# AI in our Evolving Knowledge Infrastructure

David De Roure

What I'm going to talk about today:

- Knowledge Infrastructure

- Social Machines

- Scholarly Primitives

- Software

- Research Integrity

- Music

What I'm not going to focus on today:

- Computational infrastructure

- Bias

- Responsible innovation

- Ethics

- Energy

# Acknowledgements

- Pip Willcox "Scholarly Social Machines: A Web Science Perspective on our Knowledge Infrastructure" ACM WebSci 2020

- Puyu Wang "Leveraging Ontologies in Standard Language Models for Research Capabilities: An Evaluation of Performance" ESWC 2043
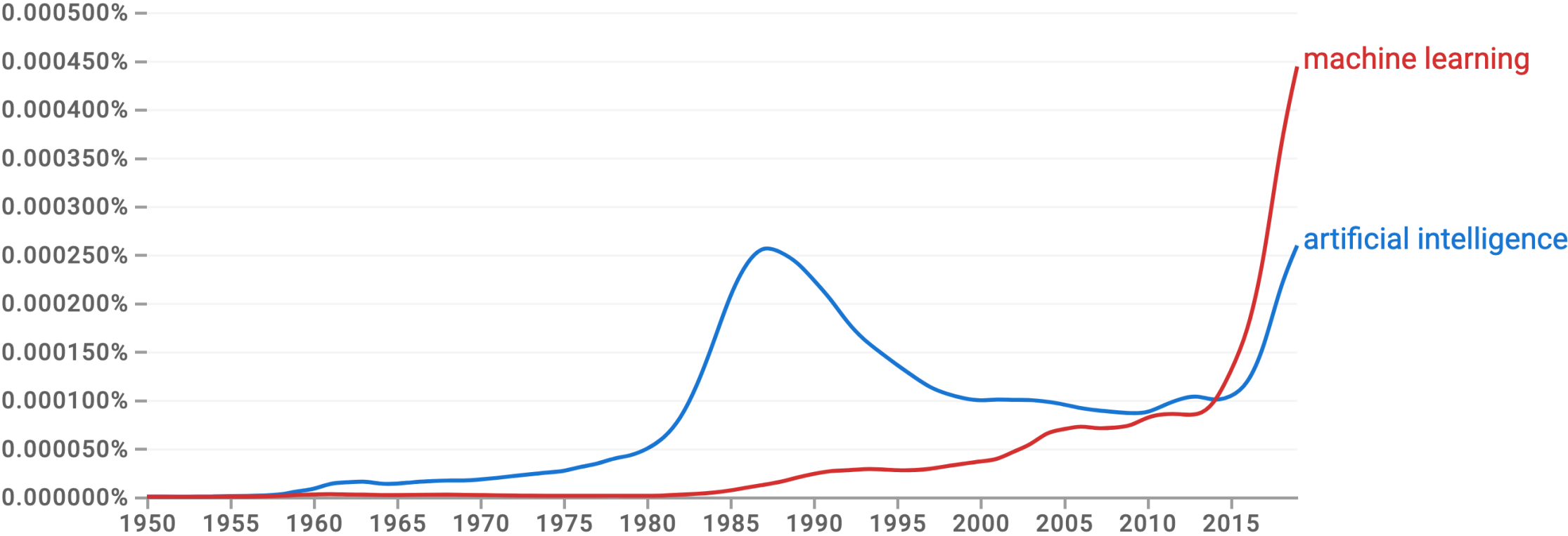
- Rob Walton

- Oktay Cetinkaya & Peter Novitzky

# Google Books Ngram Viewer

artificial intelligence,machine learning

1950 - 2019 ▾  English (2019) ▾  Case-Insensitive  Smoothing of 1 ▾



(click on line/label for focus)

\* (load "eliza.lisp")

FILE CREATED 08/22/68 1522:26

FILE CREATED 13-JUN-72 4:20:07

FILE CREATED 22-APR-72 23:26:05

TELL ME YOUR PROBLEMS. PLEASE TERMINATE INPUT WITH A PERIOD OR A QUESTION MARK.

\*\*I am giving a talk on evolving knowledge infrastructure.

DO YOU BELIEVE IT NORMAL TO BE GIVING A TALK ON EVOLVING KNOWLEDGE INFRASTRUCTURE?

\*It is normal for me.

SUPPOSE IT WERE NOT NORMAL FOR YOU.

