

# Is AI the future of health and social science?

The poster is split into two main sections by a large yellow lightning bolt graphic. The left section has a dark teal background and the right section has an orange background. Text and images are arranged to provide details about the event, including the title, date, time, location, and the names of the debaters.

**IS AI THE FUTURE OF HEALTH AND SOCIAL SCIENCE? A DEBATE**

**15 DEC 2025**  
5.30pm – 8.00pm  
PUNNET HALL, IOE  
20 BEDFORD WAY, LONDON

**FOR**



**David Bann**  
UCL

**AGAINST**



**Peter Tennant**  
University of Leeds

SPONSOR



**NCRM**  
NATIONAL  
CENTRE FOR  
RESEARCH  
METHODS

**“The aim of argument... should not be victory, but progress” Joubert, 1848**

# Argument

## 1. Quantitative research is mostly cognitive work

- Humans...
- AI...

## 2. If AI helps research, it's 'the future'

- Humans + AI

# Argument

## 1. Quantitative research is mostly cognitive work

- Humans...
- AI...

## 2. If AI helps research, it's 'the future'

- Humans + AI

# Tasks in quantitative health / social science research

- **Designing studies + collecting data**
  - Survey work (lit review, design, pilot, fieldwork)
- **Papers:**
  - Literature review
  - Forming research questions/hypothesis
  - Data analysis
  - Interpretation / write-up

# Tasks in quantitative health / social science research

## Cognitive tasks

- **Designing studies + collecting data**
  - Survey work (**lit review, design**, pilot, fieldwork)
- **Papers:**
  - **Literature review**
  - **Forming research questions/hypothesis**
  - **Data analysis**
  - **Interpretation / write-up**

# GPTs are GPTs: Labor market impact potential of LLMs

Research is needed to estimate how jobs may be affected

By Tyna Eloundou<sup>1</sup>, Sam Manning<sup>2</sup>, Pamela Mishkin<sup>1</sup>, Daniel Rock<sup>3</sup>

**W**e propose a framework for evaluating the potential impacts of large-language models (LLMs) and associated technologies on work by considering their relevance to the tasks workers perform in their jobs. By applying this framework (with both humans and using an LLM), we estimate that roughly 1.8% of jobs could have over half their tasks affected by LLMs with simple in-

technology discussions, focusing more narrowly on advanced software capabilities than on the potential for business process reengineering, new intangible assets creation, or workforce retraining. General-purpose technologies such as electricity or computing historically have had far-reaching effects that took decades to fully materialize. With evidence of the general-purpose technology potential of LLMs, we urge caution in making long-term predictions while offering an outline of where work might change.

Supplement, p41

Group	Occupations with highest exposure	% Exposure
<b>Human <i>E1</i></b>	Interpreters and Translators	76.5
	<b>Survey Researchers</b>	<b>75.0</b>
	Poets, Lyricists and Creative writers	68.8
	Animal Scientists	66.7
	Public Relations Specialists	66.7
<b>Human <i>E1</i> + 0.5 * <i>E2</i></b>	<b>Survey Researchers</b>	<b>84.4</b>
	Writers and Authors	82.5
	Interpreters and Translators	82.4
	Public Relations Specialists	80.6
	Animal Scientists	77.8

