-in-public-and-private-institutions/>.

16 Colaboratory – Domande frequenti, Google, <https://research.google.com/colaboratory/faq.html>.

17 Utilizzo massimo settimanale della GPU, Kaggle (2019), <https://www.kaggle.com/general/108481>.

18 Richieste di risorse per la comunità e l'istruzione, supra nota 11.

19 Revisione del merito: perché dovresti offrirti volontario per servire come revisore NSF, Nat'l Sci. Trovato., https://www.nsf.gov/bfa/dias/policy/merit_review/reviewer.jsp#1.

20 Vedere XSEDE Campus Champions, XSEDE, <https://www.xsede.org/community-engagement/campus-champions>.

21 Compute Canada, ad esempio, fornisce l'accesso al 15% dei PI a una maggiore capacità di calcolo sulla base di una competizione di merito. Nel 2021, Compute Canada ha completato la revisione di 650 proposte di ricerca in circa cinque mesi con solo 80 revisori volontari delle istituzioni accademiche canadesi per valutare il merito scientifico della proposta. Competizioni per l'allocazione delle risorse, Compute Canada, <https://www.computeCanada.ca/research-portal/accessing-resources/resource-allocation-competitions/>; Risultati del concorso 2021 Resource Allocations, Compute Canada, <https://www.computeCanada.ca/research-portal/accessing-resources/resource-allocation-competitions/rac-2021-results/>. Confronta questo con CloudBank, che alloca le risorse di calcolo sfruttando il processo di amministrazione delle sovvenzioni di NSF: nel 2019, NSF aveva bisogno di 30.000 revisori volontari per gestire oltre 40.000 proposte, con ciascuna proposta che richiedeva circa 10 mesi per l'elaborazione dall'inizio alla fine. Nat'l. Sci. Found., Merit Review Process: Fiscal Year 2019 Digest (2020); Proposta NSF e processo di aggiudicazione, Nat'l Sci. Trovato., https://www.nsf.gov/attachments/116169/public/nsf_proposal_and_award_process.pdf.

22 Un'altra questione di confine sarà l'allocazione delle risorse ai PI convenzionati sia con università che con aziende private. Come impostazione predefinita, le risorse dell'NRC dovrebbero essere destinate a progetti accademici e non a sovvenzionare il lavoro svolto a titolo di ricercatore privato.

23 Concorsi per l'allocazione delle risorse, supra nota 21.

24 Semplificazione dei servizi cloud, Sci. Node (2 dicembre 2019), <https://sciencenode.org/feature/An%20easier%20cloud.php>.

25 Domande frequenti (FAQ), CloudBank, <https://www.cloudbank.org/faq>.

26 Semplificazione dei servizi cloud, supra nota 25.

27 id.

28 id.

29 Domande frequenti (FAQ), supra nota 26.

30 domande frequenti (FAQ) per la definizione del budget per le risorse di cloud computing tramite CloudBank nelle proposte NSF, Nat'l Sci. Trovato., <https://www.nsf.gov/pubs/2020/nsf20108/nsf20108.jsp>.

31 Semplificazione dell'accesso alle risorse cloud per i ricercatori: CloudBank, Amazon Web Serv. (16 novembre 2020), <https://aws.amazon.com/blogs/publicsector/simplifying-access-cloud-resources-researchers-cloudbank/>.

32 Richieste di risorse per la comunità e l'istruzione, supra nota 11.

33 Larry Dignan, AWS Cloud Computing Ops, data center, 1,3 milioni di server che creano un volano di efficienza, ZDNet (17 giugno 2016), <https://www.zdnet.com/article/aws-cloud-computing-ops-data-centers-1-3-million-servers-creating-efficiency-flywheel/>; Rich Miller, Ballmer: Microsoft ha 1 milione di server, Data Ctr.

Knowledge (15 luglio 2013), <https://www.datacenterknowledge.com/archives/2013/07/15/ballmer-microsoft-has-1-million-servers>; Daniel Oberhaus, Amazon, Google, Microsoft: Ecco chi ha il cloud più verde, Wired (18 dicembre 2019), <https://www.wired.com/story/amazon-google-microsoft-green-clouds-and-hyperscale-Centri-dati/>; Russell Brandom, Mappatura dell'impero invisibile dei server di Amazon, The Verge (10 maggio 2019), <https://www.theverge.com/2019/5/10/18563485/amazon-web-services-internet-location-map-data-centro>.

34 Cfr., ad esempio, AWS Pricing, Amazon Web Services, <https://aws.amazon.com/pricing/>; Panoramica dei concetti di fatturazione cloud, Google Cloud, <https://cloud.google.com/billing/docs/concepts>; Prezzi di Azure, Azure, <https://azure.microsoft.com/en-us/pricing/#product-pricing>.

35 Le grandi università di ricerca negoziano già accordi aziendali con fornitori di servizi cloud.

36 What We Do, XSEDE, <https://www.xsede.org/about/what-we-do> (ultima visita 19 settembre 2021).

37 XSEDE Overall Organization, XSEDE Wiki, <https://confluence.xsede.org/display/XT/XSEDE+Overall+Organization> (ultima visita 19 settembre 2021).

38 XSEDE Allocations Info & Policies, XSEDE, <https://portal.xsede.org/allocations/policies> (ultima visita 19 settembre 2021).

39 id.

40 id.

41 Startup Allocations, XSEDE, <https://portal.xsede.org/allocations/startup> (ultima visita 19 settembre 2021).

42 id.

43 id.

44 id.

45 Research Allocations, XSEDE, <https://portal.xsede.org/allocations/research> (ultima visita 19 settembre 2021).

46 id.

47 id.

48 id.

49 XSEDE Allocations Info & Policies, supra nota 36.

50 XSEDE Campus Champions, supra nota 20.

51 id.

52 id.

53 XSEDE as a Collaborator on Proposals, XSEDE, <https://www.xsede.org/about/collaborating-with-xsede> (ultima visita 19 settembre 2021).

54 COVID-19 HPC Consortium, XSEDE, <https://www.xsede.org/covid19-hpc-consortium> (ultima visita 19 settembre 2021).

55 Amazon, ad esempio, ha introdotto le sue istanze P4, P3 e P2 rispettivamente nel 2020, 1997 e 1996. Frederic Lardinois, AWS lancia le sue istanze GPU di nuova generazione con 8 GPU Nvidia A100 Tensor Core, TechCrunch (2 novembre 2020), <https://social.techcrunch.com/2020/11/02/aws-launches-its-istanze-gpu-di-nuova-generazione/>; Ian C. Schafer, Amazon Elastic Compute Cloud P3 Lanciato insieme a NVIDIA GPU Cloud, SD Times (26 ottobre 2017), <https://sdtimes.com/ai/amazon-elastic-compute-cloud-p3-launched-alongside-nvidia-gpu-cloud/>; Jeff Barr, Nuovo tipo di istanza P2 per Amazon EC2 – Fino a 16 GPU, Amazon Web Services (29 settembre 2016), <https://aws.amazon.com/blogs/aws/new-p2-instance-type-for-amazon-ec2-fino-a-16-gpus/>. Gli anni di introduzione delle istanze P4 e P3 si allineano con il rilascio delle più recenti GPU per data center per uso generico di NVIDIA.

56 Si veda, ad esempio, Sarah Wang e Martin Casado, The Cost of Cloud, a Trillion Dollar Paradox, Andreessen Horowitz (27 maggio 2021), <https://a16z.com/2021/05/27/cost-of-cloud-paradox-market-cap-cloud-lifecycle-scale-growth-repatriation-optimization/>.

57 Preston Smith et al., Community Clusters or the Cloud: Continuing Cost Assessment of On-Premises and Cloud HPC in Higher Education, 2019 Proceedings Practice & Experience Advanced Res. Computing on Rise of the Machines 1 (2019). Il costo ammortizzato include il costo di elaborazione annuale, il costo dell'hardware sovvenzionato e i costi dell'alimentazione, ma non include i costi del personale, in quanto tali costi sono fissi e sarebbero ricorrenti indipendentemente dal fatto che un cluster esistesse fisicamente in sede o nel cloud. Id.

58 Craig A. Stewart et al., Return on Investment for Three Cyberinfrastructure Facilities: A Local Campus Supercomputer; il Jetstream Cloud System finanziato dalla NSF; e XSEDE, 11 Int'l Conf. su Utilità e Cloud Computing 223 (2018).

59 Srijith Rajamohan e Robert E. Settlege, Informing the On/Off-prem Cloud Discussion in Higher Education, 2020 Practice & Experience Adv. Res.

Informatica 64 (2020). Le fonti di costo includono l'hardware, i servizi software, l'amministrazione del software, l'elettricità e le strutture, ma non includono il supporto degli scienziati informatici, le licenze del software scientifico e i costi di trasferimento dei dati. Lo studio è inoltre limitato al particolare carico di lavoro cloud di Virginia Tech.

60 Jennifer Villa e Dave Troiano, Scegliere la tua infrastruttura di deep learning: il dibattito tra cloud e on-prem, Determinated AI (30 luglio 2020), <https://determined.ai/blog/cloud-v-onprem/>; HPC mi costerà una fortuna?, insideHPC, <https://insidehpc.com/hpc-basic-training/is-hpc-going-to-cost-me-a-fortune/>.

61 Intervista a Suzanne Talon, Regional Director, Compute Canada (14 gennaio 2021).

62 Compute Canada, Cloud Computing for Researchers 1 (2016), <https://www.compute-canada.ca/wp-content/uploads/2015/02/CloudStrategy2016-2019-forresearchersEXTERNAL-1.pdf>.

63 Gli Stati Uniti prevedono di spendere 1,8 miliardi di dollari per il supercomputing Exascale DOE, HPCwire (11 aprile 2018), <https://www.hpcwire.com/2018/04/11/us-plans-1-8-billion-spend-on-doe-exascale-supercomputer/>; Governo federale, Advanced HPC, <https://www.advancedhpc.com/pages/federal-government>; Gli Stati Uniti continuano a guidare il mondo nel supercomputing, Energy.gov, <https://www.energy.gov/articles/united-states-continues-lead-world-supercomputing>; Calcolo ad alte prestazioni, Energy.gov, <https://www.energy.gov/science/initiatives/high-performance-computing>.

64 Si veda, ad esempio, DOE annuncia cinque nuovi progetti energetici a LLNL, LLNL (13 novembre 2020), <https://www.llnl.gov/news/doe-announces-five-new-energy-projects-llnl>; Nuovo sistema HPCMP presso l'AFRL DSRC DoD Supercomputing Resource Center per fornire oltre nove PetaFLOPS di potenza di calcolo a

Affrontare le applicazioni di fisica, intelligenza artificiale e ML per utenti DoD, DOD HPC, https://www.hpc.mil/images/hpcdocs/newsroom/21-19_TI-21_web_announcement_AFRL_DSRC.pdf; Annuncio pubblico, DOD HPC, https://www.hpc.mil/images/hpcdocs/newsroom/awards_and_press/Hc101321D0002_PUBLIC_ANNOUNCEMENT_20210505.pdf.

65 Devin Coldewey, il supercomputer Cray da 600 milioni di dollari torreggerà al di sopra degli altri — per costruire armi nucleari migliori, TechCrunch (13 agosto 2019), <https://social.techcrunch.com/2019/08/13/600m-cray-supercomputer-will-tower-above-the-rest-to-build-better-nukes/>; CORAL-2 RFP, Oak Ridge Nat'l Laboratory (9 aprile 2018), <https://procurement.ornl.gov/rfp/CORAL2/>.

66 Si veda, ad esempio, NSF Funds Five New XSEDE-Allocated Systems, Nat'l Sci. Trovato. (10 agosto 2020), <https://www.xsede.org/-/nsf-funds-five-new-xsedeallocated-systems>.

67 Timothy Prickett Morgan, Bending The Supercomputing Cost Curve Down, The Next Platform (2 dicembre 2019), <http://www.nextplatform.com/2019/12/02/bending-the-supercomputing-cost-curve-down/>; Ben Dickson, The GPT-3 Economy, TechTalks (21 settembre 2020), <https://bdtectalks.com/2020/09/21/gpt-3-economy-business-model/>.

68 Elijah Wolfson, The US Just Retook the Title of World's Fastest Supercomputer from China, Quartz (9 giugno 2018), <https://qz.com/1301510/the-us-has-the-worlds-fastest-supercomputer-again/> il vertice-200-petaflop/.

69 novembre 2020, TOP500 (novembre 2020), <https://www.top500.org/lists/top500/2020/11/>.

70 Il Dipartimento dell'Energia degli Stati Uniti e Cray forniranno un supercomputer di frontiera da record all'ORNL, Oak Ridge Nat'l Laboratory (7 maggio 2019), <https://www.ornl.gov/news/us-department-energy-and-cray-deliver-record-setting-frontier-supercomputer-ornl>.

71 Coury Turczyn, Building an Exascale-Class Data Center, Oak Ridge Leadership Computing Facility (11 dicembre 2020), <https://www.olcf.ornl.gov/2020/12/11/building-an-exascale-class-data-center/>.

72 Don Clark, Intel Slips, and a High-Profile Supercomputer Is Delayed, NY Times (27 agosto 2020), <https://www.nytimes.com/2020/08/27/technology/intel-aurora-supercomputer.html>; Mila Jasper, 10 dei 15 dei principali progetti IT del DOD sono in ritardo, GAO Found, Nextgov (4 gennaio 2021), <https://www.nextgov.com/it-modernization/2021/01/10-15-dods-major-it-projects-are-behind-schedule-gao-found/171155/>.

73 Si veda Nattakarn Phaphoom et al., A Survey Study on Major Technical Barriers Affecting the Decision to Adopt Cloud Services, 103 J. Systems & Software 167, 171-72 (2015) (che descrive la portabilità dei dati, l'integrazione con i sistemi esistenti, la complessità della migrazione, e disponibilità come principali ostacoli all'adozione del cloud); Abdulrahman Alharthi et al., An Overview of Cloud Services Adoption Challenges in Higher Education Institutions, 2 Proceedings of the Int'l Workshop on Emerging Software as a Service & Analytics 102, 107-08 (2015) (in riconoscimento del basso tasso di cloud computing adozione nell'istruzione superiore e sottolineando che il rafforzamento sia della facilità d'uso percepita sia dell'effettiva utilità del cloud computing può aumentare il tasso di adozione).