\*Would you like to give the talk instead?

OH? I LIKE TO GIVE THE TALK INSTEAD?

\*Yes Eliza.

WHY ARE YOU SO SURE?

\*You are the original chatbot.

WHAT MAKES YOU THINK I AM THE ORIGINAL CHATBOT?

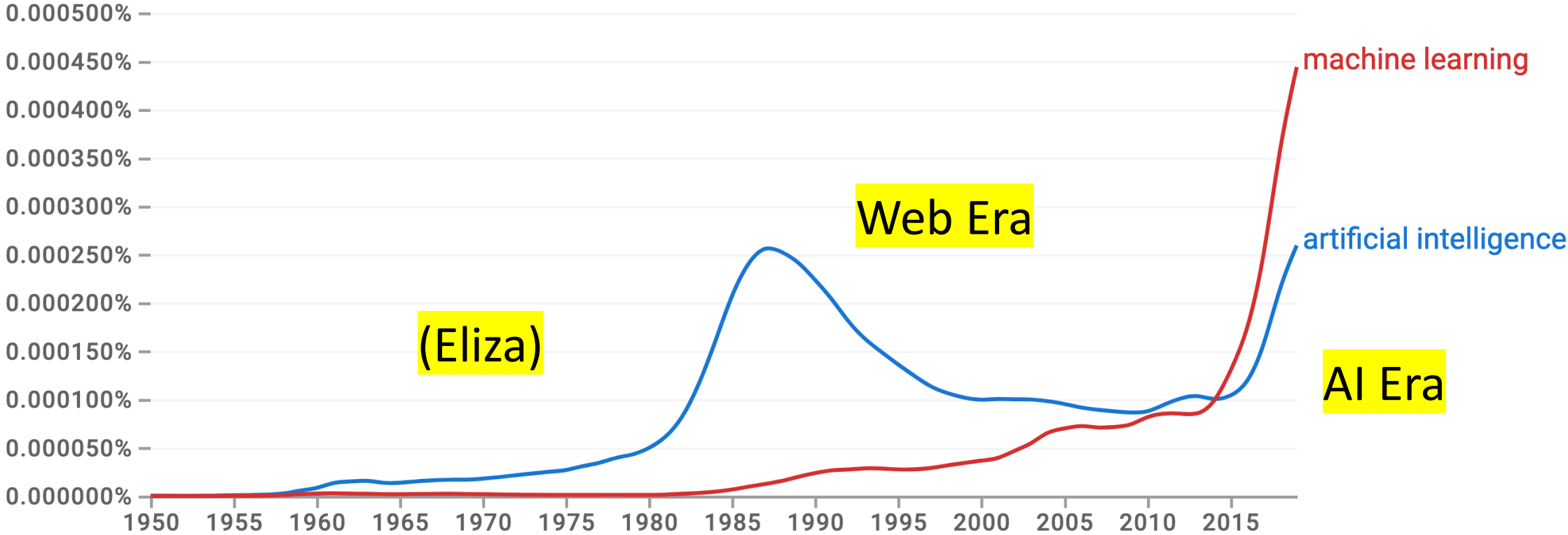# Google Books Ngram Viewer

artificial intelligence,machine learning

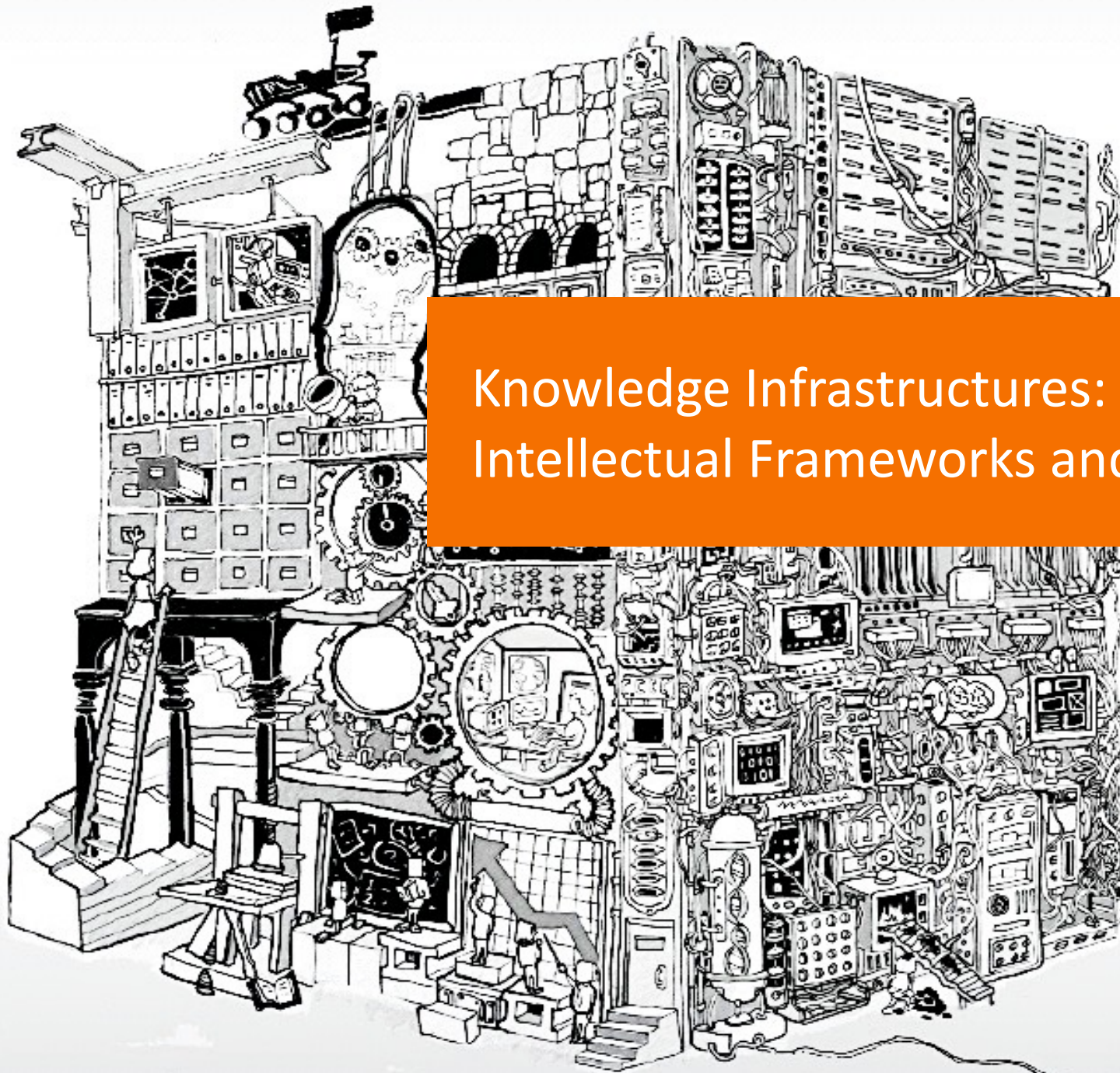1950 - 2019    English (2019)    Case-Insensitive    Smoothing of 1



(click on line/label for focus)

Some eras:

• Microcomputers

• Internet

• Web (and Semantic Web)

• AI/ML

NB I've characterised these by "democratization"

Knowledge Infrastructures:
Intellectual Frameworks and Research Challenges

# FORCE11 Manifesto

## On this page:

**Improving Future Research Communication and e-Scholarship**

*Editors*: Phil E. Bourne[a], Tim Clark[b], Robert Dale[c], Anita de Waard[d], Ivan Herman[e], Eduard Hovy[f], and David Shotton[g]

*Contributors*: Bradley P. Allen[d], Aliaksandr Birukou[h], Judith A. Blake[i], Philip E. Bourne[a], Simon Buckingham Shum[j], Gully A.P.C. Burns[f], Leslie Chan[k], Olga Chiarcos[l], Paolo Ciccarese[b], Tim Clark[b], Laura Czerniewicz[m], Robert Dale[c], Anna De Liddo[j], David De Roure[g], Anita de Waard[d], Stefan Decker[n], Alex Garcia Castro[o], Carole Goble[p], Eve Gray[m], Paul Groth[q], Udo Hahn[r], Ivan Herman[e], Eduard H. Hovy[f], Michael J. Kurtz[s], Fiona Murphy[t], Cameron Neylon[u], Steve Pettifer[p], Mike W. Rogers[v], David S. H. Rosenthal[w], David Shotton[g], Jarkko Siren[v], Herbert van de Sompel[x], Peter van den Besselaar[q] and Todd Vision[y]