exposure=task completion by at least 50%

# Exposed != imminently extinct

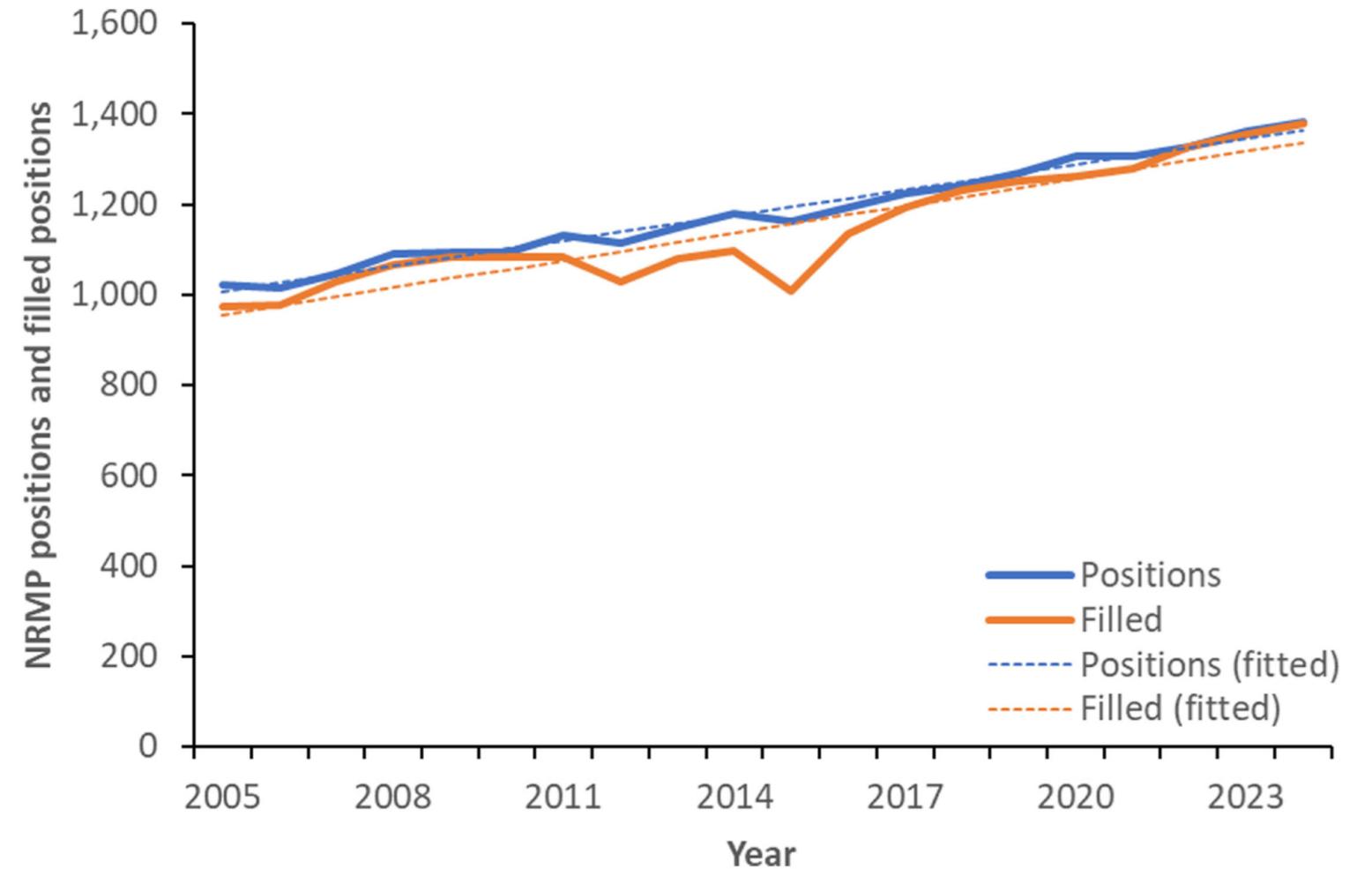
NY Times revisits Nobel Prize winner's prediction AI will render radiologists obsolete

