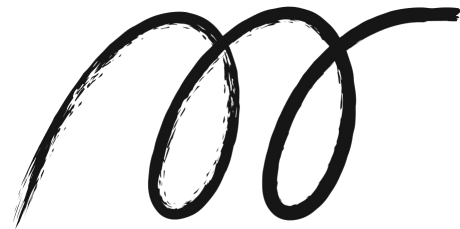


Mozilla Data Collective

Create. Curate. Control.



Your datasets, under your control



Acknowledgement of Country

Luis Mata,
Vanessa kershawi,

Royal Botanic Gardens Melbourne,
City of Melbourne, Victoria
CC-BY-NC-SA



Via <https://flic.kr/p/S3Sk7r>

Kathy Reid



Speaks with machines

Engineer. Data. Machine learning.
Voice AI.
Conscientious technologist.

Knitter.

❤️s your dog.

*Contains some replacement parts

** No I haven't finished my PhD yet, don't @ me

she/her pronouns (they/them totally fine, too)

@KathyReid in most of the places



Australian
National
University

School of
Cybernetics

<https://cybernetics.anu.edu.au/people/kathy-reid/>

Volunteers

<3

Agenda



01 **Tokenomics:**
why human-generated data is so valuable

02 **Token harvesting:**
Extractive and exclusive

03 **Is there a better way?**
YES! Mozilla Data Collective

01

Tokenomics

“
**Tokens are
building blocks
of data.**

Tokenisation example: Byte Pair

Step 0: Initial Tokenization (Characters)

t h e _ q u i c k _ b r o w n _ f o x _ j u m p e d _ o v e r _ t h e _ l a z y _ d o g

45 tokens (including spaces as '_')

Step 1: Find most frequent pair: 't' + 'h' (appears 2 times)

th e _ q u i c k _ b r o w n _ f o x _ j u m p e d _ o v e r _ th e _ l a z y _ d o g

43 tokens (merged 't' + 'h' → 'th')

Step 2: Next frequent pair: 'th' + 'e' → 'the'

the _ q u i c k _ b r o w n _ f o x _ j u m p e d _ o v e r _ the _ l a z y _ d o g

41 tokens

Step 3: Merge 'e' + 'd' → 'ed'

the _ q u i c k _ b r o w n _ f o x _ j u m p ed _ o v e r _ the _ l a z y _ d o g

40 tokens

...continue merging most frequent pairs...

Final Result: After many iterations

the _ quick _ brown _ fox _ jumped _ over _ the _ lazy _ dog

Final: 9 word tokens + 8 space tokens = 17 total (from 45 original characters)

Key:



Merging pairs



Final tokens

blue text = Merged tokens

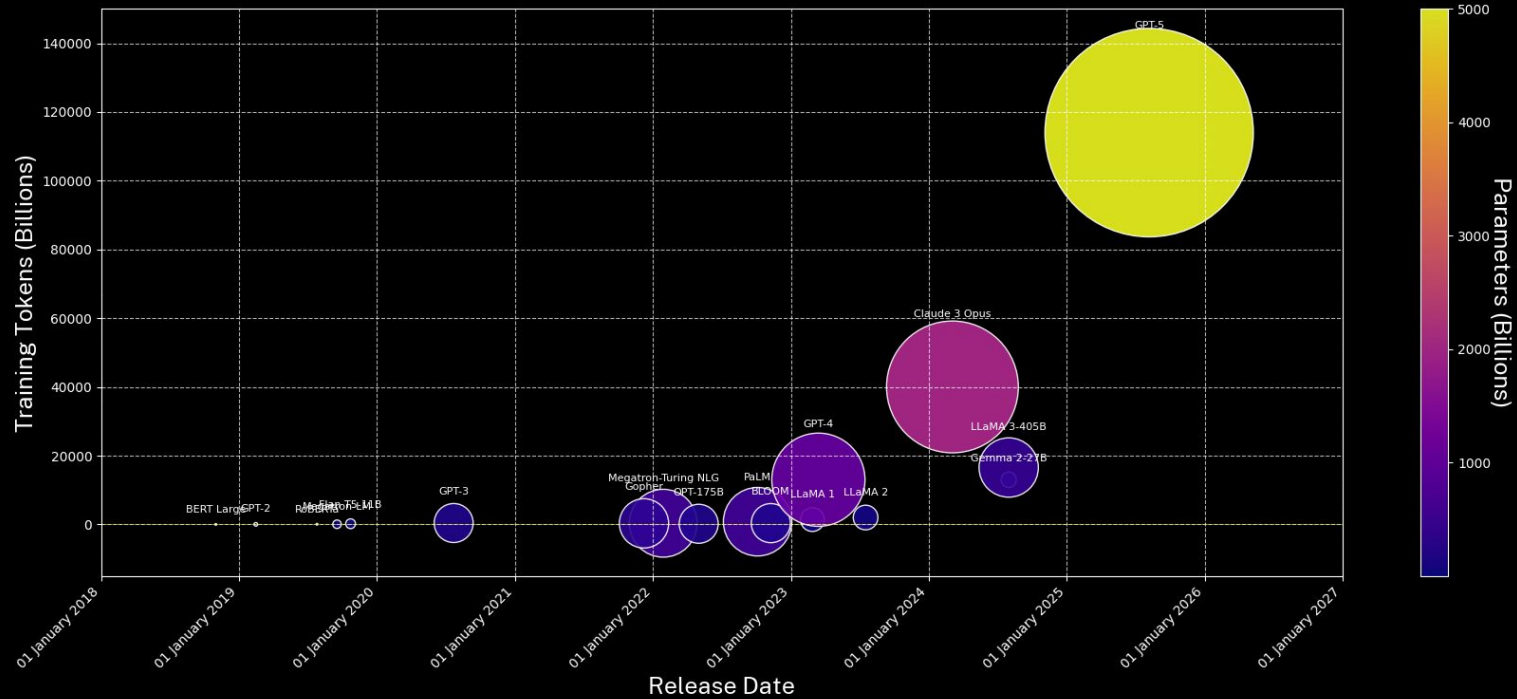
_ = space character

Relationships between tokens

the quick brown fox jumped over the lazy dog

	the	quick	brown	fox	jumped	over	the	lazy	dog
the	0.92	0.03	0.02	0.45	0.01	0.01	0.15	0.02	0.02
quick	0.05	0.85	0.15	0.42	0.02	0.01	0.02	0.03	0.02
brown	0.04	0.12	0.88	0.48	0.02	0.01	0.02	0.03	0.02
fox	0.05	0.10	0.12	0.90	0.45	0.15	0.02	0.03	0.02
jumped	0.02	0.02	0.02	0.40	0.92	0.35	0.02	0.03	0.02
over	0.02	0.02	0.02	0.10	0.38	0.88	0.02	0.03	0.02
the	0.02	0.02	0.02	0.05	0.02	0.02	0.90	0.40	0.45
lazy	0.02	0.02	0.02	0.05	0.02	0.02	0.05	0.88	0.42
dog	0.02	0.02	0.02	0.05	0.02	0.02	0.05	0.15	0.92

LLM Models: # of Training Tokens vs Release Year (Bubble size represents parameter count)



<https://github.com/KathyReid/token-wars-dataviz>

How many tokens are there in the world?

Cummins M.

[How much LLM training data is there, in the limit?](#)

Educating Silicon. May 9, 2024.

Accessed July 14, 2025.

60-
160
trillion tokens

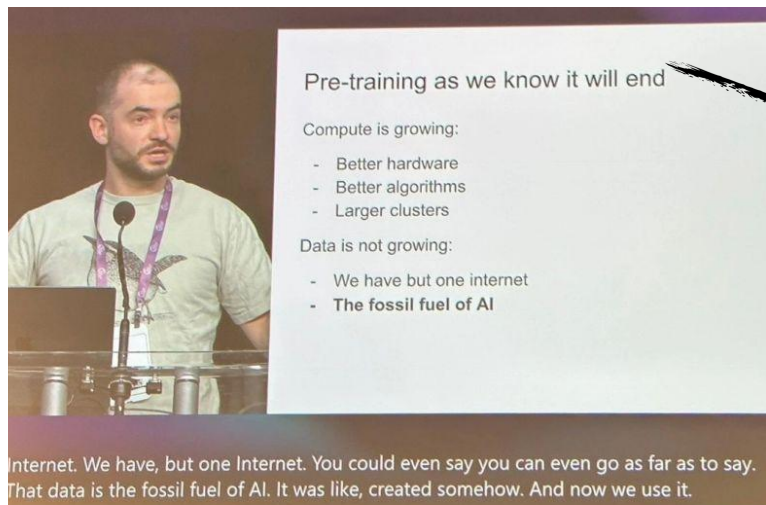
The Token Crisis

Attribution:

@JohnRushX via X/Twitter
taken at NeurIPS 2024,
Ilya Sutskever speaking

Video link:

<https://www.youtube.com/watch?v=1yvBqasHLZs>



“

Synthetic data is data that has been generated using a purpose-built mathematical model or algorithm, with the aim of solving a (set of) data science task(s).