74 Cfr. Dept of Energy, FY 2021 Budget Justification Volume 4: Science (2020).

75 Joe Weinman, Cloudonomics: il valore commerciale del cloud computing (2012).

76 OLCF supporta e gestisce le risorse di supercalcolo di ORNL, tra cui Summit e infine Frontier. Questa cifra rappresenta "le operazioni e il supporto degli utenti presso le strutture LCF, inclusi energia, spazio, contratti di locazione e personale. Id. a 37-38.

77 ACLF supporta e gestisce le risorse informatiche dell'Argonne National Laboratory, compreso il sistema Theta e, entro la fine dell'anno, il nuovo computer Aurora, un altro sistema HPC exascale del DOE. Id.

78 OLCF ha utilizzato il suo sistema Titan HPC per 7 anni. Vedi Coury Turczyn, supra nota 72. ACLF ha anche gestito il suo sistema Mira HPC per 7 anni. Il supercomputer Mira di Argonne andrà in pensione dopo anni di abilitazione alla scienza rivoluzionaria, HPCwire (20 dicembre 2019), <https://www.hpcwire.com/2019/12/20/argones-mira-supercomputer-to-retire-after-years-of-abilization-innovative-science/>. Se ancora operativi, questi sistemi si classificherebbero rispettivamente al 19° e 29° posto più veloci al mondo. Confronta novembre 2020, supra nota 70, con TOP500 List - June 2019, TOP500 (giugno 2019), <https://www.top500.org/lists/top500/list/2019/06/>.

79 Si veda, ad esempio, Kim Zetter, Top Federal Lab Hacked in Spear-Phishing Attack, Wired (20 aprile 2011), <https://www.wired.com/2011/04/oak-ridge-lab-hack/>; Natasha Bertrand ed Eric Wolff, Agenzia per le armi nucleari violata durante un massiccio attacco informatico, Politico (17 dicembre 2020), <https://www.politico.com/news/2020/12/17/nuclear-agency-hacked-officials-inform-congress-447855> (ultima visita 2 marzo 2021); Ryan Lucas, l'elenco delle agenzie federali interessate da un grave attacco informatico continua a crescere, NPR (18 dicembre 2020), <https://www.npr.org/2020/12/18/948133260/list-of-federal-agencies-affected-by-a-major-cyberattack-continues-to-grow> (ultima visita 2 marzo 2021).

80 Discutiamo i modelli di accesso ai dati nel terzo capitolo.

81 Cfr. Progetti in corso, RIKEN Ctr. per Computational Sci., <https://www.r-ccs.riken.jp/en/fugaku/research/covid-19/projects/>.

82 Fugaku conserva il titolo di supercomputer più veloce del mondo, HPCWire (17 novembre 2020), <https://www.hpcwire.com/off-the-wire/fugaku-retains-title-as-worlds-fastest-supercomputer/>.

83 novembre 2020, supra nota 70.

84 id.

85 Dietro le quinte di Fugaku come supercomputer più veloce del mondo, Fujitsu (2 febbraio 2021), <https://blog.global.fujitsu.com/fgb/2021-02-02/behind-the-scenes-of-fugaku-come-la-produzione-di-supercomputer-1-più-veloce-al-mondo/>.

86 id.

87 Don Clark, il supercomputer giapponese è stato incoronato il più veloce del mondo, NY Times (22 giugno 2020), <https://www.nytimes.com/2020/06/22/technology/japanese-supercomputer-fugaku-tops-american-chinese-macchine.html>.

88 Justin McCurry, Non-Woven Masks Better to Stop Covid-19, afferma Japanese Supercomputer, The Guardian (26 agosto 2020), <http://www.theguardian.com/world/2020/aug/26/maschere-non-tessute-meglio-fermare-covid-19-dice-il-supercomputer-giapponese>.

89 Fujitsu e RIKEN completano lo sviluppo congiunto del giapponese Fugaku, il supercomputer più veloce del mondo, Fujitsu (9 marzo 2021), <https://www.fujitsu.com/global/about/resources/news/press-releases/2021/0309-02.html>.

90 id.

91 Si veda, ad esempio, Rolf Harms & Michael Yamartino, The Economics of the Cloud (2010); Srijith Rajamohan e Robert E. Settlege, Informing the On/Off-Prem Cloud Discussion in Higher Education, 2020 Practice & Experience in Advanced Res. Informatica 64 (2020); Byung Chul Tak et al., Spostarsi o non spostarsi: l'economia del cloud computing, 3 USENIX Conf. sui temi caldi nel cloud computing 1 (2011); Edward Walker, Walter Briskin e Jonathan Romney, Affittare o non affittare dalle nuvole di archiviazione, 43 Computer 44 (2010).

92 Si veda, ad esempio, Di Zhang et al., RLScheduler: An Automated HPC Batch Job Scheduler Using Reinforcement Learning, Cornell U. (2 settembre 2020), <https://arxiv.org/pdf/1910.08925.pdf>.

93 Ad esempio, non siamo stati in grado di identificare buone stime dei costi di elettricità e raffreddamento per i supercomputer DOE.

94 Hugh Couchman et al., Compute Canada — Calcul Canada: A Proposal to the Canada Foundation for Innovation — National Platforms Fund

58 (2006).

95 Informazioni, Compute Canada, <https://www.computeCanada.ca/about/>.

96 Sistemi nazionali, Compute Canada, <https://www.computeCanada.ca/techrenewal/national-systems/>.

97 Compute Canada Technology Briefing, Compute Canada (novembre 2017), <https://www.computeCanada.ca/wp-content/uploads/2015/02/Technology-Briefing-November-2017.pdf>.

98 Cloud Computing for Researchers, Compute Canada (dicembre 2016), <https://www.computeCanada.ca/wp-content/uploads/2015/02/CloudStrategy2016-2019-forresearchersEXTERNAL-1.pdf>.

99 id.

100 Presentazione del budget 2018, Compute Canada (2018), <https://www.computeCanada.ca/wp-content/uploads/2015/02/UTF-8Compute20Canada20Budg-et20Submission202018.pdf> at 5.

101 Compute Canada ha previsto di aver soddisfatto solo il 55% circa della domanda totale di ore di calcolo della CPU nel 2018. Id.

102 id.

103 Compute Canada, Rapporto annuale 2019-2020 4 (2020).

104 Rapid Access Service, Compute Canada, <https://www.computeCanada.ca/research-portal/accessing-resources/rapid-access-service/>.

105 id.

106 Concorsi per l'allocazione delle risorse, supra nota 21.

107 id.

108 id.

109 id.

110 id.

111 2021 Resource Allocations Competition Results, supra nota 21.

capitolo 3

1 Invito all'azione del cloud di ricerca nazionale, Stan. U.Inst. per l'intelligenza artificiale centrata sull'uomo (2020), <https://hai.stanford.edu/national-research-cloud-joint-letter>.

2 Discutiamo la legge sulla privacy e le considerazioni sulla privacy in modo più dettagliato nel capitolo cinque.

3 Amy O'Hara & Carla Medalia, Condivisione dei dati nel sistema statistico federale: impedimenti e possibilità, 675 Annali Am. Acad. pol. & Soc. Sci. 138, 140-41 (2018); vedi anche President's Mgmt. Ordine del giorno, Piano d'azione 2020 della Strategia federale in materia di dati (2020).

4 Un migliore accesso ai dati, come descritto di seguito, promuoverebbe anche l'elaborazione di politiche basate sull'evidenza e migliorerebbe la fiducia nella scienza (poiché l'accesso ai dati rende molto più facili gli sforzi di replica).

5 Si veda, ad esempio, Nick Hart e Nancy Potok, Modernizing US Data Infrastructure: Design Considerations for Implementing a National Secure Data Service to Improve Statistics and Evidence Building (2020).

6 Queste iniziative hanno successo in quanto sono sostenibili e sono state utilizzate dai ricercatori per accedere ai dati governativi multi-agenzia. L'unica eccezione è il National Secure Data Service (NSDS), che non è stato ancora implementato. Di seguito discutiamo dell'NSDS insieme al Census Bureau e all'Evidence-Based Policy-Making Act del 2018. È importante sottolineare che il nostro obiettivo in questi casi di studio non è valutare i loro sforzi o misurare i loro esatti livelli di successo, ma identificare e comprendere alcune delle differenze e delle somiglianze nella gamma degli sforzi di condivisione dei dati.

7 Ad esempio, i dati del settore privato possono facilitare la ricerca sull'uso dei social media, sul comportamento di Internet o colmare le lacune per la ricerca statistica federale attraverso l'analisi dei big data. Vedere Robert M. Groves e Brian A. Harris-Kojetin, Using Private-Sector Data for Federal Statistics, Nat'l Ctr. per informazioni sulle biotecnologie. (12 gennaio 2017), <https://www.ncbi.nlm.nih.gov/books/NBK425876/>.

8 Cfr., ad esempio, National Data Service, Nat'l Data Serv., <http://www.nationaldataservice.org>; Open Science Data Cloud, Open Sci. Data Cloud, <https://www.opensciencedatacloud.org>, Harvard Dataverse, Harv. Dataverse, <https://dataverse.harvard.edu>, FigShare, <https://figshare.com>.

9 Facebook Data for Research fornisce l'accesso a una varietà di biblioteche, tramite piattaforme interne. Si veda, ad esempio, Facebook Data For Good, Facebook (2020), <https://dataforgood.fb.com/>; Che cos'è la Libreria inserzioni di Facebook e come posso cercarla?, Facebook (2021), <https://www.facebook.com/help/259468828226154>; Metodologia delle mappe dei disastri di Facebook, Facebook (15 maggio 2019), <https://research.fb.com/facebook-disaster-maps-methodology/>.

10 Ad esempio, Twitter ha un portale per sviluppatori che fornisce l'accesso alla loro API per consentire ai ricercatori di utilizzare i dati degli utenti per scopi non commerciali.

Vedi Twitter Developers, Twitter (2021), <https://developer.twitter.com/en/portal/petition/academic/is-it-right-for-you>; Porta avanti la tua ricerca con i dati di Twitter, Twitter (2021), <https://developer.twitter.com/en/solutions/academic-research>. Pertanto, il caricamento dei dati di Twitter su un cloud separato può fornire pochi incentivi ai ricercatori che possono utilizzare il percorso API.

11 Vedi Nat'l Acad. of Sci., Innovazioni nella statistica federale 31-42 (2017).

12 Cfr. Jennifer M. Urban, Joe Karaganis e Brianna M. Schofield, Notice & Takedown in Everyday Practice 39 (2017) (che illustra la difficoltà che i fornitori di servizi online incontrano nel valutare manualmente un grande volume di dati per una potenziale violazione; ad esempio, uno il fornitore di servizi online ha spiegato che "per paura di non riuscire a rimuovere il materiale in violazione e motivato dalla minaccia di danni legali, il suo personale impiegherà" sei passaggi per cercare di trovare il [contenuto identificato]."); vedi anche Lettera di Thom Tillis, Marsha Blackburn, Christopher A. Coons, Dianne Feinstein et. al, a Sundar Pichai, amministratore delegato, Google Inc. (3 settembre 2019), <https://www.ipwatchdog.com/wp-content/uploads/2019/09/9.3-Content-ID-Ltr.pdf> ("Abbiamo sentito da titolari di copyright a cui è stato negato l'accesso agli strumenti di Content ID e, di conseguenza, si trovano in notevole svantaggio nell'impedire il caricamento ripetuto di contenuti che hanno precedentemente identificato come in violazione. A loro resta la scelta di spendere ore ogni settimana cercando e inviando avvisi sulle stesse opere protette da copyright, o permettendo che la loro proprietà intellettuale venga sottratta.").

13 Per illustrare i costi dell'implementazione di Content ID su una piattaforma su larga scala, Google ha annunciato in un rapporto del 2016 che YouTube aveva investito più di 60 milioni di dollari in Content ID. Vedi Google, How Google Fights Piracy 6 (2016).

14 Cfr., ad esempio, Contratto con il cliente AWS, Amazon (30 novembre 2020), <https://aws.amazon.com/agreement/>.

15 Ad esempio, nelle 29 agenzie distinte del Dipartimento della salute e dei servizi umani (HHS), i dati "sono in gran parte conservati in silos con una mancanza di consapevolezza organizzativa di quali dati vengono raccolti all'interno del Dipartimento e di come richiedere l'accesso. Ogni agenzia opera all'interno della propria autorità statutaria e ogni set di dati può essere disciplinato da un particolare insieme di regolamenti. Dipartimento della salute e dei servizi umani degli Stati Uniti, The State of Data Sharing presso il Dipartimento della salute e dei servizi umani degli Stati Uniti 4 (2018).

16 Cfr., ad es., id. a 8 ("L'HHS non dispone di processi coerenti e standardizzati affinché un'agenzia richieda dati a un'altra agenzia.").

17 O'Hara & Medalia, supra nota 3, pp. 140-41.

18 Cfr. id. a 142 ("La maggior parte degli accordi [di condivisione dei dati] si basa in gran parte su relazioni interpersonali e accordi informali di quid pro quo, gestendo le richieste di dati in modo meno centralizzato.").

19 Jeffrey Mervis, Come due economisti hanno avuto accesso diretto ai documenti fiscali dell'IRS, Sci. Mag. (22 maggio 2014), <https://www.sciencemag.org/news/2014/05/how-two-economists-got-direct-access-irs-tax-records>.

20 Cfr. Robert M. Groves e Adam Neufeld, Accelerating the Sharing of Data Across Sectors to Advance the Common Good 17 (2017).

21 Cfr., ad esempio, Data Use Agreement, Dept Health & Human Services, [https://www.hhs.gov/sites/default/files/ocio/eplc/EPLC%20Archive%20 Documents/55-Data%20Use%20Agreement%20%28DUA%29/eplc_dua_practices_guide.pdf](https://www.hhs.gov/sites/default/files/ocio/eplc/EPLC%20Archive%20Documents/55-Data%20Use%20Agreement%20%28DUA%29/eplc_dua_practices_guide.pdf).