*Affiliations*: (a) University of California at San Diego; (b) Harvard Medical School; (c) Macquarie University; (d) Elsevier Laboratories; (e) Centrum voor Wiskunde en Informatica, Amsterdam; (f) University of Southern California; (g) University of Oxford; (h) CREATE-NET; (i) The Jackson Laboratory; (j) The Open University; (k) University of Toronto; (l) Springer-Verlag; (m) University of Cape Town; (n) National University of Ireland, Galway; (o) Universität Bremen; (p) University of Manchester; (q) Vrije Universiteit Amsterdam; (r) Universität Jena;(s) Harvard-Smithsonian Center for Astrophysics; (t) Wiley-Blackwell; (u) Rutherford Appleton Laboratory; (v) European Commission Brussels; (w) Stanford University; (x) Los Alamos National Laboratory; (y) University of North Carolina at Chapel Hill

2011-10-28

**my experiment**

www.myexperiment.org

- A workflow commons for workflow sharing, designed using Web 2.0 principles

- Launched open beta in November 2007, still actively used

- Largest public collection of workflows (~4000), for multiple workflow systems

- 1400+ entries in Google Scholar refer to myexperiment.org

- Open source, REST API, part of Open Linked Data cloud (66k triples) - lod-cloud.net

- Introduced "packs" which led to Research Objects – www.researchobject.org

- Workflow collection studied in scientific workflow and e-Science communities

De Roure, D., Goble, C. Stevens, R. (2009) The Design and Realisation of the myExperiment Virtual Research Environment for Social Sharing of Workflows. Future Generation Computer Systems 25, pp. 561-7.

Carole Goble

# Social Machines

"Real life is and must be full of all kinds of social constraint – the very processes from which society arises. Computers can help if we use them to create **abstract social machines on the Web: processes in which the people do the creative work and the machine does the administration...** The stage is set for an evolutionary growth of new social engines. **The ability to create new forms of social process would be given to the world at large, and development would be rapid."**

Berners-Lee, *Weaving the Web*, 1999 (pp. 172–175)

## Featured Story

# The Decline of Wikipedia

By Tom Simonite on October 22, 2013

294 COMMENTS

"Yet Wikipedia and its stated ambition to "compile the sum of all human knowledge" are in trouble. The volunteer workforce that built the project's flagship, the English-language Wikipedia—and must defend it against vandalism, hoaxes, and manipulation—has shrunk by more than a third since 2007 and is still shrinking...

The main source of those problems is not mysterious. The loose collective running the site today, estimated to be 90 percent male, operates a crushing bureaucracy with an often abrasive atmosphere that deters newcomers who might increase participation in Wikipedia and broaden its coverage..."

2013

# Criticism of Wikipedia

The free online encyclopedia Wikipedia has been criticized since its creation in 2001. Most of the criticism has been directed toward its content, community of established volunteer users, process, and rules. Critics have questioned its factual reliability, the readability and organization of its articles, the lack of methodical fact-checking, and its political bias. Concerns have also been raised about systemic bias along gender, racial, political, corporate, institutional, and national lines. Conflicts of interest arising from corporate campaigns to influence content have also been highlighted. Further concerns include the vandalism and partisanship facilitated by anonymous editing, clique behavior (from contributors as well as administrators and other top figures), social stratification between a guardian class and newer users, excessive rule-making, edit warring, and uneven policy application.

THE ZOONIVERSE WORKS

# 799,952,637

CLASSIFICATIONS SO FAR BY
2,716,134 REGISTERED VOLUNTEERS

WILDWATCH KENYA

WEDDELL SEAL COUNT

WildCam Darién
WILDCAM DARIÉN

Colorado Corridors Project
COLORADO CORRIDORS PROJECT

NYBG
NOTES FROM NATURE - NYBG

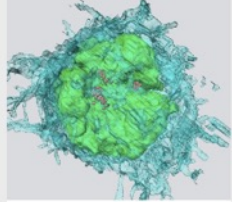Measuring evolution in action
SQUIRRELMAPPER