Jordon, J., Szpruch, L., Houssiau, F., Bottarelli, M., Cherubin, G., Maple, C., ... & Weller, A. (2022). Synthetic Data--what, why and how?. [arXiv preprint arXiv:2205.03257](https://arxiv.org/abs/2205.03257).

Model Collapse



Shumailov, I., Shumaylov, Z., Zhao, Y., Papernot, N., Anderson, R., & Gal, Y. (2024). AI models collapse when trained on recursively generated data. *Nature*, 631(8022), 755-759.

“

Peak token is the point in time when the most human-generated, authentic data is available on the web.

Kathy Reid, 2025

Authentic, human-generated data is endangered



1

Scarce

models need
trillions of
tokens to train
on, none left

2

Rare

human-created
tokens are
polluted with AI
slop

3

Contested

because tokens
are endangered,
they're highly
sought after

02

Token harvesting

Token harvesting is extractive

Behind 'miracle' AI is an army of 'ghost workers' — and they're speaking out about Appen

By technology reporter Ariel Bogle

ABC Science AI

Fri 14 Oct 2022



Companies building artificial intelligence rely on an army of global workers. (Getty Images: The Good Brigade)

Australian authors challenge Productivity Commission's proposed copyright law exemption for AI

By Nicola Heath ABC Arts AI

Wed 13 Aug



Authors Danielle Clode, Rhett Davis and Kate Kruijink are concerned about the effects of proposed changes to copyright law. (ABC Arts: Christian Harimanow)

Small group of players shape AI innovation



+ GOOGLE + AI + BUSINESS

Google cut a deal with Reddit for AI training data

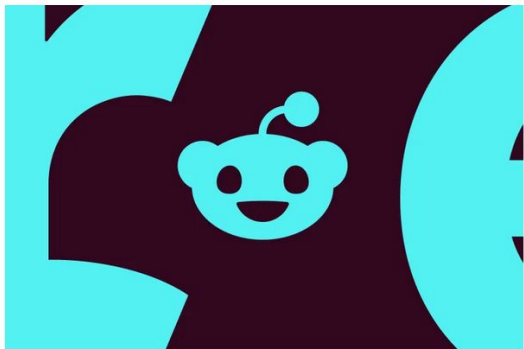


Image: The Verge

/ A deal reportedly worth \$60 million per year will give Google real-time access to Reddit's data and use Google AI for Reddit's search.

by + Emma Roth
Feb 23, 2024, 5:24 AM GMT-11



S&P Global

Press Releases

S&P Global and Anthropic Announce Integration of S&P Global's Trusted Financial Data into Claude



- New cutting-edge MCP server developed by Kensho, S&P Global's AI Innovation Hub, enables seamless access to S&P Global unrivaled datasets through Claude by Anthropic
- Integration expands how customers, from hedge fund managers to private equity analysts, can access S&P Global's data across the GenAI ecosystem

NEW YORK, July 15, 2025 /PRNewswire/ -- S&P Global (NYSE: SPGI) today announced a collaboration with Anthropic to bring S&P Global's trusted financial data into Claude by Anthropic. This integration enables financial professionals to answer complex financial questions and get fast, reliable answers grounded in trusted data from S&P Global via Claude.

Billions of people are not represented

WIRED

SECURITY POLITICS THE BIG STORY BUSINESS SCIENCE CULTURE REVIEWS

PARESH DAVE

BUSINESS MAY 31, 2023 7:00 AM

ChatGPT Is Cutting Non-English Languages Out of the AI Revolution

AI chatbots are less fluent in languages other than English, threatening to amplify existing bias in global commerce and innovation.