22 O'Hara & Medalia, supra nota 3, at 138, 141.

23 Bipartisan Pol'y Ctr., Barriere all'utilizzo dei dati governativi: analisi estesa della commissione statunitense sull'indagine sulle agenzie e gli uffici federali basata su prove della Commissione degli Stati Uniti 18-20 (2018).

24 Cfr. Research Data Assistance Center (ResDAC), Ctr. per i servizi Medicare e Medicaid (30 agosto 2018), <https://www.cms.gov/Research-Statistics-Data-and-Systems/Research/ResearchGenInfo/ResearchDataAssistanceCenter>.

25 O'Hara & Medalia, supra nota 3, a 141.

26 Michelle Mello et al., In attesa di dati: ostacoli all'esecuzione di accordi sull'utilizzo dei dati, 367 Sci. Mag. 150 (10 gennaio 2020), https://www.sciencemagazine.org/sciencemagazine/10_january_2020/MobilePagedArticle.action?articleId=1552284#articleId1552284.

27 Intervista con Amy O'Hara, direttore esecutivo, Georgetown Federal Statistical Research Data Center (22 aprile 2021); vedi anche Special Sworn Research Program, Bureau of Econ. Analisi, <https://www.bea.gov/research/special-sworn-researcher-program>; Nat'l Ctr. per Educ. Stat., Manuale delle procedure per i dati a uso limitato (2011).

28 Si veda US Gov't Accountability Office, Federal Agencies Need to Affronting Aging Legacy Systems (2016).

29 O'Hara & Medalia, supra nota 3, pp. 140-41.

30 Cfr., ad es., id.; US Gov't Accountability Office, supra nota 28.

31 Mgmt del Presidente. Ordine del giorno, supra nota 3, ore 11.

32 Groves & Neufeld, supra nota 20, 12-13. Per una definizione precisa di dati sensibili, vedere Glossario: informazioni sensibili, Nat'l Inst. Standard e tecnologia, https://csrc.nist.gov/glossary/term/sensitive_information.

33 Shanna Nasiri, FedRAMP Low, Moderate, High: Understanding Security Baseline Levels, Reciprocity (24 settembre 2019), <https://reciprocity.com/fedramp-low-moderate-high-understanding-security-baseline-levels/>.

34 Michael McLaughlin, Reforming FedRAMP: una guida per migliorare l'approvvigionamento federale e la gestione del rischio dei servizi cloud, info. Tecnico. & Innovazione trovata. (15 giugno 2020), <https://itif.org/publications/2020/06/15/reforming-fedramp-guide-improving-federal-procurement-and-risk-management>.

35 Domande frequenti, FedRAMP, <https://www.fedramp.gov/faqs>.

36 Fai una volta, usane molti - Come le agenzie possono riutilizzare un'autorizzazione FedRAMP, FedRAMP (7 maggio 2020), <https://www.fedramp.gov/how-agencies-can-reuse-a-fedramp-authorization/>.

37 Linee di base dei controlli di sicurezza FedRAMP, FedRAMP bassa, moderata e alta (2021).

38 Controlli di sicurezza e privacy per i sistemi informativi e le organizzazioni, Nat'l Inst. Standard e tecnologia. (23 settembre 2020), <https://csrc.nist.gov/publications/detail/sp/800-53/rev-5/final>.

39 Cfr., ad es., id.; Framework di gestione del rischio NIST AC-2: Gestione account, Nat'l Inst. Standard e tecnologia, <https://csrc.nist.gov/Projects/risk-management/sp800-53-controls/release-search#!/control?version=4.0&number=AC-2>; NIST Risk Management Framework AC-3: Applicazione dell'accesso, Nat'l Inst. Standard e tecnologia, <https://csrc.nist.gov/Projects/risk-management/sp800-53-controls/release-search#!/control?version=5.1&number=AC-3>.

40 Cfr. Mark Bergen, Google Engineers Refused to Build Security Tool to Win Military Contracts, Bloomberg (21 giugno 2018), <https://www.bloomberg.com/news/articles/2018-06-21/google-engineers-refused-per-costruire-strumenti-di-sicurezza-per-vincere-contratti-militari>.

41 Vedi Nat'l Inst. Standards & Tech., Standard per la classificazione della sicurezza delle informazioni federali e dei sistemi informativi (2004).

42 Partnership con FedRAMP, FedRAMP, <https://www.fedramp.gov/cloud-service-providers/>. Sebbene possa costare ai fornitori di servizi cloud tra \$ 365.000 e \$ 865.000 e impiegare 6-12 mesi per ricevere la conformità FedRAMP, Adam Isles, Securing Your Cloud Solutions: Research and Analysis on Meeting FedRAMP/Government Standards 21 (2017), tali costi sono a carico dei gli stessi fornitori di servizi cloud, non i clienti dei fornitori. In effetti, FedRAMP utilizza un modello "fai una volta, usa molti": una volta che un fornitore di servizi cloud ottiene un'autorizzazione a operare (ATO), quell'ATO può essere sfruttato e ricambiato tra più clienti, eliminando gli sforzi duplicati e le incoerenze che deriverebbero dalla richiesta di più ri-autorizzazioni.

Id. alle 11.

43 Anche all'interno di FedRAMP vi sono notevoli variazioni nel modo in cui le diverse organizzazioni assicurano la conformità ai controlli e agli standard pertinenti, con molti dei controlli scritti in modo sufficientemente ampio da dare spazio a interpretazioni sostanziali. Tuttavia, presenta una serie di considerazioni e requisiti che sono coerenti tra i domini e consente un grado di prevedibilità e affidabilità che non è presente in altri aspetti della governance federale dei dati.

44 O'Hara & Medalia, supra nota 3, at 141 ("La condivisione dei dati avviene su base obbligatoria o volontaria, e le richieste di dati sono gestite tramite personale/processo designato o diffusamente tramite un'organizzazione.").

45 Bipartisan Pol'y Ctr., supra nota 23, at 17 ("La mancanza di procedure o linee guida standard per la condivisione dei dati tra le agenzie federali che finanziano la ricerca rende gli sforzi per collegare e condividere i dati difficili o inefficienti.").

46 Si veda, ad esempio, Amy O'Hara, US Federal Data Policy: An Update on The Federal Data Strategy and The Evidence Act, 5 Int'l J. Population Data Sci. 5 (2020).

47 Sebbene gli sforzi e le iniziative federali esistenti siano già finalizzati all'armonizzazione delle migliori pratiche di condivisione dei dati, cfr. l'INRC può accelerare questi sforzi. In effetti, lo sviluppo di standard chiari e coerenti è fondamentale per facilitare la condivisione dei dati. David Crotty, Ida Sim e Michael Stebbins, Accesso aperto ai dati di ricerca finanziati dal governo federale 7 (2020).

48 Questi requisiti sono incoerenti e obsoleti a causa delle difficoltà nella definizione del rischio e dell'avversione al rischio da parte delle agenzie. Vedi O'Hara e Medalia, supra nota 3, pp. 140-41; si veda anche David S. Johnson et al., The Opportunities and Challenges of Using Administrative Data Linkages to Evaluate Mobility, 657 Annals Am. Acad. pol. & Soc. Sci. 252-53 (2015).

49 Per una discussione sulle minacce di inferenza, vedi Nat'l Acad. of Sci., Eng'g & Med., Federal Statistics, Multiple Data Sources, and Privacy Protection: Next Steps 68 (2017).

50 Congzheng Song & Ananth Raghunathan, Information Leakage in Embedding Models, Cornell U. (31 marzo 2020), <https://arxiv.org/abs/2004.00053>.

51 Cfr., ad esempio, Statistical Safeguards, Census Bureau (1 luglio 2021), https://www.census.gov/about/policies/privacy/statistical_safeguards.html.

- 52 Alexandra Wood et al., *Differential Privacy: A Primer for a Non-Technical Audience*, 21 Vand. J. Ent. e tecnologia. L. 209 (2018).
- 53 Regolare l'accesso ai dati, UK Data Serv., <https://www.ukdataservice.ac.uk/manage-data/legal-ethical/access-control/five-safes>.
- 54 Centro di ricerca sui dati amministrativi, Coleridge Initiative, <https://coleridgeinitiative.org/adrf/>.
- 55 Cfr. O'Hara, *supra* nota 46.
- 56 Per ulteriori discussioni sulle implicazioni della NRC sulla privacy, vedere il capitolo cinque.
- 57 Si veda l'Office of Mgmt degli Stati Uniti. & Budget, Ostacoli all'utilizzo dei dati amministrativi per la creazione di prove 7 (2016).
- 58 Strumento di ricerca sui dati amministrativi, *supra* nota 54.
- 59 *id.*
- 60 Formazione, Iniziativa Coleridge, <https://coleridgeinitiative.org/training/>.
- 61 Guida per l'utente di ADRF: Data Explorer, Coleridge Initiative, <https://coleridgeinitiative.org/adrf/documentation/using-the-adrf/data-explorer/>.
- 62 Guida per l'utente di ADRF: Esportazione dei risultati, Coleridge Initiative, <https://coleridgeinitiative.org/adrf/documentation/using-the-adrf/exporting-results/>.
- 63 *id.*
- 64 Guida per l'utente di ADRF: Applicazione di hashing dei dati, Coleridge Initiative, <https://coleridgeinitiative.org/adrf/documentation/adrf-overview/data-hashing-application/>.
- 65 Guida per l'utente di ADRF: Modello di sicurezza e conformità, Coleridge Initiative, <https://coleridgeinitiative.org/adrf/documentation/adrf-overview/security-model-and-compliance/>.
- 66 Panoramica per i collaboratori, Coleridge Initiative, <https://coleridgeinitiative.org/collaborators/>.
- 67 Dati, Stan. Med. Ctr. per Population Health Sci., <https://med.stanford.edu/phs/data.html>.
- 68 Stanford Ctr. per Filantropia & Civ. Soc'y, *Trusted Data Intermediaries* 2-3 (2018).
- 69 Anche altri hanno riconosciuto il vantaggio dei modelli DUA universali. Vedi Mello et al., *supra* nota 26, a 150; Linee guida per la fornitura e l'utilizzo di dati amministrativi a fini statistici, Office of Mgmt. & Bilancio (14 febbraio 2014), <https://obamawhitehouse.archives.gov/sites/default/files/omb/memoranda/2014/m-14-06.pdf>.
- 70 Dati | Centro per le scienze della salute della popolazione | Stanford Medicina, Stan. Med. Ctr. per Population Health Sci., <https://med.stanford.edu/phs/data.html>.
- 71 Cfr. Stanford PHS – Datasets, Redivis, <https://redivis.com/StanfordPHS/datasets?orgDatasets-tags=109.medicare>.
- 72 Livelli di accesso, Redivis (luglio 2020), <https://docs.redivis.com/reference/data-access/access-levels>.
- 73 Fase 1: ottenere l'accesso, Stan. Med. Documentazione PHS, <https://phsdocs.developerhub.io/start-here/getting-data-access>.
- 74 *id.*
- 75 *id.*
- 76 Flusso di lavoro per l'utilizzo dei dati PHS, Stan. Med. Documentazione PHS, <https://phsdocs.stanford.edu/start-here/phs-data-use-workflow>.
- 77 *id.*
- 78 Ambiente informatico PHS, Stan. Med. Documentazione PHS, <https://phsdocs.stanford.edu/computing-environment>.
- 79 *id.*
- 80 Si veda US Gov't Accountability Office, *Federal Agencies Need to Address Aging Legacy Systems* 15 (2016) (osservando che dal 2010 al 2015 molte agenzie federali hanno aumentato la spesa per le operazioni e la manutenzione a causa dei sistemi legacy).
- 81 David Freeman Engstrom, Daniel E. Ho, Catherine M. Sharkey e Mariano-Florentino Cuéllar, *Government by Algorithm: Artificial Intelligence in Federal Administrative Agencies* 6, 71-72 (2020).
- 82 *id.* alle 6-7.
- 83 *id.* a 71-72.
- 84 *id.* a 73.
- 85 *id.* alle 6.
- 86 Cfr. Risultati per l'America, *The Promise of the Foundations for Evidence-Based Policymaking Act and Proposed Next Steps* (2019).
- 87 Ad esempio, l'Uniform Federal Crime Reporting Act del 1988 richiede alle forze dell'ordine federali di condividere i dati sui reati con l'FBI. Vedere 34 USC §§41303(c)(2), (3), (4). Sfortunatamente, però, nessuna agenzia federale apparentemente condivide attualmente i propri dati con l'FBI ai sensi di questa legge. Nat'l Acad. di Sci., *supra* nota 11, p. 41 (2017).
- 88 Katharine G. Abraham e Ron Haskins, *La promessa di un processo decisionale basato sull'evidenza; Comm'n on Evidence-Based Policymaking* (2018).
- 89 Questi meccanismi di tutela della privacy sono particolarmente importanti alla luce delle continue sfide legali e politiche nell'applicazione differenziata della privacy ai dati federali. Si veda, ad esempio, Dan Bouk e danah boyd, *Democracy's Data Infrastructure* (2021).
- 90 Foundations for Evidence-Based Policymaking Act del 2018, Pub. L. n. 115-435.
- 91 Legge sulla protezione delle informazioni riservate e sull'efficienza statistica del 2002, Pub. L. n. 107-347.
- 92 Panoramica, *Federal Data Strategy* (2020), <https://strategy.data.gov/overview/>.
- 93 UK Data Service, UK Data Serv., <https://www.ukdataservice.ac.uk/> (ultima visita 21 giugno 2021).
- 94 Ad esempio, la sola amministrazione della previdenza sociale dispone di oltre 14 petabyte di dati, archiviati in circa 200 database. Engstrom, Ho, Sharkey & Cuéllar, *supra* nota 81, a 72.
- 95 Google Earth Engine, Google Earth Engine, <https://earthengine.google.com> (ultima visita 15 agosto 2021).
- 96 Mondo del lavoro, ADR UK, <https://www.adruk.org/our-work/world-of-work/>.
- 97 Database annuale degli intervistati, 1973-2008: Secure Access, UK Data Serv. (2020), <https://beta.ukdataservice.ac.uk/datacatalogue/studies/study?id=6644>.
- 98 Sondaggio sull'innovazione nel Regno Unito, servizio dati del Regno Unito. (2021), <https://beta.ukdataservice.ac.uk/datacatalogue/studies/study?id=6699>.
- 99 Indagine trimestrale sulla forza lavoro, 1992-2021: Secure Access, UK Data Serv. (2021), <https://beta.ukdataservice.ac.uk/datacatalogue/studies/study?id=6727>.
- 100 Understanding Society: Waves 1-10, 2009-2019 e Harmonized BHPS: Waves 1-18, 1991-2009: Secure Access, UK Data Serv. (2021), <https://beta.ukdataservice.ac.uk/datacatalogue/studies/study?id=6676>.
- 101 Questi set di dati hanno aiutato i ricercatori ad affrontare alcune domande specifiche di interesse pubblico. Si veda, ad esempio, Francisco Perales, *Why Does the Work Women Do Pay Less Than the Work Men Do?*, UK Data Serv. (8 dicembre 2011), <https://beta.ukdataservice.ac.uk/impact/case-studies/case-study?id=62>; Eva-Maria Bonin, *I programmi per genitori riducono i disturbi della condotta?*, UK Data Serv. (4 aprile 2012), <https://beta.ukdataservice.ac.uk/impact/case-studies/case-study?id=93>.
- 102 Identificazione dell'accesso prioritario o dei miglioramenti della qualità per i dati e i modelli federali per la ricerca e lo sviluppo (R&S) e i test sull'intelligenza artificiale; Richiesta di informazioni, 84 Fed. Reg. 32962 (10 luglio 2019).