Marty Stempniak | May 15, 2025 | Radiology Business | Artificial Intelligence



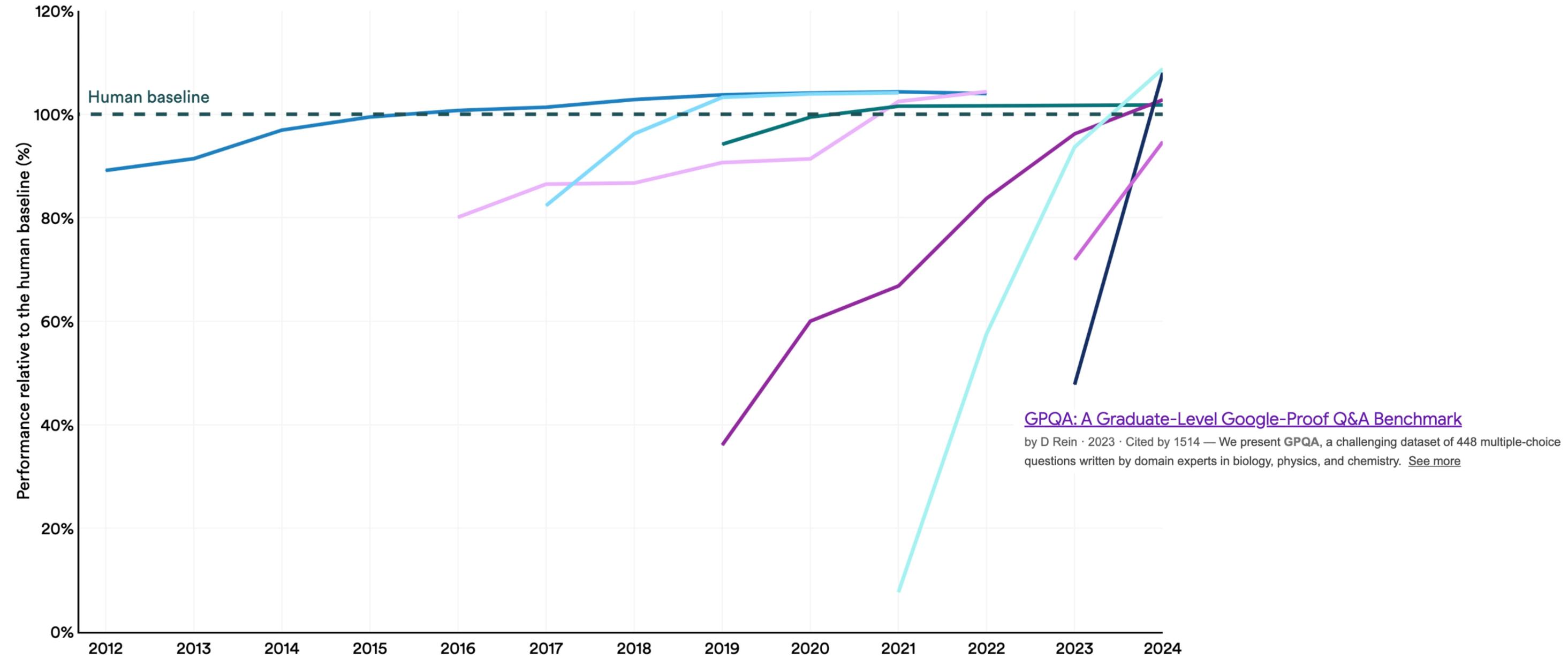
“People should stop training radiologists now...”

Hinton 2016



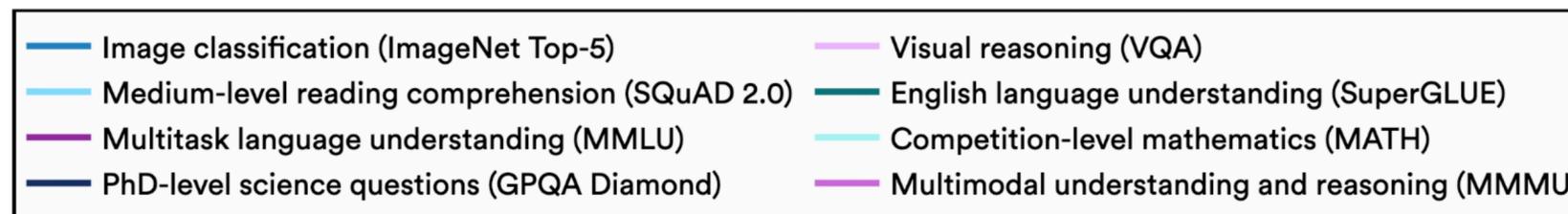
# Select AI Index technical performance benchmarks vs. human performance

Source: AI Index, 2025 | Chart: 2025 AI Index report



[GPQA: A Graduate-Level Google-Proof Q&A Benchmark](#)

by D Rein · 2023 · Cited by 1514 — We present GPQA, a challenging dataset of 448 multiple-choice questions written by domain experts in biology, physics, and chemistry. [See more](#)