QUARTZ

Search Free Newsletters Editions

HOME LATEST BUSINESS NEWS MONEY & MARKETS TECH & INNOVATION A.I. LIFESTYLE LEADERSHIP EMAILS PODCASTS

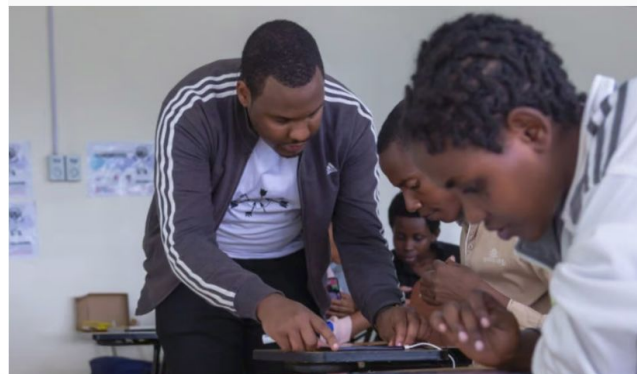
TECH & INNOVATION

Siri and Alexa still don't support African languages

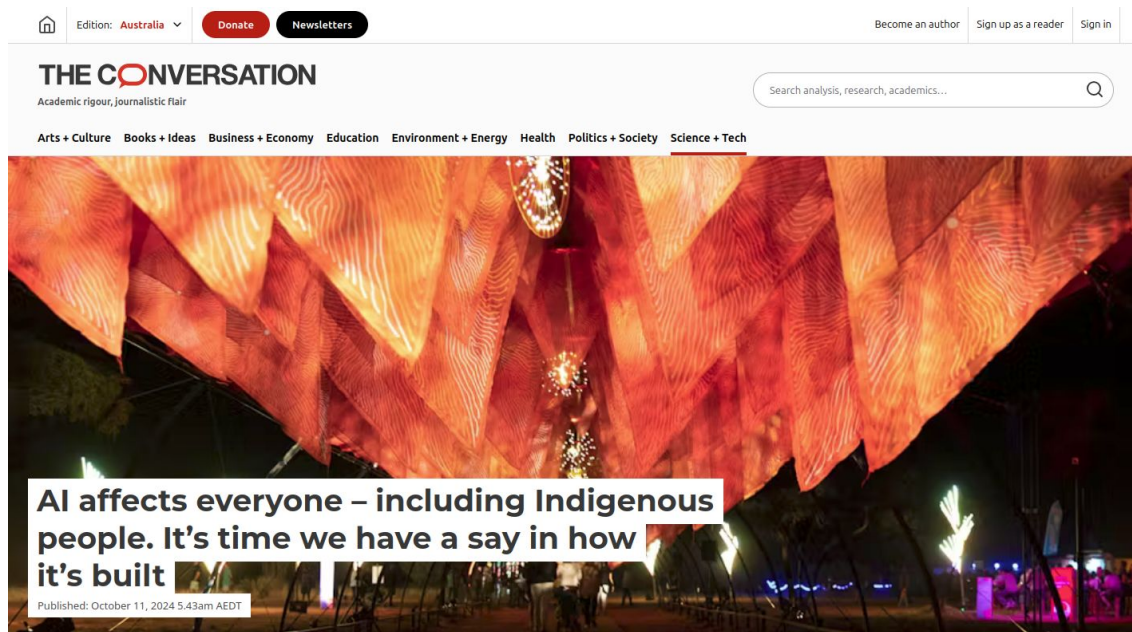
Despite speech increasingly becoming one of the main ways people interact with devices, voice technology remains largely closed off to Africa's languages, accents, and speech patterns. Case in point: The world's most popular voice assistants, Siri, Alexa and Google Assistant, still don't support any African languages. The continent has more than 1,000 languages.

By Carlos Mureithi and Carlos Mureithi Updated July 21, 2022

Twitter Facebook Messenger Email Print



Current approaches bulldoze rights



The 2022 Parrtjima Festival, held on Mparntwe (Alice Springs). Tamati Smith/Getty Images