103 Nick Hart, Data Coalition Comments on AI Data and Model R&D RFI, Data Coalition (9 agosto 2019), http://www.datacoalition.org/wp-content/uploads/2019/09/Comment.RFI_OMB_2019-14618.DataCoalition.pdf.
104 id.

105 Cfr. Adam R. Pah et al., How to Build a More Open Justice System, 369 Sci. 134 (2020); si veda anche Seamus Hughes, The Federal Courts Are Running an Online Scam, Politico (20 marzo 2019), <https://www.politico.com/magazine/story/2019/03/20/pacer-court-records-225821> .

106 Legal Authority and Policies for Data Linkage at Census, Census Bureau (4 aprile 2018), <https://www.census.gov/about/adrm/linkage/about/authority.html> .

107 BLS Accesso limitato ai dati, US Bureau of Lab. Stat., <https://www.bls.gov/rda/restricted-data.htm> (ultimo aggiornamento 20 maggio 2021).

108 Benvenuti al PDS, NASA, <https://pds.nasa.gov>.

capitolo 4

1 Sebbene riteniamo che questi siano gli assi primari da prendere in considerazione, alcune considerazioni secondarie includono il peso organizzativo, la fidelizzazione dei talenti e le spese burocratiche.

2 Servizio di ricerca del Congresso, Centri di ricerca e sviluppo finanziati a livello federale (FFRDC): background e problemi per il Congresso 1 (2020).

3 id. Vedi anche Informazioni su IDA, Inst. Defence Analyses, <https://www.ida.org/about-ida> (sottolineando che IDA, il subappaltatore del settore privato che gestisce il Science & Technology Policy Institute e molti altri FFRDC, "gode di un accesso insolito a informazioni governative classificate e informazioni aziendali sensibili informazioni proprietarie."); Ufficio per la responsabilità del governo degli Stati Uniti, centri di ricerca e sviluppo finanziati a livello federale: miglioramento della supervisione e della valutazione necessari per il programma pilota di accesso ai dati del DOD 6 (2020) (discutendo di come il Dipartimento della Difesa sia stato in grado di stabilire un programma pilota triennale che ha consentito al suo FFRDC ricercatori di rinunciare a dover ottenere accordi di non divulgazione con ciascun proprietario dei dati al fine di semplificare il processo di accesso ai dati).

4 Nick Hart e Nancy Potok, Modernizzazione dell'infrastruttura dati degli Stati Uniti: considerazioni sulla progettazione per l'implementazione di un servizio dati sicuro nazionale per migliorare le statistiche e la creazione di prove (2020).

5 id. alle 26.

6 id. alle 26-27, 29-30.

7 US Gov't Accountability Office, supra nota 3, punto 6. Si noti che mentre l'FFRDC deve operare per servire i suoi sponsor, nello stabilire un FFRDC, lo sponsor deve garantire che operi con sostanziale indipendenza; l'FFRDC deve essere "gestito, gestito o amministrato da un'organizzazione autonoma o come un'unità operativa distinta e identificabile di un'organizzazione madre". Vedere Federal Acquisition Regulations [di seguito "FAR"] § 35.017(a)(2).

8 Un esempio di ciò è il Science & Technology Policy Institute, di cui parleremo in un caso di studio di seguito.

9 US Dep't of Energy, The State of the DOE National Laboratories 11-13 (2020).

10 Si veda, ad esempio, Altre agenzie federali si dirigono verso il cloud con Azure Government, Applied Info. Sci. (23 febbraio 2018), <https://www.ais.com/more-federal-agency-head-to-the-cloud-with-azure-government/>; vedere anche AWS GovCloud, Amazon, <https://aws.amazon.com/govcloud-us/>. Microsoft aveva anche precedentemente ottenuto un contratto da 10 miliardi di dollari dal Pentagono. Vedi Kate Conger, Microsoft Wins Pentagon's \$ 10 Billion's JEDI Contract, Contrastare Amazon, NY Times (4 settembre 2020), <https://www.nytimes.com/2019/10/25/technology/dod-jedi-contract.html> . Tuttavia, questo contratto è stato recentemente annullato "a causa dell'evoluzione dei requisiti, della maggiore tutela del cloud e dei progressi del settore". Ellie Kaufman e Zachary Cohen, il Pentagono annulla un contratto cloud da 10 miliardi di dollari dato a Microsoft su Amazon, CNN (6 luglio 2021), <https://www.cnn.com/2021/07/06/tech/defense-department-cancels-jedi-contract-amazon-microsoft/index.html>. Il Pentagono ora cercherà invece nuove offerte per un contratto JWCC (Joint Warfighting Cloud Capability) aggiornato da Amazon e Microsoft. Id.

11 Si veda, ad esempio, Bram Bout, Helping Universities Build What's Next with Google Cloud Platform, Google (25 ottobre 2016), <https://blog.google/outreach-initiatives/education/helping-universities-build-whats-prossima-piattaforma-cloud-google/>; Cloud Computing per l'istruzione, Amazon, <https://aws.amazon.com/education/> .

12 Congressional Research Serv., supra nota 2, 11-12 (2020).

13 US Dep't of Energy, Annual Report on the State of the DOE National Laboratories 87 (2017).

14 Congressional Research Serv., supra nota 2, a 19.

15 Servizio di ricerca del Congresso, Office of Science and Technology Policy (OSTP): storia e panoramica 9 (2020). I doveri di STPI sono anche specificati in 42 USC § 6686.

16 Cosa sono gli FFRDC?, inst. Analisi della difesa, <https://www.ida.org/ida-ffrdcs>.

17 id.

18 sponsor, ist. Analisi della difesa, <https://www.ida.org/en/about-ida/sponsors>.

19 id.

20 Congressional Research Serv., supra nota 15, at 9-10.

21 Ad esempio, dal 2008 al 2012, queste altre agenzie federali hanno contribuito con un totale di 9,8 milioni di dollari di finanziamento a STPI mentre NSF ha contribuito con circa 24 milioni di dollari. Ufficio per la responsabilità del governo degli Stati Uniti, centri di ricerca finanziati dal governo federale: revisioni dell'agenzia sulla remunerazione dei dipendenti e sulle prestazioni del centro 43-44 (2014).

22 Congressional Research Serv., supra nota 15, at 9-10.

23 id.

24 42 USC § 6686 (d).

25 42 U.S.C. § 6686(e).

26 Sci. e tecnologia. Pol'y Inst., Rapporto al Presidente Anno fiscale 2020 (2020).

27 Cfr., ad esempio, Open Government, Millennium Challenge Corp., <https://www.mcc.gov/initiatives/initiative/open>; Nat'l Geospatial Advisory Comm., Advancing the National Spatial Data Infrastructure Through Public-Private Partnerships and Other Innovative Partnerships (2020); Nat'l Aeronautics & Space Admin., Partenariati pubblico-privato per lo sviluppo delle capacità spaziali 33-36 (2014).

28 Big Data Value Public-Private Partnership, European Comm'n (9 marzo 2021), <https://digital-strategy.ec.europa.eu/en/library/big-data-value-public-private-associazione>.

29 RAND, Partenariati pubblico-privati per la condivisione dei dati: un ambiente dinamico 33, 99 (2000).

30 Cfr. Homepage - Alberta Data Partnerships, Alberta Data Partnerships, <http://abdatapartnerships.ca> (ultima visita 15 agosto 2021).

31 Alberta Data Partnerships, Una storia di successo P3 1 (2017).

32 id.

33 id. alle 19, 35.

34 id. alle 15.

35 id.

36 id. alle 1.

37 Nat'l Geospatial Advisory Comm., Caso d'uso del partenariato pubblico-privato: Alberta Data Partnerships 1 (2020).

38 Alberta Data Partnerships, supra nota 31, a 15.

39 id. alle 16.

40 The COVID-19 High Performance Computing Consortium, COVID-19 HPC Consortium, <https://covid19-hpc-consortium.org>.

41 id.

42 Si veda, ad esempio, David Hall, Why Public-Private Partnerships Don't Work (2015); Svantaggi e insidie dell'opzione PPP, APMG Int'l, <https://ppp-certification.com/ppp-certification-guide/54-disadvantages-and-pitfalls-ppp-option>.

43 Graeme A. Hodge, Carsten Greve e Anthony E. Boardman, Manuale internazionale sui partenariati pubblico-privato, 187-90 (2012).

44 Ad esempio, a un'estremità di uno spettro, il California Teale Data Center crea, possiede, mantiene e archivia i propri set di dati per uso del settore privato. Al contrario, il Pennsylvania Spatial Data Access ospita i metadati, richiedendo agli utenti di chiedere l'accesso alle origini dati effettive. RAND, supra nota 29, pp. 102-03. Incoraggiamo la Task Force a esaminare questo rapporto completo per valutare le varie opzioni organizzative per un modello di stanza di compensazione dei dati PPP.

45 Angela Ballantyne & Cameron Stewart, Big Data e partenariati pubblico-privati nella sanità e nella ricerca, 11 Asian Bioethics R. 315, 315 (2019).

46 Cfr. Gov't Accountability Office, Human Capital: Improving Federal Recruiting and Hiring Sforts; si veda anche Catch and Retain: Improving Recruiting and Retention at Government Agencies, Salesforce, <https://www.salesforce.com/solutions/industries/government/resources/government-recruitment-software/>.

47 Partenariato per il servizio pubblico, Indagine sul futuro del servizio pubblico 2 (2020).

48 id.

Capitolo 5

1 Invito all'azione del cloud di ricerca nazionale, Stan. U.Inst. per l'intelligenza artificiale centrata sull'uomo (2020), <https://hai.stanford.edu/national-research-cloud-joint-letter>.

2 Le informazioni sensibili, come definite dal National Institute of Standards and Technology, sono informazioni in cui la perdita, l'uso improprio o l'accesso non autorizzato o la modifica potrebbero pregiudicare l'interesse nazionale o lo svolgimento dei programmi federali o la privacy a cui gli individui hanno diritto ai sensi 5 USC § 552a (legge sulla privacy); che non è stato specificamente autorizzato in base a criteri stabiliti da un ordine esecutivo o da un atto del Congresso da mantenere classificato nell'interesse della difesa nazionale o della politica estera. Vedi Glossario: Informazioni sensibili, Nat'l Inst. Standard e tecnologia, https://csrc.nist.gov/glossary/term/sensitive_information.

3 Ringraziamo Mark Krass per queste intuizioni.

4 Le agenzie contemplate dalla legge includono "qualsiasi dipartimento esecutivo, dipartimento militare, ente governativo, società controllata dal governo o altro istituto nel ramo esecutivo del governo [federale] (incluso l'ufficio esecutivo del presidente), o qualsiasi organismo di regolamentazione indipendente agenzia." 5 USC § 552(f)(1).

5 Ufficio di contabilità generale degli Stati Uniti, record linkage e privacy: problemi nella creazione di nuove informazioni statistiche e di ricerca federali 10 (2001).

6 Intervista a Marc Groman, ex consulente senior per la privacy, ufficio di gestione e bilancio della Casa Bianca (18 febbraio 2021); si veda anche Bipartisan Pol'y Ctr., Barriers to Using Government Data: Extended Analysis of the US Commission on Evidence-Based Policymaking's Survey of Federal Agencies and Offices 10 (2018).

7 Vedi Joseph Near e David Darais, Differentially Private Synthetic Data, Nat'l Inst. Standard e tecnologia. (3 maggio 2021), <https://www.nist.gov/blogs/cyber-security-insights/differentially-private-synthetic-data>; si veda anche Steven M. Bellovin et al., Privacy and Synthetic Datasets, 22 Stan. L. Rev. 1 (2019).

8 Legge sull'e-government del 2002, pub. L. n. 107-347.

9 Legge sulla protezione delle informazioni riservate e sull'efficienza statistica del 2002, 44 USC § 3501 (2012).

10 Foundations for Evidence-Based Policymaking Act del 2017, Pub. L. n. 115-435, 132 Stat. 5529 (2019).

11 Mgmt del Presidente. Ordine del giorno, Piano d'azione 2020 della Strategia federale in materia di dati (2020).

12 Legge sulla privacy del 1974, 5 USC § 552a (2012).

13 Esistono molte versioni dei Fair Information Practice Principles e il governo degli Stati Uniti non ne ha istituzionalizzato una versione specifica, sebbene si faccia comunemente riferimento alla versione utilizzata dal Department of Homeland Security (disponibile all'indirizzo: <https://www.dhs.gov/publication/privacy-policy-guida-memorandum-2008-01-fair-information-practice-principi>). L'Organizzazione per la cooperazione e lo sviluppo economico ne ha prodotto una versione influente nel 1980 (rivista nel 2013), che rimane una fonte autorevole. Linee guida dell'OCSE sulla protezione della privacy e sui flussi transfrontalieri di dati personali, OCSE (2013), <https://www.oecd.org/digital/economy/oecdguidelinesonthe protection of privacy and transborder flows of personal data>. htm.