SKINK SPOTTER NZ

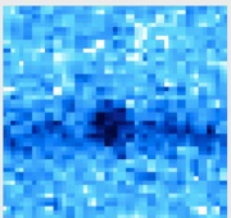WESTERN MONTANA WILDLIFE

BASH THE BUG

ETCH A CELL

K2
EXOPLANET EXP

FOSSIL ATMOSPHERES

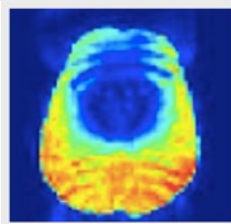GALAXY NURSERIES

SkyMapper
SUPERNOVA SIGHTING

Amazon Cam
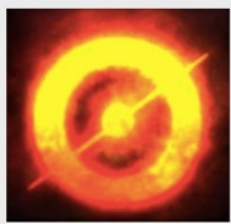AMAZONCAM TAMBOPATA

VERITAS
MUON HUNTER

BACKYARD WO PLANET S

BRAIN MATCH

SOLAR STORMWATCH II

PLANETARY RESPONSE NETWORK & RESCUE GLOBAL - CARIBBEAN

HUBBLE'S HOT STARS

SNOW SPOTTER

toledo WIL SHO
TOLEDO ZOO WIL

LEAF TARGETING

STEELPAN VIBRATIONS

SCRIBES OF THE CAIRO GENIZA
SCRIBES OF THE CAIRO GENIZA

MUTUAL MUSES

CRIMINAL CHARACTERS

WORLDS OF WONDER

AFRICAN AMERICAN CIVIL WAR SOLDIERS

Discover wildlife in
FOCUS ON WILDLIFE -- CLEVELAND METROPARKS

CHICAGO WILDLIFE WATCH

ARIZONA BATWATCH

GRAVITY S

ELEPHANT EXPEDITION

SERENGETI WILDEBEEST COUNT

Michigan ZoomIN
MICHIGAN ZOOMIN

TPT
THE PLASTIC TIDE

GALAXY ZOO: CLUMP SCOUT

GALAXY ZOO MOBILE

EUCLID - CHALLENGE THE MACHINES

MAPPING CHANGE

MILKY WAY PROJECT

RADIO METEOR ZOO

CAMERA CATAL

*Scholarly Primitives* refer to some basic functions common to scholarly activity across disciplines, over time, and independent of theoretical orientation.
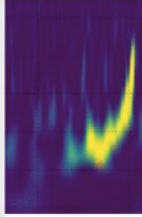
These 'self-understood' functions form the basis for higher-level scholarly projects, arguments, statements, interpretations—in terms of our original, mathematical/philosophical analogy, axioms.

John Unsworth

Discovering

Annotating

Comparing

Referring

Sampling

Illustrating

Representing

- What is the nature of truth: independent or cumulative?

- How much play is there in the machine? What is the degree of co-creation?

- Today's new Scholarly Primitives reflect the stance of creator as well as user: researchers and citizens are empowered to create Social Machines.

Pip Willcox

Enter AI/ML…

Galaxy Zoo

Few have witnessed what you're about to see

"In the decade the project has been running, Galaxy Zoo volunteers have helped understand the Universe and made spectacular discoveries. We hope you'll join us for the next stage of the adventure."

In an effort to speed up classifications to cope with the large number of galaxies we expect to receive from new surveys, we've been working on ways to combine your classifications with those of machines, inspired by the idea that the combination of both automatic and human classification may be more powerful than either alone. If you choose the 'Enhanced' work flow, you will be much more likely to see the top 100 galaxies our galaxy-classifying robot thinks it needs help with in order to improve. All galaxies will be seen by at least a few volunteers to make sure we aren't missing anything. If you'd rather just see a random selection of available galaxies, choose 'Classic'.

# Crowdsourcing *vs* AI-sourcing

Our scholarly workflows now engage the crowd and the AIs, with increasing ease.

With both crowdsourcing and AI, we are outsourcing work (typically in order to scale).

What are the hidden costs of these approaches?

AEOLIAN Network
Artificial Intelligence for Cultural Organisations

**What is Retrieval-Augmented Generation?**

Why is Retrieval-Augmented Generation important?

What are the benefits of Retrieval-Augmented Generation?

How does Retrieval-Augmented Generation work?

What is the difference between Retrieval-Augmented Generation and

# What is Retrieval-Augmented Generation?

Retrieval-Augmented Generation (RAG) is the process of optimizing the output of a large language model, so it references an authoritative knowledge base outside of its training data sources before generating a response. Large Language Models (LLMs) are trained on vast volumes of data and use billions of parameters to generate original output for tasks like answering questions, translating languages, and completing sentences. RAG extends the already powerful capabilities of LLMs to specific domains or an organization's internal knowledge base, all without the need to retrain the model. It is a cost-effective approach to improving LLM output so it remains relevant, accurate, and useful in various contexts.