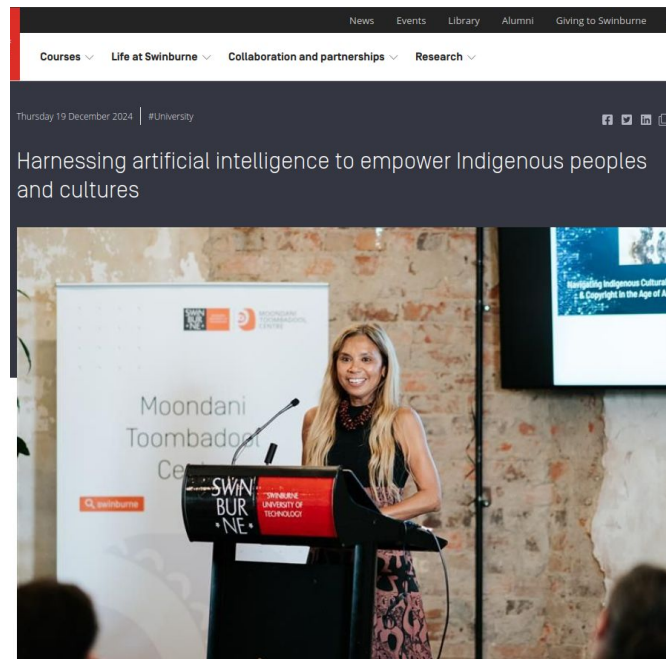


Since artificial intelligence (AI) became mainstream over the past two years, many of the risks it poses have been widely documented. As well as fuelling deep fake porn, threatening personal privacy and accelerating the climate crisis, some people believe the emerging technology could even lead to human extinction.

Author



Tamari Worrell
Senior Lecturer in the Department of Critical Indigenous Studies, Macquarie University



Dr Terri Janke, a Wuthathi, Yathaigana and Meriam woman, and an international authority on Indigenous cultural and intellectual property gave the 2024 Barak Wonga Oration

Home > News > 2024 > December > harnessing-artificial-intelligence-to-empower-indigenous-peoples-and-cultures

03

A better way: Mozilla Data Collective



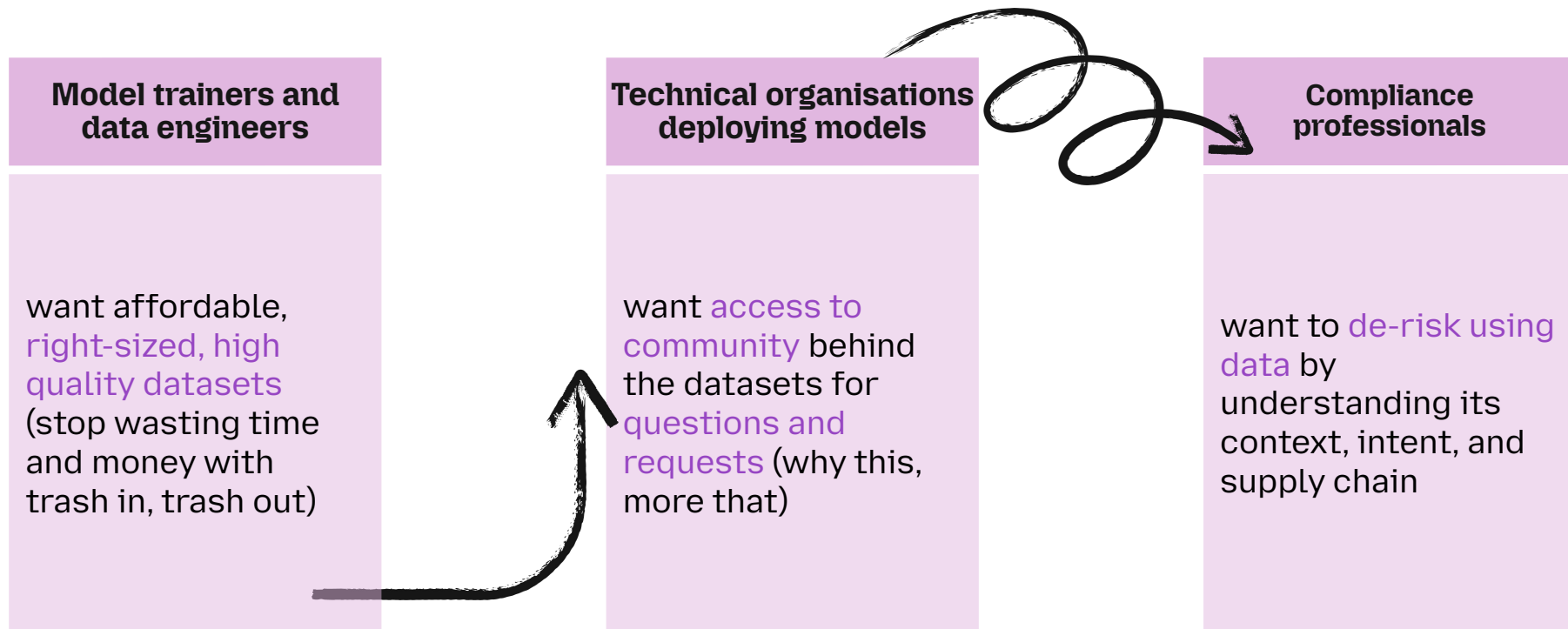
In an industry that relies on extractive practices, Mozilla Data Collective is **rebuilding the AI data ecosystem - with communities at the centre.**