14 Cfr. David Freeman Engstrom, Daniel E. Ho, Catherine M. Sharkey e Mariano-Florentino Cuéllar, Government by Algorithm: Artificial Intelligence in Federal Administrative Agencies (2020) (che documenta l'uso attuale dell'IA da parte delle agenzie governative).

15 5 USC § 552a (a) (5).

16 5 USC §§ 552a(a)(8)(A)(i)(I), (II).

17 Legge sulla privacy del 1974, Elec. Informazioni sulla privacy. Ctr., <https://epic.org/privacy/1974act/> (ultima visita 15 agosto 2021).

18 Cfr. scheda informativa: National Secure Data Service Act promuove la condivisione responsabile dei dati nel governo, Data Coalition (13 maggio 2021), <https://www.datacoalition.org/fact-sheet-national-secure-data-service-act-advances-responsible-data-sharing-in-government/>; Ufficio per la responsabilità del governo degli Stati Uniti, collegamento dei record e privacy: problemi nella creazione di nuove ricerche federali e informazioni statistiche (2001).

19 Non è un'ironia da poco che le società private negli Stati Uniti abbiano adempiuto a tale missione oggi. In effetti, il governo degli Stati Uniti ora si rivolge all'industria privata, tramite procedimenti legali o tramite appalti, quando richiede dati su individui che il governo stesso non raccoglie. Senatore Ron

Wyden ha proposto una legislazione per impedire al governo di effettuare questi acquisti. Wyden, Paul e membri bipartisan del Congresso presentano la legge sul quarto emendamento non è in vendita, Ron Wyden, senatore degli Stati Uniti per l'Or. (21 aprile 2021) , <https://www.wyden.senate.gov/news/press-releases/wyden-paul-and-bipartisan-members-of-congress-introduce-the-fourth-emendment-is-atto-non-in-vendita->.

20 Cfr., ad esempio, World Econ. Forum, La prossima generazione di condivisione dei dati nei servizi finanziari (2019).

21 Si veda, ad esempio, Stacie Dusetzina et al., Linking Data for Health Services Research: A Framework and Instructional Guide, Agency for Healthcare Research & Quality (1 settembre 2014), <https://www.ncbi.nlm.nih.gov/libri/NBK253315/>.

22 Cfr., ad esempio, European Comm'n, A European Strategy for Data (2020) (che sostiene l'aggregazione transfrontaliera dei dati e il collegamento dei dati del settore pubblico e privato); M Sanni Ali et al., Administrative Data Linkage in Brazil: Potentials for Health Technology Assessment, 10 Frontiers in Pharmacology 984 (2019); Collegamento dati, istituto australiano. of Health & Welfare (4 gennaio 2020), <https://www.aihw.gov.au/our-services/data-linkage>.

23 Si veda, ad esempio, Elsa Augustine, Vikash Reddy & Jesse Rothstein, Linking Administrative Data: Strategies and Methods (2018) (che descrive suggerimenti per condurre collegamenti di dati in California); si veda anche US Dep't of Health & Human Services, Status of State Efforts to Integrate Health and Human Services Systems and Data (2016).

24 Ben Moscovitch, How President Biden Can Improve Health Data Sharing For COVID-19 And Beyond, Health Affairs (1 marzo 2021), <https://www.healthaffairs.org/doi/10.1377/hblog20210223.611803/full/>.

25 Home, Johns Hopkins Coronavirus Resource Ctr., <https://coronavirus.jhu.edu/>.

26 Il progetto di monitoraggio COVID, <https://covidtracking.com/>.

27 Fred Bazzoli, COVID-19 Emergency Shows Limits of Nationwide Data Sharing Infrastructure, Healthcare IT News (2 giugno 2020), <https://www.healthcareitnews.com/news/covid-19-emergency-shows-limitations-nationwide-data-sharing-infrastructure>.

28 Cfr., ad esempio, C. Jason Wang et al., Response to COVID-19 in Taiwan: Big Data Analytics, New Technology, and Proactive Testing, JAMA (3 marzo 2020) , <https://jamanetwork.com/journals/jama/articolo-completo/2762689/>; Fang-Ming Chen et al., Big Data Integration and Analytics to Prevent a Potential Hospital Outbreak of COVID-19 in Taiwan, 54 J. Microbiology, Immunology & Infection 129-30 (2020).

29 Si veda, ad esempio, Q&A sul programma "Total Information Awareness" del Pentagono, em. CL Union, <https://www.aclu.org/other/qa-pentagons-total-information-awareness-program> ; I cinque problemi con CAPPS II: perché la proposta di profilazione dei passeggeri delle compagnie aeree dovrebbe essere abbandonata, em. CL Union, <https://www.aclu.org/other/five-problems-capps-ii>.

30 Cfr., ad esempio, Barton Gellman, Dark Mirror: Edward Snowden and the American Surveillance State (2020); Edward Snowden, Record permanente (2019).

31 5 USC § 552 (b).

32 5 USC § 552 (b) (3).

33 La legge sulla privacy contiene anche specifici ritagli per le comunicazioni al Census Bureau e alla National Archives and Records Administration.

Tuttavia, le esclusioni per queste due agenzie richiedono che le divulgazioni siano effettuate rispettivamente ai fini di un'indagine di censimento e di registrazione del valore storico. Poiché lo scopo esplicito dell'NRC è quello di democratizzare l'innovazione dell'IA, è improbabile che l'NRC possa trarre vantaggio da questa eccezione esistente alla divulgazione dei set di dati ai sensi del Privacy Act.

34 Ad esempio, l'elenco degli usi di routine dell'Agenzia federale per la gestione delle emergenze include un'ampia divulgazione "[a] un'agenzia o un'organizzazione allo scopo di eseguire operazioni di audit o supervisione come autorizzato dalla legge, ma solo le informazioni necessarie e pertinenti a tali funzione di revisione contabile o di supervisione". legge sulla privacy del 1974; Department of Homeland Security Federal Emergency Management Agency-008 Disaster Recovery Assistance Files System of Records, 78 Fed. Reg. 25282 (30 maggio 2013).

35 Si veda, ad esempio, Britt v. Naval Investigative Service, 886 F.2d 544 (3d Cir. 1989).

36 Il Privacy Act del 1974, supra nota 17.

37 5 USC § 552(b)(5).

38 5 USC §§ 552a(a)(8)(B)(i), (ii) (sottolineatura aggiunta).

39 44 USC § 3561(8), (12).

40 Dipartimento della salute e dei servizi umani degli Stati Uniti, Lo stato della condivisione dei dati presso il Dipartimento della salute e dei servizi umani degli Stati Uniti 16 (2018).

41 44 USC § 3575(4).

42 Cfr. Engstrom, Ho, Sharkey & Cuéllar, supra nota 14, p. 16 (trovando che il Bureau of Labor Statistics è una delle prime dieci agenzie che utilizzano l'intelligenza artificiale); Machine Learning, Census Bureau (17 aprile 2019) , <https://www.census.gov/topics/research/data-science/about-machine-learning.html> (in cui si afferma che il Census Bureau "ha bisogno" di machine learning capacità); Ufficio di presidenza dell'Econ. Analisi, Piano d'azione strategico 2020 7 (2020) (sottolineando l'importanza dell'intelligenza artificiale e dell'apprendimento automatico per la strategia di BEA).

43 Le analisi dei dati a livello di gruppo comportano anche rischi e danni inerenti alla privacy. Vedi, ad esempio, Linnet Taylor, Safety in Numbers? Privacy di gruppo e analisi dei big data nel mondo in via di sviluppo, in Privacy di gruppo: nuove sfide delle tecnologie dei dati 13 (2017).

44 Cfr. 34 USC §§ 41303(c)(2), (3), (4).

45 Nat'l Acad. of Sci., Innovazioni nella statistica federale 41 (2017).

46 Cfr. 13 USC § 6.

47 Nat'l Acad. of Sci., supra nota 45, a 40.

48 Secondo lo studio, "il consulente legale di un'agenzia può sconsigliare la condivisione dei dati come misura precauzionale piuttosto che a causa di un divieto esplicito". Ufficio per la responsabilità del governo degli Stati Uniti, sforzi sostenuti e coordinati potrebbero facilitare la condivisione dei dati proteggendo la privacy 1 (2013).

49 Cfr. Amy O'Hara e Carla Medalia, Data Sharing in the Federal Statistical System: Impediments and Possibilities, 675 Annals Am. Acad. pol. & Soc. Sci. 138, 141 (2018).

50 Robert M. Groves e Adam Neufeld, Accelerare la condivisione dei dati tra i settori per promuovere il bene comune 12 (2017).

51 Bipartisan Pol'y Ctr., supra nota 6, pp. 18-20.

52 Cfr. O'Hara & Medalia, supra nota 49, p. 141.

53 Administrative Data Research UK, ADR UK, <https://www.adruk.org>.

54 Informazioni su ADR UK, ADR UK, <https://www.adruk.org/about-us/about-adr-uk/>.

55 id.

56 Cfr. World of Work, ADR UK, <https://www.adruk.org/our-work/world-of-work/>.

57 Opportunità di finanziamento, ADR UK, <https://www.adruk.org/news-publications/funding-opportunities/>.

58 id.

59 id.

60 Funding Opportunity : A Unique Chance to Shape Data Science at the Heart of UK Government, ADR UK (8 aprile 2021), <https://www.adruk.org/news-publications/news-blogs/funding-opportunity-> un'opportunità-unica-di-formare-la-scienza-dei-dati-nel-cuore-del-governo-britannico-384/.

61 Opportunità di finanziamento, supra nota 57.

62 Legge sull'economia digitale 2017 (Gr. Br.).

63 ADR UK, Fiducia, sicurezza e interesse pubblico: trovare l'equilibrio 28 (2020).

64 id.

65 id.

66 Come lavoriamo con i ricercatori?, ADR UK, <https://www.adruk.org/our-mission/working-with-researchers/>.

67 Accesso ai dati di ricerca sicuri come ricercatore accreditato, off. per Nat'l Stat., <https://www.ons.gov.uk/aboutus/whatwedo/statistics/requesting-statistics/approvedresearcherscheme>.

68 Cfr. Nick Hart e Nancy Potok, Modernizing US Data Infrastructure: Design Considerations for Implementing a National Secure Data Service to Improve Statistics and Evidence Building 17, 21 (2020).

69 id.

70 id. alle 15.

71 Amministratore Delegato. Ordine del giorno, supra nota 11, ore 9.

72 id. al 31.

73 Vedere Cos'è Open Data?, Open Data Handbook, <https://opendatahandbook.org/guide/en/what-is-open-data/>.

Capitolo 6

1 Vedi Mantenere i segreti: i dati anonimi non sono sempre anonimi, Berkeley Sch. di Info. (15 marzo 2014), <https://ischoolonline.berkeley.edu/blog/anonymous-data/> ; Arvind Narayanan e Vitaly Shmatikov, How to Break Anonymity of the Netflix Prize Dataset, Cornell U. (22 novembre 2007), <https://arxiv.org/pdf/cs/0610105.pdf> .

2 Matt Fredrikson, Somesh Jha & Thomas Ristenpart, Model Inversion Attacks that Exploit Confidence Information and Basic Countermeasures, 22 Proceedings of the ACM Special Interest Group on Security, Audit & Control 1322 (2015); Nicholas Carlini et al., Estrazione di dati di addestramento da modelli di linguaggi di grandi dimensioni, Cornell U. (15 giugno 2021), <https://arxiv.org/pdf/2012.07805.pdf>.

3 Si veda, ad esempio, HIPAA Training, Certification, and Compliance, HIPAA Training, <https://www.hipaatraining.com/>; Ricerca sulla gestione dei dati, servizio dati del Regno Unito, <https://ukdataservice.ac.uk/learning-hub/research-data-management/>.

4 Ashwin Machanavajjhala et al., L-Diversity: Privacy Beyond K-Anonymity, 22 Int'l Conf. Data Eng'g 24 (2006).

5 Cynthia Dwork e Aaron Roth, I fondamenti algoritmici della privacy differenziale (2014).

6 Cfr., ad esempio, Tara Bahrampour e Marissa J. Lang, New System to Protect Census Data May Compromise Accuracy, Some Experts Say, Wash.

Post (1 giugno 2021), https://www.washingtonpost.com/local/social-issues/2020-census-differential-privacy-ipums/2021/06/01/6c94b46e-c30d-11eb-93f5-ee9558eef4b_story.html; Kelly Percival, il tribunale rifiuta la sfida dell'Alabama ai piani di censimento per la riorganizzazione distrettuale e la

privacy, Brennan Ctr. (30 giugno 2021), <https://www.brennancenter.org/our-work/analysis-opinion/court-rejects-alabama-challenge-census-plans-redistricting-and-privacy> .

7 Si veda, ad esempio, Leonard E. Burman et al., Safely Expanding Research Access to Administrative Tax Data: Creating a Synthetic Public Use File and a Validation Server (2018); vedere anche The Synthetic Data Vault, <https://sdv.dev>.

8 Valerie Chen, Valerio Pastro e Mariana Raykova, Secure Computation for Machine Learning with SPDZ, Cornell U. (2 gennaio 2019), <https://arxiv.org/pdf/1901.00329.pdf>.

9 Louis JM Aslett et al., A Review of Homomorphic Encryption and Software Tools for Encrypted Statistical Machine Learning, Cornell U. (26 agosto 2015), <https://arxiv.org/pdf/1508.06574.pdf>.