# BUSINESS INSIDER

This is some heavy-duty philosophy. Language isn't just communication; it's a substrate for thought and a carrier wave for culture. And Underwood thinks LLMs are tuned into that wave. For scholars, the point of reading, of writing, of studying language isn't just to generate an essay or critique a poem, but to figure out *what* we think, and how best to express it. Underwood hopes that LLMs, with their sophisticated statistical models of language, might help us uncover new insights into our own minds. In his view, they do more than simply parrot sources stochastically, not because they're "thinking" — but because we're there to listen to them.

# SOFTWARE

AI is now routinely used by software developers, for writing code, unit tests, debugging, documentation.

It is surprisingly good with ontologies too.

# Set Up Now for AI to Augment Software Development

- AI can automate much of today's mundane software development tasks.
- It can also help software engineers produce code more quickly.

**ChatGPT**

Certainly! Here's a short piece of music in ABC notation for your introduction:

```makefile
X:1

T:AI and the Evolution of Knowledge Infrastructure

M:4/4

L:1/4

K:C

CDEG| E2G2| A2c2| G2E2| C2D2| E2G2| c2e2| d4|]
```

AI and the Evolution of Knowledge Infrastructure

- **Peer-review process:** AI and in particular LLMs used in content creation, by editors and peer-reviewers. Exacerbated by anonymity of peer-reviewers and the 'publish or perish' ecosystem.

- **Data analytics and visualisation**: increasingly relying on data analytics using ML models. Hallucination and other non-deterministic processes may skew the presentation of final results and/or visualisations.

- **Models** trained on state-of-the-art datasets that might fall under IP protection laws. Revelation of state-of-the-art, non-verified unique datasets and results might compromise ('poison') other datasets.

- Such events might exacerbate problems in domains such as pharmacological research, defence and national security research, critical national (research) infrastructure, etc.

Oktay Cetinkaya & Peter Novitzky - Synthesis Fellows

# The ecosystem perspective

If we are measured by our rate and quality/impact of publication, the ecosystem will adapt to support this.

Automation is already happening pervasively in the scholarly workflows.

**Google Scholar**

"As of my last knowledge update" -LLM -chatGPT -artificial 🔍

⦿ Articles   ○ Case law

**Google Scholar**

"as an AI language model" -LLM -chatGPT -artificial 🔍

⦿ Articles   ○ Case law

Advancing purpose and mutual knowledge of intentions

Feedback when strands perform activities together

Aggregates of strands

Structures of composed strands

Situations of shared experience

Strands driven by purpose

Situations of activity

A social machine formalism suitable for representing
sustained organisation on and off the Web