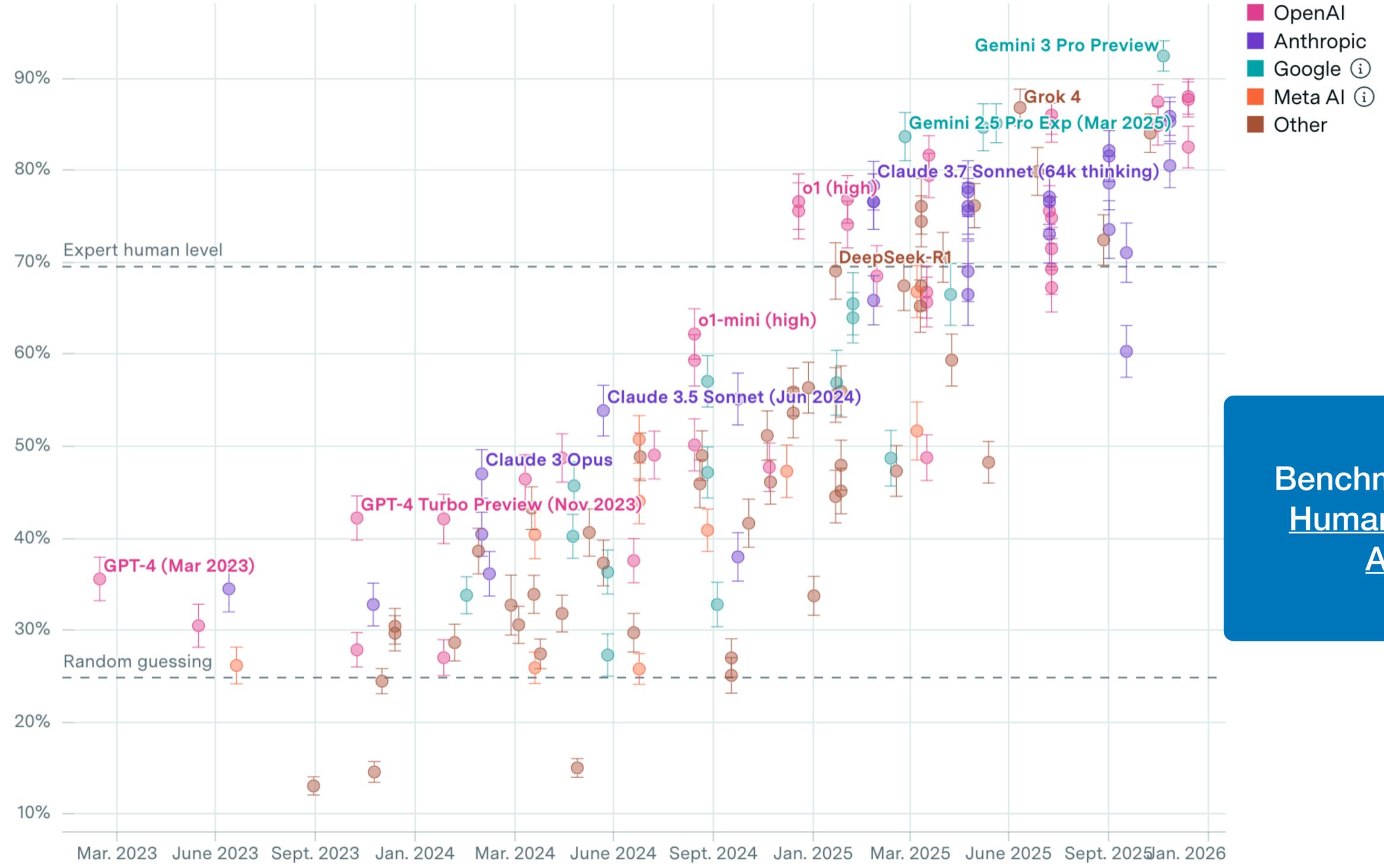


# AI performance on a set of Ph.D.-level science questions

Graph Settings

GPQA Diamond accuracy ⓘ

142 Results Organization



Benchmark saturation?  
Humanity's last exam  
ARC-AGI...

# GPT 5.2 - 3 days ago, ~30 day since 5.1 (!)



OpenAI

Run with maximum available reasoning effort.