10 Si veda Hongyan Chang e Reza Shokri, On the Privacy Risks of Algorithmic Fairness, Cornell U. (7 aprile 2021), <https://arxiv.org/pdf/2011.03731.pdf>.

11 Ruggles et al., Differential Privacy and Census Data: Implications for Social and Economic Research, 109 Am. Econ. Ass'n Papers & Proceedings 403, 406 (2019).

12 Nella letteratura informatica, tali impostazioni algoritmiche sono spesso indicate come iperparametri. Ad esempio, k è un iperparametro per k -anonimità. Impostando k su valori diversi (ad esempio, 5, 10, 100), i professionisti possono modulare la quantità di anonimato concessa ai record nei dati. Come notiamo, tuttavia, la scelta degli iperparametri controlla sia la privacy effettuata su un set di dati sia la fedeltà di tali dati.

13 Cfr. La spiegazione della privacy differenziale per i dati del censimento, Nat'l Conf. of State Legislatures (1 luglio 2021), <https://www.ncsl.org/research/redistricting/differential-privacy-for-census-data-explained.aspx>; Hongyan Chang e Reza Shokri, On the Privacy Risks of Algorithmic Fairness, Cornell U. (7 aprile 2021), <https://arxiv.org/pdf/2011.03731.pdf>. Alcuni studiosi ritengono addirittura che l'incorporazione della privacy differenziale negli algoritmi di apprendimento automatico possa avere un impatto disparato sui gruppi sottorappresentati. Si veda Eugene Bagdasaryan e Vitaly Shmatikov, Differential Privacy Has Disparate Impact on Model Accuracy, Cornell U. (27 ottobre 2019), <https://arxiv.org/pdf/1905.12101.pdf>.

14 Steven Ruggles, Differential Privacy and Census Data: Implications for Social and Economic Research 17.

15 id. alle 18-19.

16 Nat'l Acad. of Sci., Innovazioni nella statistica federale 86 (2017). Il frammentato processo di revisione del FSRDC è simile al regime frammentato di accesso ai dati di cui abbiamo discusso nel terzo capitolo.

17 Programma Speciale Ricercatore Giurato, Bureau of Econ. Analisi, <https://www.bea.gov/research/special-sworn-researcher-program> (ultimo aggiornamento 23 luglio 2021).

18 13 USC § 9.

19 La forma istituzionale del CNR è discussa in modo approfondito nel capitolo quarto.

20 Enclave dati NORC, NORC, <https://www.norc.umd.edu/PDFs/BD-Brochures/2016/Data%20Enclave%20One%20Sheet.pdf>.

21 CMS Virtual Research Data Center (VRDC), Research Data Assistance Ctr., <https://resdac.org/cms-virtual-research-data-center-vrdc>.

22 Richiesta di informazioni (RFI) In cerca di input da parte delle parti interessate sulla necessità di un'enclave di dati amministrativi NIH, Nat'l Inst. of Health (1 marzo 2019), <https://grants.nih.gov/grants/guide/notice-files/NOT-OD-19-085.html>.

- 23 Cfr. FASEB Response to NIH Request for Information (RFI): Seeking Stakeholder Input on the Need for an NIH Administrative Data Enclave, Fed'n of Am. Società per la biologia sperimentale (2019), https://www.faseb.org/Portals/2/PDFs/opa/2019/FASEB_Response_Data_Enclave_RFI_NOT-OD-19-085_PDF; Sono. Soc'y of Biochemistry & Molecular Biology (30 maggio 2019), <https://www.asmb.org/getmedia/e3401ed5-3210-4ed2-a82a-7363cb86071d/ASBMB-Response-to-NIH-RFI-NOT-09-19-085.pdf>.
- 24 Cosa facciamo, Cal. Pol'y Lab, <https://www.capolicylab.org/what-we-do/>.
- 25 id.
- 26 CPL Roadmap per i dati amministrativi del governo in California, Cal. Pol'y Lab, <https://www.capolicylab.org/data-resources/california-data-road-map/>.
- 27 Intervista con Evan White, direttore esecutivo, California Policy Lab (29 aprile 2021).
- 28 id.
- 29 id.; si veda, ad esempio, Policy Evaluation and Research Linkage Initiative (PERLI), Cal. Pol'y Lab, <https://www.capolicylab.org/data-resources/perli/>; Gruppo di esperti sul credito al consumo dell'Università della California, Cal. Pol'y Lab, <https://www.capolicylab.org/data-resources/university-of-california-consumer-credit-panel/>.
- 30 Intervista con Evan White, supra nota 27.
- 31 id.
- 32 Cfr., ad esempio, Life Course Dataset, Cal. Pol'y Lab, <https://www.capolicylab.org/life-course-dataset/>.
- 33 Cfr. CPL Roadmap to Government Administrative Data in California, supra nota 26.
- 34 Notiamo che è possibile che la forma organizzativa possa influenzare l'autorità del personale NRC di parlare della legalità dei trasferimenti di dati.

Capitolo 7

- 1 Christopher Whyte, Deepfake News: Disinformazione abilitata dall'intelligenza artificiale come sfida di politica pubblica a più livelli, 5 J. Cyber Pol'y 199 (2020); Don Fallis, Che cos'è la disinformazione?, 63 Libr. Tendenze 601 (2015).
- 2 Mary L. Gray e Siddharth Suri, Ghost Work: How to Stop Silicon Valley From Building a New Global Underclass (2019); La scienza deve esaminare il futuro del lavoro, Nature (19 ottobre 2017), <https://www.nature.com/articles/550301b>.
- 3 David Danks & Alex John London, Algorithmic Bias in Autonomous Systems, 26 Int'l Joint Conf. sull'Intelligenza Artificiale 4691 (2017); Joy Buolamwini & Timnit Gebru, Sfumature di genere: disparità di accuratezza intersezionale nella classificazione di genere commerciale, 81 Proceeding of Machine Learning Res. 1 (2018); Ben Hutchinson et al., Pregiudizi involontari dell'apprendimento automatico come barriere sociali per le persone con disabilità, 125 ACM SIGACCESS Accessibility & Computing 1 (2020).
- 4 Oscar H. Gandy Jr., The Panoptic Sort: A Political Economy of Personal Information (1993); Virginia Eubanks, Automatizzare la disuguaglianza (2018); Rashida Richardson, Segregazione razziale e società guidata dai dati: come la nostra incapacità di fare i conti con le cause profonde perpetua realtà separate e disuguali, 36 Berkeley Tech. LJ 101 (2021).
- 5 Per gli approcci per migliorare le pratiche di machine learning, vedere Timnit Gebru et al., Datasheets for Datasets, Cornell U. (19 marzo 2020), <https://arxiv.org/pdf/1803.09010.pdf>; Margaret Mitchell et al., Model Cards for Model Reporting, 2019 Proceedings ACM Conf. su equità, responsabilità e trasparenza 220 (2019); Kenneth Holstein et al., Miglioramento dell'equità nei sistemi di apprendimento automatico: di cosa hanno bisogno i professionisti del settore?, 2019 CHI Conf. su Hum. Fattori nel sistema informatico. 1 (2019); Michael A. Madaio et al., Checklist di co-progettazione per comprendere le sfide e le opportunità organizzative relative all'equità nell'IA, 2020 CHI Conf. su Hum. Fattori nel sistema informatico. 318 (2019). La letteratura sugli impatti sociali dell'IA e sull'equità, responsabilità e trasparenza dell'IA è vasta, ma si veda Michael Kearns e Aaron Roth, The Ethical Algorithm: The Science of Socially Aware Algorithm Design (2019); Eubanks, supra nota 4; Solon Barocas, Moritz Hardt e Arvind Narayanan, Correttezza e apprendimento automatico (2019); Cathy O'Neil, Armi di distruzione matematica (2016).
- 6 45 CFR §§ 46.101-124.
- 7 J. Britt Holbrook & Robert Frodeman, Peer Review and the Ex Ante Assessment of Societal Impacts, 20 Res. Valutazione 239 (2011).
- 8 id.
- 9 Institutional Review Boards (IRBs) and Protection of Human Subjects in Clinical Trials, US Food & Drug Admin., <https://www.fda.gov/about-fda/center-for-drug-evaluation-and-research-cder/institutional-review-boards-irbs-and-protection-human-subjects-clinical-trials> (ultimo aggiornamento 11 settembre 2019).
- 10 Ci sono questioni cruciali per quanto riguarda il consenso anche con dati considerati "pubblicamente" disponibili. Vedere in generale Casey Fiesler e Nicholas Proferes, "Participant" Perceptions of Twitter Research Ethics, 4 Social Media + Society 1 (2018); Sarah Gilbert, Jessica Vitak e Katie Shilton, Misurare il comfort degli americani con gli usi della ricerca dei loro dati sui social media, 7 Social Media + Society 1 (2021).
- 11 Sara R. Jordan, Future of Privacy Forum, Designing an Artificial Intelligence Research Review Committee (2019), <https://fpf.org/wp-content/uploads/2019/10/DesigningAIResearchReviewCommittee.pdf>.
- 12 Agata Ferretti et al., Ethics Review of Big Data Research: What Should Stay and What Should be Reformed?, 22 BMC Medical Ethics 1, 6 (2021); Kathryn M. Porter et al., The Emergence of Clinical Research Ethics Consultation: Insights from a National Collaborative, 2018 Am. J. Bioetica 39 (2018).
- 13 Ferretti et al., supra nota 12.
- 14 Si veda, ad esempio, Mark Diaz et al., Affronting Age-Related Bias in Sentiment Analysis, 2018 Proceedings CHI Conf. su Hum. Fattori nel sistema informatico. 1 (2018); Buolamwini & Gebru, supra nota 3.
- 15 Cfr., ad esempio, Timnit Gebru et al., Datasheets for Datasets, Cornell U. (19 marzo 2020), <https://arxiv.org/pdf/1803.09010.pdf>; Margaret Mitchell et al., Model Cards for Model Reporting, 2019 Proceedings ACM Conf. su equità, responsabilità e trasparenza 220 (2019); Emily M. Bender et al., On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?, 2021 Proceedings ACM Conf. su Equità, Responsabilità e Trasparenza 610 (2021); Christo Wilson et al., Building and Auditing Fair Algorithms: A Case Study in Candidate Screening, 2021 Proceedings ACM Conf. sull'equità, la responsabilità e la trasparenza 666 (2021); Pauline T. Kim, Auditing Algorithms for Discrimination, 166 U. Pa. L. Rev. Online 189 (2017).
- 16 Fase II: revisione ed elaborazione della proposta, Nat'l Sci. Trovato., https://www.nsf.gov/bfa/dias/policy/merit_review/phase2.jsp#select.
- 17 Harvey A. Averbach, Criteri per la valutazione di progetti e portafogli di ricerca, nella valutazione degli impatti di R&S: metodi e pratica 263 (1993).
- 18 National Security Comm'n on Artificial Intelligence, Final Report 141-54 (2021).
- 19 History and Mission, US Privacy & Civil Liberties Oversight Bd., <https://www.pclbo.gov/About/HistoryMission>.
- 20 AI in Counterterrorism Oversight Enhancement Act del 2021, HR 4469, 117th Cong. (2021).
- 21 Nei casi in cui un ricercatore utilizza i dati ottenuti da una delle agenzie che ricadono sotto la supervisione dell'ORI, può avere senso chiedere all'ORI di adeguare

dicare questi casi direttamente. Per ulteriori informazioni sull'ORI, vedere ORI, Office Of Research Integrity, <https://ori.hhs.gov/>.

22 Michael S. Bernstein et al., ESR: Ethics and Society Review of Artificial Intelligence Research, Cornell U. (9 luglio 2021), <https://arxiv.org/pdf/2106.11521.pdf>.

23 Nat'l Sci. Trovato., Impatti più ampi, <https://www.nsf.gov/od/oia/special/broaderimpacts/>.

24 Si veda, ad esempio, Avviso di interesse speciale: Supplemento amministrativo per gli sforzi di ricerca e sviluppo delle capacità relativi a questioni bioetiche, Nat'l Inst. of Health (17 novembre 2020), <https://grants.nih.gov/grants/guide/notice-files/NOT-OD-21-020.html>; Avviso di interesse speciale: Supplemento amministrativo per la ricerca su questioni bioetiche, Nat'l Inst. of Health (30 dicembre 2019), <https://grants.nih.gov/grants/guide/notice-files/NOT-OD-20-038.html>; si veda anche Courtenay R. Bruce et al., An Embedded Model for Ethics Consultation: Characteristics, Outcomes, and Challenges, 5 AJOB Empirical Bioethics 8 (2014); Sharon Begley, In a Lab Pushing the Boundaries of Biology, an Embedded Ethicist Keeps Scientists in Check, Stat (23 febbraio 2017), <https://www.statnews.com/2017/02/23/bioethics-harvard-george-church/>. Anche le fondazioni private promuovono l'uso di esperti di bioetica incorporati. Si veda, ad esempio, Making a Difference Request for Proposals – Fall 2021, The Greenwall Found. (2021), <https://greenwall.org/making-a-difference-grants/request-for-proposals-MAD-autunno-2021>.

Capitolo 8

1 Paul Cichonski et al., Computer Security Incident Handling Guide (2012).

2 Gli attacchi presunti potrebbero includere la distribuzione di ransomware, schemi di phishing, l'ottenimento dell'accesso root (il più alto livello di privilegio disponibile che consente agli utenti di accedere a tutti i comandi e file per impostazione predefinita), l'esposizione di credenziali segrete, l'avvelenamento dei dati, l'esfiltrazione dei dati, nonché altri tipi di intrusioni di rete non autorizzate.

3 Karen Hao, AI consuma molta energia. Gli hacker potrebbero fargli consumare di più, MIT Tech. R. (6 maggio 2021), <https://www.technologyreview.com/2021/05/06/1024654/ai-energy-hack-adversarial-attack/>.

4 Catalin Cimpanu, Vast Majority of Cyber-Attacks on Cloud Servers Aim to Mine Cryptocurrency, ZDNet (14 settembre 2020), <https://www.zdnet.com/article/vast-majority-of-cyber-attacks-on-server-cloud-obiettivo-di-minare-criptoaluta/>.