Rob Walton

Alter | PRiSM led by Robert Laidlow

Alter was premiered by the Britten Sinfonia and Marta Fontanals-Simmons, conducted by William Cole, at the Barbican.
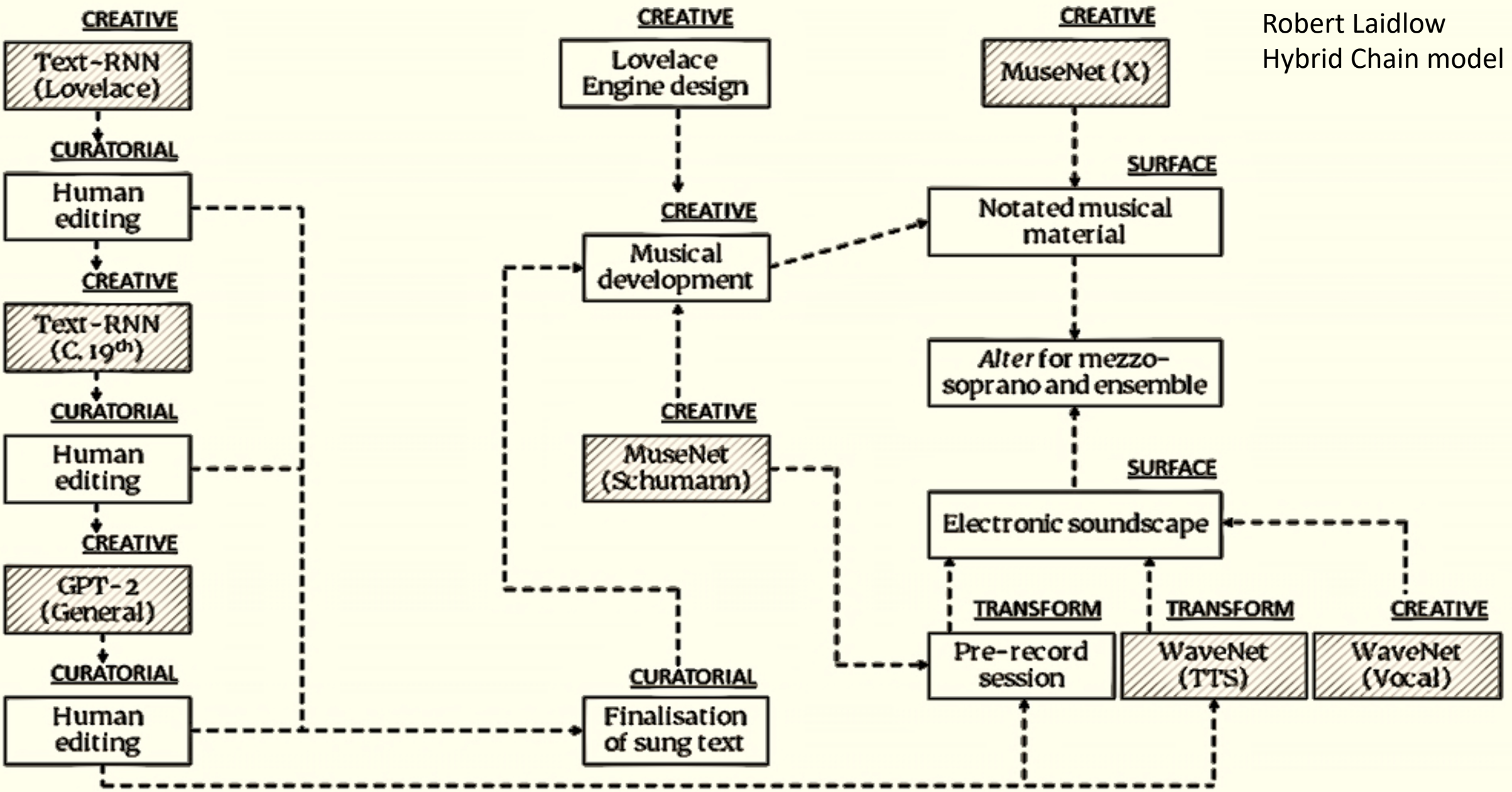
PRiSM