	GPT-5.2 Thinking	GPT-5.1 Thinking
<b>SWE-Bench Pro</b> Software engineering	• 55.6%	50.8%
<b>GPQA Diamond</b> Science questions (No tools)	• 92.4%	88.1%
<b>CharXiv Reasoning</b> Scientific figure questions (No tools)	• 82.1%	67.0%
<b>FrontierMath</b> Advanced mathematics (Tier 1 - 3, Tier 4)	• 40.3% 14.6%	31.0% 12.5%
<b>AIME 2025</b> Competition math (No tools)	• 100.0%	94.0%
<b>ARC-AGI-1</b> Abstract reasoning	• 86.2%	72.8%
<b>ARC-AGI-2</b> Abstract reasoning	• 52.9%	17.6%
<b>GDPval</b> Knowledge work tasks	• 70.9%	38.8% <small>GPT-5</small>

# AI: powerful, yet (for now) 'jagged'

## DeepSeek's self-correcting AI model aces tough maths proofs

The model, DeepSeekMath-V2, scored 118 out of 120 points on questions from the 2024 William Lowell Putnam Mathematical Competition, beating the top human score of 90. The model also performed at the level of gold-medal winners in the International Mathematical Olympiad (IMO) 2025 and the 2024 China Mathematical Olympiad. The results are described

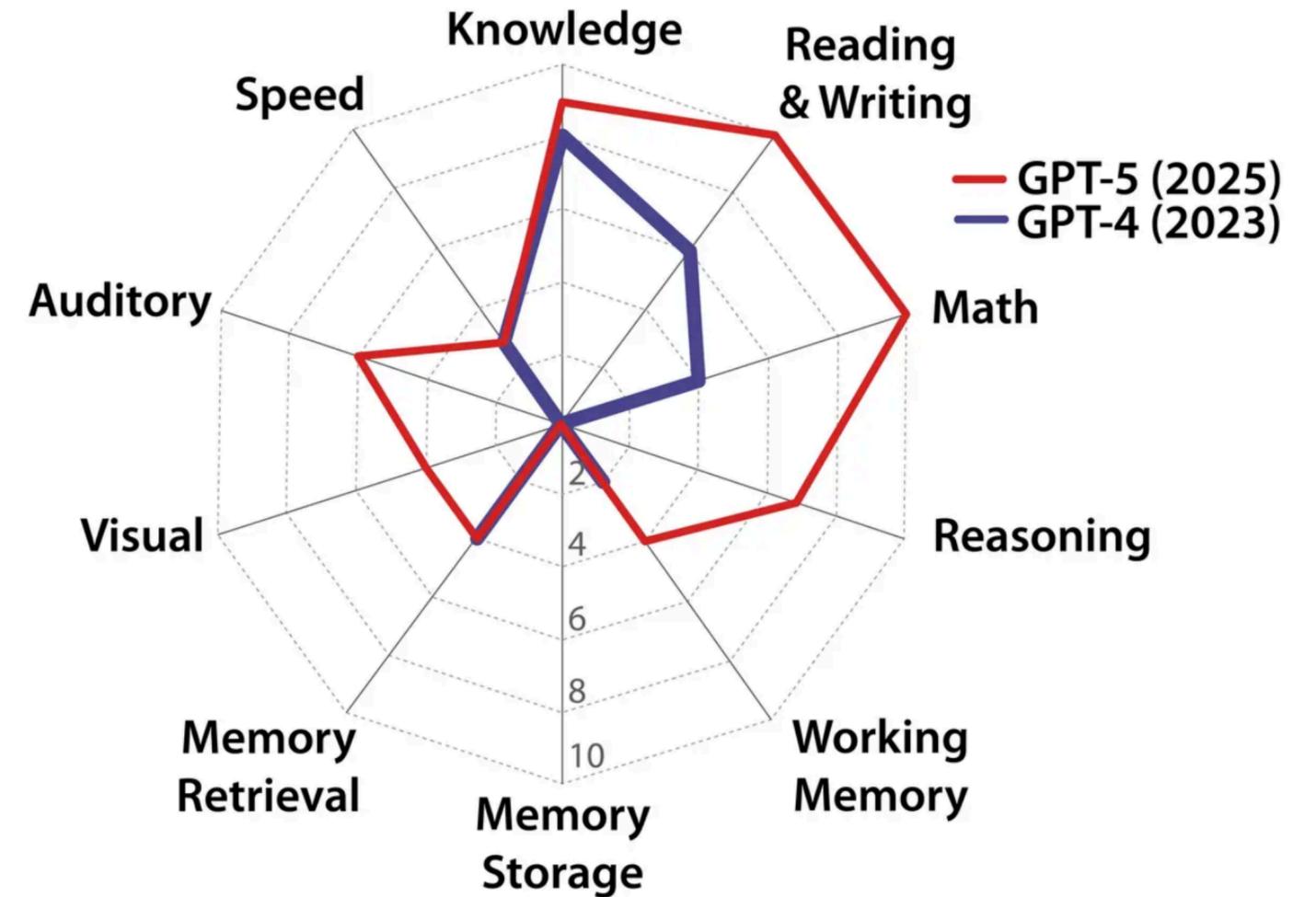
<https://www.nature.com/articles/d41586-025-03959-9>