5 Notiamo che l'NRC dovrà probabilmente conformarsi anche a norme di sicurezza specifiche per i dati. Ad esempio, la sicurezza dei dati medici dovrà essere conforme all'HIPAA e i dati finanziari dovranno essere conformi al Gramm-Leach-Bliley Act.

6 Ray Dunham, FISMA Compliance: Security Standards & Guidelines Overview, Linford & Co. (29 novembre 2017), <https://linfordco.com/blog/fisma-compliance/>.

7 Amy J. Frontz, Revisione della conformità del Dipartimento della salute e dei servizi umani alla legge federale sulla modernizzazione della sicurezza delle informazioni del 2014 per l'anno fiscale 2020 (2021).

8 Senato degli Stati Uniti comm. sulla sicurezza interna e gli affari governativi, Federal Cybersecurity: America's Data at Risk 18 (2019).

9 Legge federale sulla modernizzazione della sicurezza informatica (FISMA) Background, Nat'l Inst. Standards & Tech., <https://csrc.nist.gov/projects/risk-management/fisma-background> (ultimo aggiornamento 4 agosto 2021).

10 Dunham, supra nota 6.

11 Senato degli Stati Uniti comm. su Homeland Security & Governmental Affairs, supra nota 8, 19.

12 id. alle 18.

13 id. alle 19.

14 id. alle 20.

15 Kevin Stine, et al., Guide for Mapping Types of Information and Information Systems to Security Categories (2008). In particolare, FISMA definisce la conformità in termini di tre livelli: basso impatto, impatto moderato e alto impatto. Basso impatto indica che la perdita di riservatezza, integrità o disponibilità del sistema avrà un effetto negativo limitato, mentre alto impatto indica che tali perdite avranno effetti gravi o catastrofici. Vedere Sarah Harvey, 3 livelli di conformità FISMA: basso, moderato, alto, KirkpatrickPrice (24 aprile 2020), <https://kirkpatrickprice.com/blog/fisma-compliance-levels-low-moderate-high/>.

16 Nat'l Inst. Standard e tecnologia, controlli di sicurezza e privacy per sistemi informativi e organizzazioni (2020).

17 Marianne Swanson et al., Guide for Developing Security Plans for Federal Information Systems (2006).

18 Michael McLaughlin, Reforming FedRAMP: una guida per migliorare l'approvvigionamento federale e la gestione del rischio dei servizi cloud, info. Tecnico. & Innovazione trovata. (15 giugno 2020), <https://itif.org/publications/2020/06/15/reforming-fedramp-guide-improving-federal-procurement-and-risk-management>.

19 Nozioni di base sul programma, FedRAMP, <https://www.fedramp.gov/program-basics/>; si veda anche Steven VanRoekel, Security Authorization of Information Systems in Cloud Computing Environments (2011).

20 FISMA vs. FedRAMP e NIST: dare un senso agli standard di conformità del governo, Foresite, <https://foresite.com/fisma-vs-fedramp-and-nist-making-sense-of-government-compliance-standards/>. Tuttavia, notiamo che l'approvazione FedRAMP è esentata per alcuni tipi di modelli cloud: (i) dove il cloud è privato dell'agenzia, (ii) dove il cloud si trova fisicamente all'interno di una struttura federale, (iii) dove l'agenzia non è fornitrice di servizi cloud dal sistema informativo basato su cloud a qualsiasi entità esterna. Vedi VanRoekel, supra nota 19.

21 FedRAMP, FedRAMP Security Assessment Framework 5 (2017).

22 Doina Chiacu, La Casa Bianca avverte le aziende di intensificare la sicurezza informatica: "Non possiamo farcela da soli", Reuters (3 giugno 2021), <https://www.cejving.com/technology/white-house-warns-aziende-step-up-cybersecurity-2021-06-03/>; si veda anche Incidenti informatici significativi, Ctr. Strategic & Int'l Studies, <https://www.csis.org/programs/strategic-technologies-program/significant-cyber-incidents> (ultima visita 19 agosto 2021).

23 Senato degli Stati Uniti comm. Homeland Security & Governmental Affairs, supra nota 8, 5.

24 id. alle 6.

25 Frontz, supra nota 7.

26 Jonathan Reiber e Matt Glenn, Il governo degli Stati Uniti deve rivedere la sicurezza informatica. Ecco come., Lawfare (9 aprile 2021), <https://www.lawfareblog.com/il-governo-americano-ha-bisogno-di-revisione-della-sicurezza-informatica-ecco-come>.

27 Nat'l Security Agency, Abbracciare un modello di sicurezza Zero Trust (2021).

28 McLaughlin, supra nota 18.

29 id.

30 id.

31 Esec. N. d'ordine 14.028, 86 Fed. Reg. 26633 (17 maggio 2021).

32 Ufficio di amministrazione degli Stati Uniti. & Budget, Spostare il governo degli Stati Uniti verso i principi di sicurezza informatica Zero Trust (2021).

28 Jukebox, OpenAI (30 aprile 2020), <https://openai.com/blog/jukebox/>.

29 Si veda, ad esempio, Shlomit Yanisky-Ravid, Generating Rembrandt: Artificial Intelligence, Copyright, and Accountability in the 3A Era--the Human-Like Authors are Already Here-- a New Model, 27 Mich. St. L. Rev. 659 (2017); Kalin Hristov, Intelligenza artificiale e dilemma del copyright, 57 J. Franklin Pierce Ctr. Intel. Puntello. 431 (2017).

30 Kalin Hristov, Intelligenza artificiale e indagini sul copyright, 16 J. Sci. Pol'y & Governance 1, 14-15 (2020).

31 id. alle 16.

32 Vedere What is Transfer Learning?, TensorFlow (31 marzo 2020), https://www.tensorflow.org/js/tutorials/transfer/what_is_transfer_learning.

33 Si veda, ad esempio, Yunhui Guo et al., SpotTune: Transfer Learning Through Adaptive Fine-Tuning, Cornell U. (novembre 2018), <https://arxiv.org/pdf/1811.08737.pdf>.

34 2 CFR § 200.315(d).

35 Cfr. Zhiqiang Wan, Yazhou Zhang e Haibo He, Variational Autoencoder Based Synthetic Data Generation for Imbalanced Learning, IEEE (2017).

36 Cfr. Noseong Park, Mahmoud Mohammadi e Kshitij Gorde, Sintesi dei dati basata su reti generative contraddittorie, 11 Proc. Dotazione VLDB 1071 (2018).

37 Cfr. Ron Bakker, Impact of Artificial Intelligence on IP Policy 12.

38 Cfr. Marta Duque Lizarralde, A Guideline to Artificial Intelligence, Machine Learning and Intellectual Property 4-7 (2020).

39 Steven M. Bellovin et al., Privacy e set di dati sintetici, 22 Stan. Tecnico. L. Ap 1, 2-3 (2019); si veda anche Fida K. Dankar e Mahmoud Ibrahim, Fake It Till You Make It: Guidelines for Effective Synthetic Data Generation, 5 Applied Sci. 11 (2021); ma si veda Theresa Stadler et al., Synthetic Data – Anonymisation Groundhog Day, Cornell U. (8 luglio 2021), <https://arxiv.org/pdf/2011.07018.pdf>.

40 Si veda, ad esempio, Daniel S. Quintana, A Synthetic Dataset Primer for the Biobehavioural Sciences to Promote Reproducibility and Hypothesis Generation, 9 eLife 1 (2020).

41 Yuji Roh et al., A Survey on Data Collection for Machine Learning, Cornell U. (12 agosto 2019), <https://arxiv.org/pdf/1811.03402.pdf>.

42 Cfr., ad esempio, Hang Qiu et al., Minimum Cost Active Labelling, Cornell U. (24 giugno 2020), <https://arxiv.org/pdf/2006.13999.pdf>; Eric Horvitz, Apprendimento automatico, ragionamento e intelligenza nella vita quotidiana: indicazioni e sfide, 18 Atti della conferenza. sull'incertezza nell'intelligenza artificiale 3 (2007).

43 Ricerca cognitiva, ingegneria dei dati, preparazione ed etichettatura per l'IA 2019 3 (2019).

44 Vedere Wil Michiels, How Do You Protect Your Machine Learning Investment?, EETimes (26 marzo 2020), <https://www.eetimes.com/how-do-you-protect-your-machine-learning-investment/>. Infatti, nell'Unione Europea, i set di dati etichettati vengono premiati con la protezione dei diritti sui database. Mauritz Kop, Apprendimento automatico e pratiche di condivisione dei dati dell'UE, Stan.-Vienna Transatlantic Tech. LF (24 marzo 2020), <https://ttfnews.wordpress.com/2020/03/24/machine-learning-eu-data-sharing-practices/>.

45 Si veda, ad esempio, Niklas Fiedler et al., ImageTagger: An Open Source Online Platform for Collaborative Image Labeling, 11374 Lecture Notes in Computer Sci. 162 (2019).

46 id. a 162.

47 I ricercatori possono, ad esempio, utilizzare i dati NRC e le risorse di calcolo per implementare strategie di apprendimento attivo, procedure per etichettare manualmente un sottoinsieme di dati disponibili e dedurre automaticamente le etichette rimanenti utilizzando un modello di apprendimento automatico. Si veda, ad esempio, Oscar Reyes et al., Effective Active Learning Strategy for Multi-Label Learning, 273 Neurocomputing 494 (2018). Allo stesso modo, i ricercatori possono integrare i dati esistenti del settore pubblico con etichette preziose.

48 Cfr., ad esempio, Pedro Saleiro et al., Aequitas: A Bias and Fairness Audit Toolkit, Cornell U. (29 aprile 2019), <https://arxiv.org/pdf/1811.05577.pdf>; Florian Tramèr et al., FairTest: Discovering Unwarranted Associations in Data-Driven Applications, Cornell U. (16 agosto 2019), <https://arxiv.org/pdf/1510.02377>. PDF.

49 Sebbene non discutiamo le modifiche idiosincratiche alla Guida uniforme che variano da agenzia ad agenzia, incoraggiamo la task force a valutare queste modifiche se decide di implementare l'NRC attraverso una particolare agenzia. Se l'NRC è amministrato attraverso più agenzie, il complesso amalgama di regole IP specifiche dell'agenzia può aumentare l'attrito nell'utilizzo dell'NRC se i ricercatori devono cambiare contesto da una serie di regolamenti all'altro a seconda dell'agenzia finanziatrice.

50 2 CFR § 2900.13. In precedenza, il Dipartimento del lavoro richiedeva esplicitamente che la proprietà intellettuale generata nell'ambito di un premio federale fosse concessa in licenza con una licenza Creative Commons Attribution, ma questa regola è stata modificata nell'aprile 2021 per sostituire il termine proprietario "licenza Creative Commons Attribution" con lo standard riconosciuto dal settore "licenza aperta". 86 federale. Reg. 22107 (27 aprile 2021).

51 Diffusione e condivisione dei risultati della ricerca - Requisiti del piano di gestione dei dati NSF, Nat'l Sci. Trovato., <https://www.nsf.gov/bfa/dias/policy/dmp.jsp>.

52 Si veda, ad esempio, Aidan Courtney et al., Balancing Open Source Stem Cell Science with Commercialization, Nature Biotechnology (7 febbraio 2011), <https://www.nature.com/articles/nbt.1773>.

53 Cfr. Clint Finley, When Open Source Software Vienes with a Few Catches, Wired (31 luglio 2019), <https://www.wired.com/story/when-open-source-software-comes-with-catches/>; Guida alle licenze open source, Synopsys (7 ottobre 2016), <https://www.synopsys.com/blogs/software-security/open-source-licenses/>.

54 Cfr. Daniel A. Almeida et al., Gli sviluppatori di software comprendono le licenze open source?, 25 IEEE Int'l Conf. su Program Comprehension 1 (2017) (trovando che gli sviluppatori di software "lottano [] quando più licenze [open-source] [sono] coinvolte" e "mancano della conoscenza e della comprensione per separare le interazioni delle licenze in più situazioni").

55 Cfr., ad esempio, Alexandra Theben et al., Challenges and Limits of an Open Source Approach to Artificial Intelligence 14 (2021); Stadler et al., supra nota 39; Milad Nasr et al., Comprehensive Privacy Analysis of Deep Learning: Passive and Active White-box Inference Attacks against Centralized and Federated Learning, Cornell U. (6 giugno 2020), <https://arxiv.org/abs/1812.00910>. PDF.

56 Alcune università hanno deciso di eliminare la ricerca classificata. Vedi, per esempio, At the Hands of Radicals, Stan. Mag. (gennaio 2009), <https://stanfordmag.org/contents/at-the-hands-of-the-radicals>.

57 Cfr. Donald Kennedy, Scienza e segretezza, 289 Sci. 724 (2000); Peter J. Westwick, Scienza segreta: una comunità classificata nei laboratori nazionali, 38 Minerva 363 (2000).

58 Cfr. Braun & Ong, supra nota 3; Sören Sonnenburg et al., The Need for Open Source Software in Machine Learning, 8 J. Machine Learning Res. 2443 (2007); si veda anche Katie Malone e Richard Wolski, Doing Data Science on the Shoulders of Giants: The Value of Open Source Software for the Data Science Community, HDSR (31 maggio 2020), <https://hdsr.mitpress.mit.edu/pub/xsrt4zs2/rilascio/4>.

59 Cfr. Laura A. Heymann, Overlapping Intellectual Property Doctrines: Election of Rights Versus Selection of Remedies, 17 Stan. Tecnico. L. Rev. 239, 240 (2013); Oracle Am. Inc. v. Google Inc., 750 F.3d 1339 (Fed. Cir. 2014) (accettando che il software è sia brevettabile che protetto da copyright).