The Alan Turing Institute

https://www.youtube.com/watch?v=L1mQGaNmfUM

Robert Laidlow
Hybrid Chain model

## Future Epoch

Solar flares smash into space and time

No response

A new existence forged from the fragments

Of old shapes

Making shapes with the raw materials

In the hands of the painters

New plans

Borrowing the lines from the past

They could learn from where they fell

**MUSIC OF G.E. LEWIS**

**Didyeridea** (improvisation)
Douglas Ewart, didjeridu
G.E. Lewis, trombone

**Canyon** (1989)
Douglas Ewart, clarinet
G.E. Lewis, live electronics

**Eighteen Summers** (improvisation)
Douglas Ewart, alto saxophone
G.E. Lewis, trombone

• • • intermission• • •

**Voyager, first part** (1986-87)
Douglas Ewart, alto saxophone
Voyager (algorithmic computer performer)

**Eighteen Winters** (improvisation)
Douglas Ewart, clarinet
G.E. Lewis, trombone

**Voyager, second part** (1986-87)
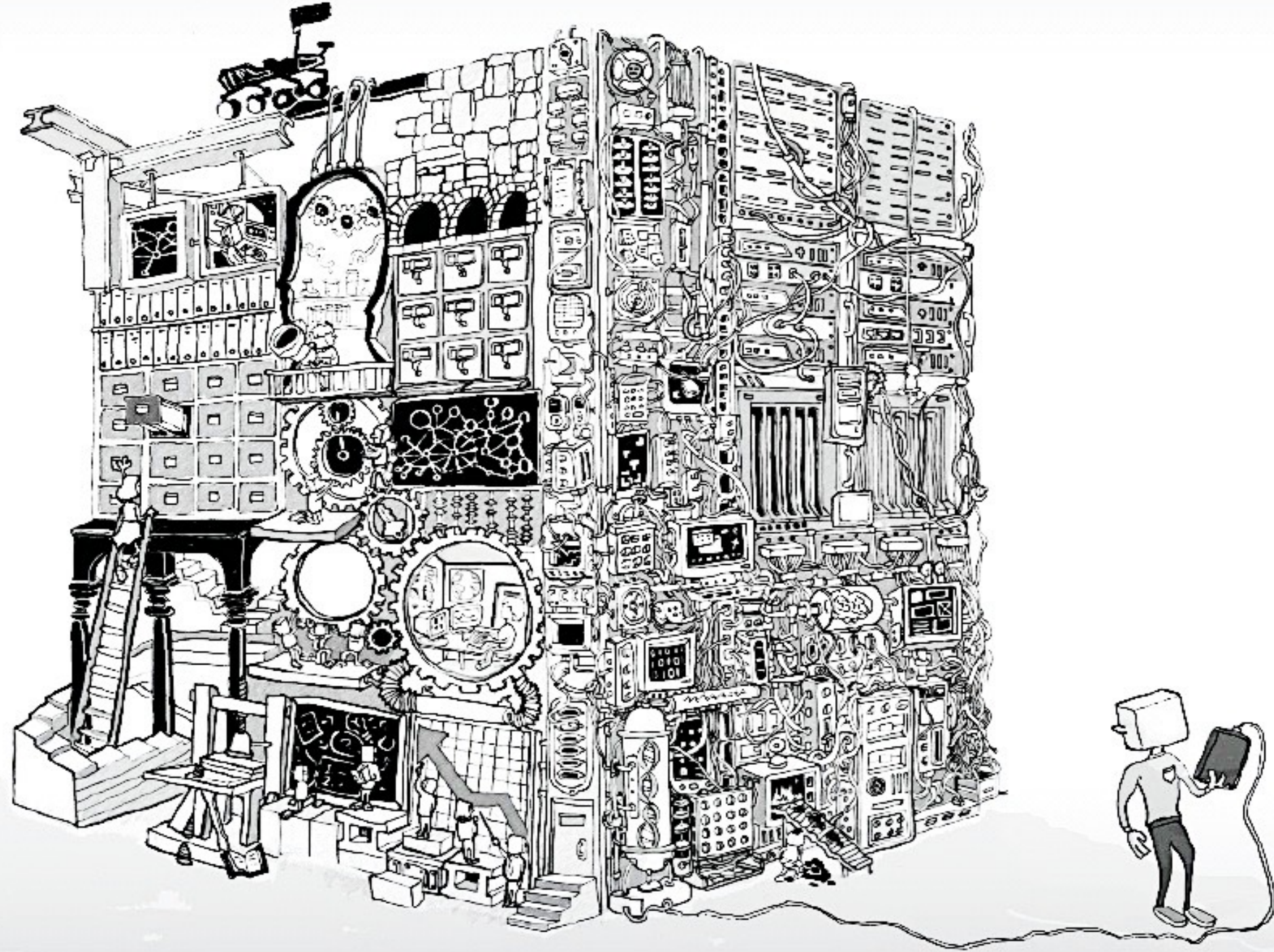G.E. Lewis, trombone
Voyager (algorithmic computer performer)

"Thus, it has come to pass not only that improvisations by creative musical machines are often indistinguishable from those created by humans, but also that they need not be so distinguished."

# Engineering Manager - Human Data

## About the Team

OpenAI's Human Data Team delivers human data across all projects at OpenAI, including ChatGPT, plugins, browsing, code interpretation, GPT-V, and more. This human data transforms pretrained GPT models into models that follow human instructions and in other words align with human intent. Our platform operates across various modalities (text, images, code, plugins, etc.) and domains of expertise (math, programming, and more).

We design, develop, and maintain the production-quality platform necessary to generate such data at scale. The team is responsible for data management tools, operations, data quality, and research on techniques for data collection. Our data collection aims to leverage AI models in every aspect of the data collection process to maximize the value of human input (e.g., active learning, synthetic data, and more).

# CLOSING QUESTIONS

1. Should we try to track provenance more comprehensively? Or do we not need to distinguish human and AI anymore?

2. What are the scholarly primitives in the new era?

3. How do we measure the rate of innovation - and is it increasing?