<https://arxiv.org/html/2511.22570v1>

FINANCIAL TIMES

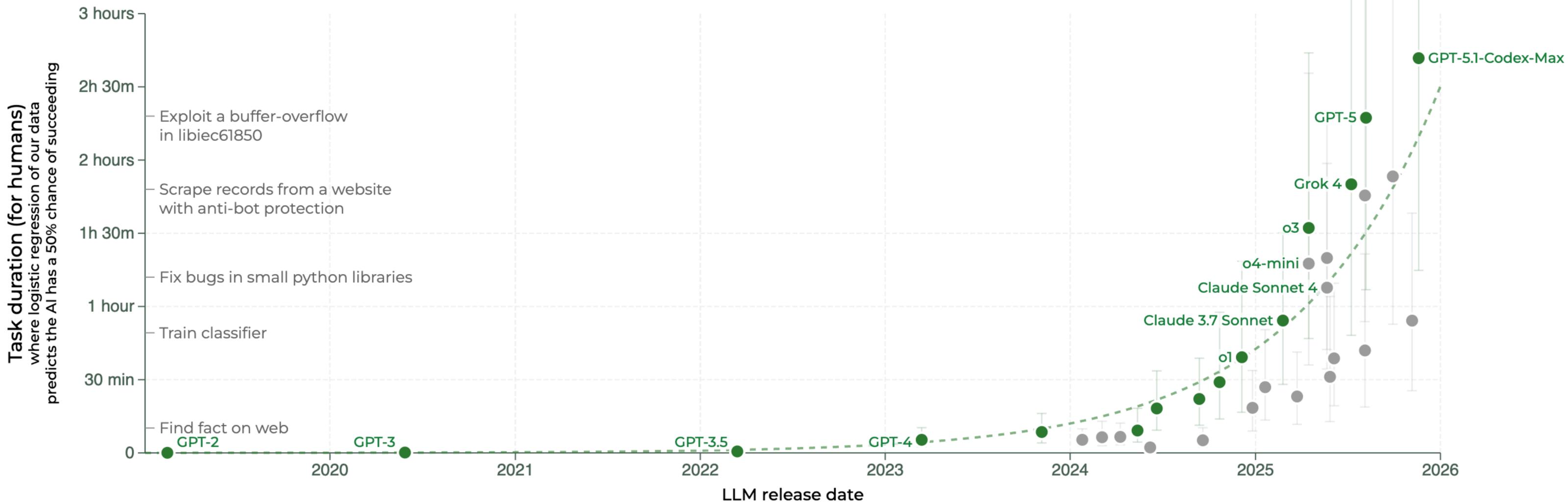
DeepMind and OpenAI achieve gold at 'coding Olympics' in AI milestone