Data contributors can share their data on their own terms; data users can access high quality data in a transparent and ethical way.

Data contributors



Data consumers



“

We are launching a
platform that unlocks **data**
abundance by giving **people**
and **communities** control
over their datasets.

Features for data contributors



1

Promote datasets

Store, host and list datasets with clear **stewardship** & robust **datasheet**

2

Governance dashboards

Share data in line with **your values and needs** - under particular licenses, or only for particular purposes

3

New value from existing data

By **opting in** to cross-walking, blending or augmenting their data with other datasets

Datasheets

common-voice/**cv-**
datasheets



13

Contributors

1

Issue

4

Stars

13

Forks



Features for data consumers



1

Discover datasets

With **transparent documentation** for licensing, intended task and source

2

Connect with community

to **contextualise** data and enable more **collaboration** and creation (data flywheels)

3

Supply chain transparency

enables **compliance in shifting legal landscape**

ISO 42001:

International standard for Artificial Intelligence Management System

(also Australian
standard)

[Standards](#)[Sectors](#)[About ISO](#)[Insights & news](#)[Taking part](#)[Store](#)[Read sample](#)

ISO/IEC 42001:2023

Information technology — Artificial intelligence —
Management system

Published (Edition 1, 2023)

What is ISO/IEC 42001?

ISO/IEC 42001 is an international standard that specifies requirements for establishing, implementing, maintaining, and continually improving an Artificial Intelligence Management System (AIMS) within organizations. It is designed for entities providing or utilizing AI-based products or services, ensuring responsible development and use of AI systems.

Takeaways



- 1 Authentic tokens are scarce, getting more so, and are highly sought after**
- 2 Token harvesting from the web is exploitative in many ways**

- 3 There's a better way to solve the needs of data producers & consumers**

Next steps



<https://discord.gg/4TjgEdq25Y>



mozdatacollective.bsky.social



mozilladatacollective@fosstodon.org

Come chat with us!
We run regular office hours

Send us feedback!

Tell us about a dataset!

mozilladatacollective@mozillafoundation.org

<https://datacollective.mozillafoundation.org/>

Thank you

Thank you to **EM Lewis-Jong**,
Jessica Rose and **Justin Grant**
for feedback, iteration and
assistance with this deck

04

Additional information

Further reading



- Associated Press. [AI startup Anthropic agrees to pay \\$1.5bn to settle book piracy lawsuit](#). The Guardian. September 5, 2025.
- Cummins M. [How much LLM training data is there, in the limit?](#) Educating Silicon. May 9, 2024.
- Global Indigenous Data Alliance. [CARE Principles for Indigenous Data Governance](#).
- Hao K. [Artificial intelligence is creating a new colonial world order](#). MIT Technology Review. April 19, 2022.
- Heath N. [Authors warn AI copyright exception a “free pass” for Big Tech to steal work](#). August 13, 2025.
- Jones PL, Mahelona K, Duncan S, Leoni G. Kaitiaki: closing the door on open Indigenous data. *International Journal on Digital Libraries*. 2025;26(1):1. doi:10.1007/s00799-025-00410-2
- Jordon, J., Szpruch, L., Houssiau, F., Bottarelli, M., Cherubin, G., Maple, C., ... & Weller, A. (2022). Synthetic Data--what, why and how?. [arXiv preprint arXiv:2205.03257](#).
- OpenAI. [OpenAI and Reddit Partnership](#). May 16, 2024.
- Shumailov, I., Shumaylov, Z., Zhao, Y., Papernot, N., Anderson, R., & Gal, Y. (2024). AI models collapse when trained on recursively generated data. *Nature*, 631(8022), 755-759.