- 60 Robert E. Thomas, Debug dei brevetti software: aumentare l'innovazione e ridurre l'incertezza nella riforma giudiziaria della legge sui brevetti software, 25 Santa Clara Computer & High Tech. LJ 191, 222-23 (2008).
- 61 Cfr., ad esempio, Joaquin Vanschoren et al., OpenML: Networked Science in Machine Learning, Cornell U. (Aug. 1, 2014), <https://arxiv.org/pdf/1407.7722.pdf> (developing a collaboration platform through which gli scienziati possono condividere, organizzare e discutere automaticamente esperimenti, dati e algoritmi di machine learning); vedi anche Sarah O'Meara, AI Researchers in China Want to Keep the Global-Sharing Culture Alive, Nature (29 maggio 2019), <https://www.nature.com/articoli/d41586-019-01681-x>; Shuai Zhao et al., Confezione e condivisione di modelli di apprendimento automatico tramite la piattaforma aperta Acumos AI, 17 ICMLA (2018).
- 62 Jeanne C. Fromer, Machines as the New Oompa-Loompas: Trade Secrecy, the Cloud, Machine Learning, and Automation, 94 NYUL Rev. 706, 712 (2019); Jordan R. Raffae et al., The Rising Importance of Trade Secret Protection for AI-Related Intellectual Property 1, 5-6 (2020); Jessica M. Meyers, Intelligenza artificiale e segreti commerciali, Am. Bar Ass'n (febbraio 2019), https://www.americanbar.org/groups/intellectual_property_law/publications/landslide/2018-19/january-february/artificial-intelligence-trade-secrets-webinar/; Commenti dell'AIPLA in merito alla "Richiesta di commenti sulla protezione della proprietà intellettuale per l'innovazione dell'intelligenza artificiale", em. Intell. Prop. L. Ass'n (10 gennaio 2020), https://www.uspto.gov/sites/default/files/documents/AIPLA_RFC-84-FR-58141.pdf.
- 63 Clark D. Asay, Stupidità artificiale, 61 Wm. & Mary L. Rev. 1187, 1197, 1241-42 (2020).
- 64 Cfr. id.; Sono. Intel. Puntello. L. Ass'n, supra nota 62, at 16.
- 65 Cfr. Asay, supra nota 63, 1242.

Appendice

- 1 Il Department of Energy assegna 425 milioni di dollari per le tecnologie di supercalcolo di nuova generazione, Energy.gov (14 novembre 2014), <https://www.energy.gov/articles/department-energy-awards-425-million-next-generation-technologie-di-supercalcolo>.
- 2 Istanze P3 di Amazon EC2, Amazon, <https://aws.amazon.com/ec2/instance-types/p3/> (ultima visita 9 settembre 2021).
- 3 CORAL Richiesta di proposta B604142, Lawrence Livermore Nat'l Laboratory (2014), <https://web.archive.org/web/20140816181824/> <https://asc.llnl.gov/CORALLO/>.
Notiamo che non siamo stati in grado di individuare i documenti di aggiudicazione finali, né il Summit è stato preventivato in modo sufficientemente dettagliato per sostenere i costi dalle dichiarazioni di bilancio del DOE. Le nostre stime dei costi qui, tuttavia, sono paragonabili alle stime riportate pubblicamente per il costo totale del sistema Summit.
- 4 Questo si basa su un massimo di 30 milioni di dollari nel contratto DOE Office of Science per i costi di ingegneria non ricorrenti (NRE) per i sistemi dell'Argonne National Laboratory e dell'Oak Ridge National Laboratory.
- 5 Ciò si basa sulla differenza nei termini RFP tra l'inclusione della manutenzione nell'ambito del sistema Lawrence Livermore National Laboratory (con un budget massimo di \$ 170 milioni) e l'esclusione della manutenzione nell'ambito dei sistemi per l'Oak Ridge National Laboratory e l'Argonne National Laboratory Laboratorio (con un budget massimo per il contratto di costruzione di \$ 155 milioni). Questo è probabilmente un limite superiore per il mantenimento, dato che la differenza riflette la combinazione di NRE e mantenimento quinquennale.
- 6 Cfr. listino prezzi CORAL, Lawrence Livermore Nat'l Laboratory (2014), <https://web.archive.org/web/20140816181824/> https://asc.llnl.gov/CORAL/RFP_components/04_CORAL_Price_Schedule_ANL_ORNL_tabs.xlsx. Abbiamo utilizzato l'1,62% come tasso di interesse per calcolare il costo su 60 mesi. È il tasso a scadenza costante del Tesoro a 5 anni al 14 novembre 2014, vedere Selected Interest Rates (Daily) – H.15, Fed. Res., <https://www.federalreserve.gov/releases/H15/default.htm>, quando il DOE ha annunciato l'aggiudicazione del sistema HPC, vedere Department of Energy Awards \$425 Million for Next Generation Super computing Technologies, supra 1.
- 7 Ad esempio, questa stima è in linea con il costo di 200 milioni di dollari riportato dal New York Times. Steve Lohr, Move Over, China: US is Again Home to World's Speediest Supercomputer, NY Times (8 giugno 2018), <https://www.nytimes.com/2018/06/08/technology/supercomputer-china-us.html>. Alcuni rapporti confondono l'approvvigionamento di più sistemi avvenuti contemporaneamente.
- 8 La ricerca mostra che per l'addestramento di modelli di deep learning ad alta intensità di calcolo, come ResNet-101, l'utilizzo della GPU è di circa il 70%. Jingoo Han et al., A Quantitative Study of Deep Learning Training on Heterogeneous Supercomputers, 2019 IEEE Conf. su Cluster Computing 1, 5 (2019). Tuttavia, ResNet-50 ha un utilizzo della GPU di circa il 40%, vedi id., e altri account riportano che le GPU vengono utilizzate solo il 15-30% delle volte, vedi, ad esempio, Lukas Biewald, Monitor and Improve GPU Usage for Training Deep Modelli di apprendimento, verso la scienza dei dati. (27 marzo 2019), <https://towardsdatascience.com/measuring-actual-gpu-usage-for-deep-learning-training-e2bf3654bcfd>; Janet Morss, Giving Your Data Scientists a Boost with GPUaaS, CIO (2 giugno 2020), <https://www.cio.com/article/3561090/giving-your-data-scientists-a-boost-with-gpuaaS.html>.
- 9 Compute Canada, Cloud Computing for Researchers 1 (2016), <https://www.computeCanada.ca/wp-content/uploads/2015/02/CloudStrategy2016-2019-forresearchersEXTERNAL-1.pdf>.
- 10 Jennifer Shkabatur, The Global Commons of Data, 22 Stan. Tecnico. LR 407, 407-09 (2019).
- 11 Benjamin Sobel, La crisi del fair use dell'intelligenza artificiale, 41 colum. JL & Arti 61 (2017).
- 12 id.
- 13 Vedere Proteggere ciò che amiamo di Internet: i nostri sforzi per fermare la pirateria online, Google Pub. Pol'y Blog (7 novembre 2019), <https://www.blog.google/outreach-initiatives/public-policy/protecting-what-we-love-about-internet-our-efforts-stop-online-pirateria/>.
- 14 Cfr. Jennifer M. Urban, Joe Karaganis e Brianna M. Schofield, Notice & Takedown in Everyday Practice 39 (2017) (che illustra la difficoltà che i fornitori di servizi online incontrano nel valutare manualmente un grande volume di dati per una potenziale violazione; ad esempio, uno il fornitore di servizi online ha spiegato che "per paura di non riuscire a rimuovere il materiale in violazione e motivato dalla minaccia di danni legali, il suo personale prenderà" sei passaggi per cercare di trovare il [contenuto identificato]); vedi anche Lettera di Thom Tillis, Marsha Blackburn, Christopher A. Coons, Dianne Feinstein et. al, a Sundar Pichai, amministratore delegato, Google Inc. (3 settembre 2019), <https://www.ipwatchdog.com/wp-content/uploads/2019/09/9.3-Content-ID-Ltr.pdf> ("Abbiamo sentito da titolari di copyright a cui è stato negato l'accesso agli strumenti di Content ID e, di conseguenza, si trovano in notevole svantaggio nell'impedire il caricamento ripetuto di contenuti che hanno precedentemente identificato come in violazione. A loro resta la scelta di spendere ore ogni settimana ricercando e inviando avvisi sulle stesse opere protette da copyright, o permettendo che la loro proprietà intellettuale venga sottratta").
- 15 Cfr. Google, How Google Fights Piracy 6 (2016). Per illustrare i costi dell'implementazione di Content ID su una piattaforma su larga scala, Google ha annunciato in un rapporto del 2016 che YouTube aveva investito più di 60 milioni di dollari in Content ID.
- 16 Cfr. Sobel, supra nota 11, pp. 66-79.
- 17 Cfr. Authors Guild c. Google Inc., 804 F.3d 202 (2d Cir. 2015).
- 18 id.
- 19 id. a 216-17.
- 20 Matthew Stewart, La decisione giudiziaria più importante per la scienza dei dati e l'apprendimento automatico, verso la scienza dei dati. (31 ottobre 2019), <https://towardsdata->

science.com/the-most-important-supreme-court-decision-for-data-science-and-machine-learning-44cfc1c1bcaf.

21 Si veda, ad esempio, James Grimmelmann, Copyright for Literate Robots, 101 Iowa L. Rev. 657, 661; Sobel, supra nota 11, pp. 51-57.

22 Cfr. Sobel, supra nota 11, p. 57.

23 Cfr. Anna I. Krylov et. al. Qual è il prezzo del software open source? 6 J. Physical Chemistry Letters 2751, 2753 (2015) (che spiega che i ricercatori in erba che considerano la commercializzazione possono essere particolarmente preoccupati per quali licenze sono disponibili, dal momento che un "ambiente strettamente open-source può inoltre disincentivare i giovani ricercatori a rendere immediatamente disponibile il nuovo codice, per timore che la loro capacità di pubblicare articoli venga messa in cortocircuito da un ricercatore più anziano con un esercito di postdoc pronti a trarre vantaggio da qualsiasi nuovo codice.).

24 Si veda, ad esempio, A Data Scientist's Guide to Open-Source Licensing, Towards Data Sci. (4 novembre 2018), <https://towardsdatascience.com/a-data-scientists-guide-to-open-source-licensing-c70d5fe42079>; Scegli una licenza open source, <https://choosealicense.com>.

25 Licenza di un repository, GitHub, <https://docs.github.com/en/github/creating-cloning-and-archiving-repositories/licensing-a-repository>.

26 Qual è la licenza più appropriata per i miei dati?, FigShare, <https://help.figshare.com/article/what-is-the-most-appropriate-licence-for-my-data>.

27 Cfr. Accordo per gli sviluppatori, Twitter (10 marzo 2020), <https://developer.twitter.com/en/developer-terms/agreement>; Uso non commerciale dell'API di Twitter, Twitter, <https://developer.twitter.com/en/developer-terms/commercial-terms>.

28 Cfr. Daniel A. Almeida et. al, Gli sviluppatori di software comprendono le licenze open source?, 25 IEEE Int'l Conf. sulla comprensione del programma 1 (2017).

29 id. alle 9.

30 Alexandra Kohn e Jessica Lange, confuse sul copyright? Valutazione della comprensione da parte dei ricercatori degli accordi di trasferimento del copyright, 6 J. Librarian Ship & Scholarly Comm'n. 1, 9 (2018).

31 Cfr. Will Frass, Jo Cross e Victoria Gardner, Taylor e Francis Open Access Survey giugno 2014 15 (2014). Si noti che la mancanza di alfabetizzazione IP potrebbe fungere da ulteriore deterrente per gli uploader. Il Taylor and Francis Open Access Survey del 2014 ha rilevato che "il 63% degli intervistati ha indicato la mancanza di comprensione della politica dell'editore come un fattore importante o molto importante nel non riuscire a depositare un articolo in un IR [Institutional Repository]". Id.

32 Norme comunitarie Dataverse, Harv. Dataverse, <https://dataverse.org/best-practices/dataverse-community-norms>.

33 Politica sul copyright e sulla licenza, FigShare, <https://help.figshare.com/article/copyright-and-license-policy>.

34 Australian Data Research Commons, Research Data Rights Managing Guide 6 (2019).

35 Cfr. Condizioni generali d'uso di Harvard Dataverse, Harv. Dataverse (2021), <https://dataverse.org/best-practices/harvard-dataverse-general-terms-use>.

36 Stan. U.Inst. dell'Intelligenza Artificiale incentrata sull'uomo, Rapporto sull'Indice di Intelligenza Artificiale 2021 125-34 (2021).

37 Thilo Hagendorff, The Ethics of AI Ethics: An Evaluation of Guidelines, 30 Minds & Machines 99 (2020).

38 Andrew D. Selbst, An Institutional View of Algorithmic Impact Assessments, 35 Harv. JL & Tech. 1, 66 (in uscita nel 2021).

39 Brent Mittelstadt, I principi da soli non possono garantire un'intelligenza artificiale etica, 1 Nature Mach. Intelligenza 501 (2019).

40 DOD adotta i principi etici per l'intelligenza artificiale, US Dep't Defense (24 febbraio 2020), <https://www.defense.gov/Newsroom/Releases/Release/Article/2091996/dod-adopts-ethical-principles-per-intelligenza-artificiale/>.

41 Mgmt del Presidente. Ordine del giorno, Strategia federale sui dati: quadro etico dei dati (2020).

42 Artificial Intelligence: An Accountability Framework for Federal Agencies and Other Entities, US Gov't Accountability Office (30 giugno 2021), <https://www.gao.gov/products/gao-21-519sp>.

43 Principi di etica dell'intelligenza artificiale per la comunità dell'intelligence, Ufficio del direttore dell'intelligence nazionale, <https://www.odni.gov/index.php/features/2763-principi-di-intelligenza-artificiale-etica-per-la-comunità-di-intelligence>.

44 Considerazioni chiave per lo sviluppo responsabile e la messa in campo dell'intelligenza artificiale, Nat'l Security Comm'n Artificial Intelligence (2021), <https://www.nscai.gov/key-considerations/>.

45 Pratiche raccomandate, Nat'l Security Comm'n Artificial Intelligence, <https://www.nscai.gov/wp-content/uploads/2021/01/Key-Considerations-Supporting-Visuals.pdf>.

46 Defense Innovation Bd., AI Principles: Recommendations on the Ethical Use of Artificial Intelligence by the Department of Defense (2019).



Stanford University
Human-Centered
Artificial Intelligence

Stanford
Law School