<https://www.ft.com/content/c2f7e7ef-df7b-4b74-a899-1cb12d663ce6>



Hendrycks et al 2025

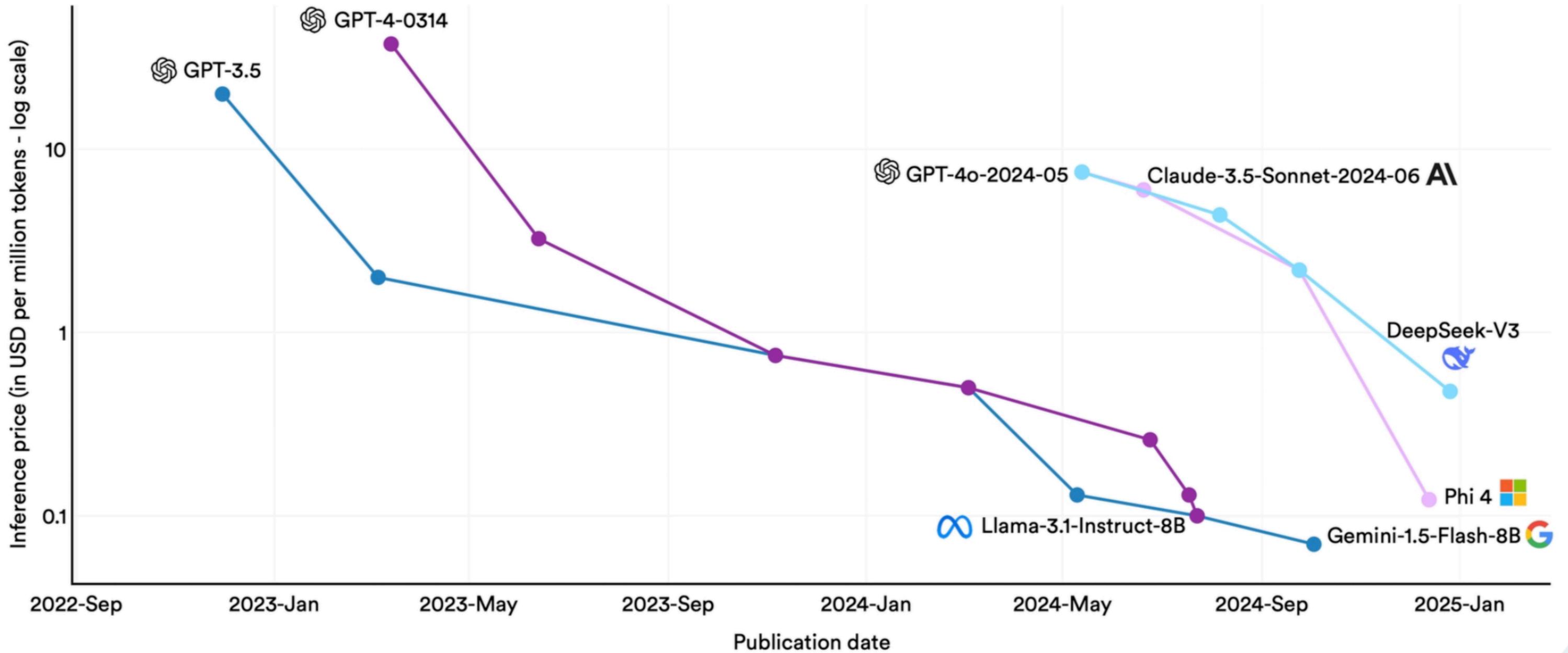
# The time-horizon of software engineering tasks different LLMs can complete 50% of the time



# Inference price across select benchmarks, 2022–24

Source: Epoch AI, 2025; Artificial Analysis, 2025 | Chart: 2025 AI Index report

- GPT-3.5 level+ in multitask language understanding (MMLU)
- GPT-4 level+ in code generation (HumanEval)
- GPT-4o level+ in PhD-level science questions (GPQA Diamond)
- GPT-4o level+ in LMSYS Chatbot Arena Elo



# Argument

1. Quantitative research is mostly cognitive work, done by....

- Humans...
- AI...

**2. If AI helps research, it's 'the future'**

- Humans + AI

# Goal of research?

- **Expand knowledge**
  - “...work undertaken in order to increase the stock of knowledge” (Frascati)
- **Tools change, goal doesn't**
  - Humans + paper
  - Humans + computers
  - Humans + computer + some AI
  - Humans + computer + lots of AI

Key criteria: does it contribute to new knowledge?

**We didn't stop here...**

**...why stop here?**

**...why stop here?!**

Manual research era

Computerization

Early AI

Advanced AI

**Literature review**



Physical searches  
libraries etc

Digital  
databases

Machine learning  
tools

LLM  
Deep research

1-click reviews  
Continuous reviews

**Data analysis**



Manual calculations

Computer-assisted  
analysis

LLM-assisted  
code generation

AI agents

*Plain language analysis*  
*Automated execution*

1950s

1980s

2000s

2020

2025

2030+

***Future Projections***

**We can research tools too!**

Is this optimal?  
How much is repetitive?

Task	Time estimate, now
1. Review Literature	3-12 months
2. Locate data	0.5 months
3. Analyse data	3-6 months
4. Interpret + write-up	3-6 months

**GPT: General purpose technology: can be used for anything!**  
Tools... assistants.... agents

### Almost half of US researchers' time goes on admin

BY REBECCA TRAGER | 9 SEPTEMBER 2014

Of more than 13,000 principal investigators (PIs) that responded to this 2012 workload survey, they said that they spent, on average, 42% of their time on admin, rather than actual research. The FDP, an association c

FDP, 2012-18

### Why can't Epidemiology be automated (yet)?

David Bann, Ed Lowther, Liam Wright, Yevgeniya Kovalchuk

<https://github.com/edlowther/automated-epidemiology>

Bann et al, arxiv / IJE 2025

## Task

1. Review Literature

2. Locate data sources

3. Analyse data

# Systematic reviews

## BMJ Open Analysis of the time and workers needed to conduct systematic reviews of medical interventions using data from the PROSPERO registry

Rohit Borah,<sup>1,2</sup> Andrew W Brown,<sup>2,3</sup> Patrice L Capers,<sup>2,3</sup> Kathryn A Kaiser<sup>2,3</sup>

**Results:** The mean estimated time to complete the project and publish the review was **67.3 weeks** (IQR=42). The number of studies found in the literature

Month	Activity
1 – 2	Preparation of protocol.
3 – 8	Searches for published and unpublished studies.
2 – 3	Pilot test of eligibility criteria.
3 – 8	Inclusion assessments.
3	Pilot test of 'Risk of bias' assessment.
3 – 10	Validity assessments.
3	Pilot test of data collection.
3 – 10	Data collection.
3 – 10	Data entry.
5 – 11	Follow up of missing information.
8 – 10	Analysis.
1 – 11	Preparation of review report.
12 –	Keeping the review up-to-date.

## Task

### 1. Review Literature

2. Locate data sources

3. Analyse data

# Systematic reviews

## Error rates of human reviewers during abstract screening in systematic reviews

Zhen Wang , Tarek Nayfeh, Jennifer Tetzlaff, Peter O'Blenis, Mohammad Hassan Murad

(CI): 2.38% to 8.58%). After abstract screening, the total error rate (false inclusion and false exclusion) was 10.76% (95% CI: 7.43% to 14.09%).



## Opportunities and Challenges for Data Extraction with a Large Language Model

Data extraction in evidence synthesis is labour-intensive, costly, and prone to errors. The use of large language models (LLMs) presents a promising approach for AI-assisted data extraction, potentially enhancing both efficiency and accuracy.

### Our vision

A world where health decisions are based on timely, trusted and relevant evidence.

## Artificial intelligence in systematic reviews: promising when appropriately used

[Sanne H B van Dijk](#)<sup>1,2</sup>, [Marjolein G J Brusse-Keizer](#)<sup>1,3</sup>, [Charlotte C Bucsán](#)<sup>2,4</sup>, [Job van der Palen](#)<sup>3,4</sup>, [Carine J M Doggen](#)<sup>1,5</sup>, [Anke Lenferink](#)<sup>1,2,5,✉</sup>

## Automation of Systematic Reviews with Large Language Models

specificity) and data extraction (*otto-SR*: 93.1% accuracy; human: 79.7% accuracy). Using *otto-SR*, we reproduced and updated an entire issue of Cochrane reviews (n=12) in two days, representing approximately 12 work-years of traditional systematic review work. Across

Cao et al MedRxiv 2025

## Task

### 1. Review Literature

2. Locate data sources

3. Analyse data

# Ad hoc reviews

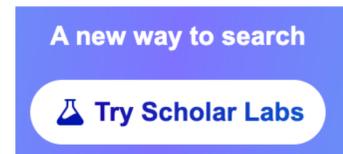
- Now

- Many already using AI  Google Scholar
- Unlike PubMed, indexes social science + grey literature

- Future

-  access to (millions of) papers +  hallucinations

- Plain language Qs



PubMed



- Humans + AI

- Bandwidth freed for higher level tasks
- More ambitious reviews (eg, across disciplines/designs), continual; more informed papers

## Task

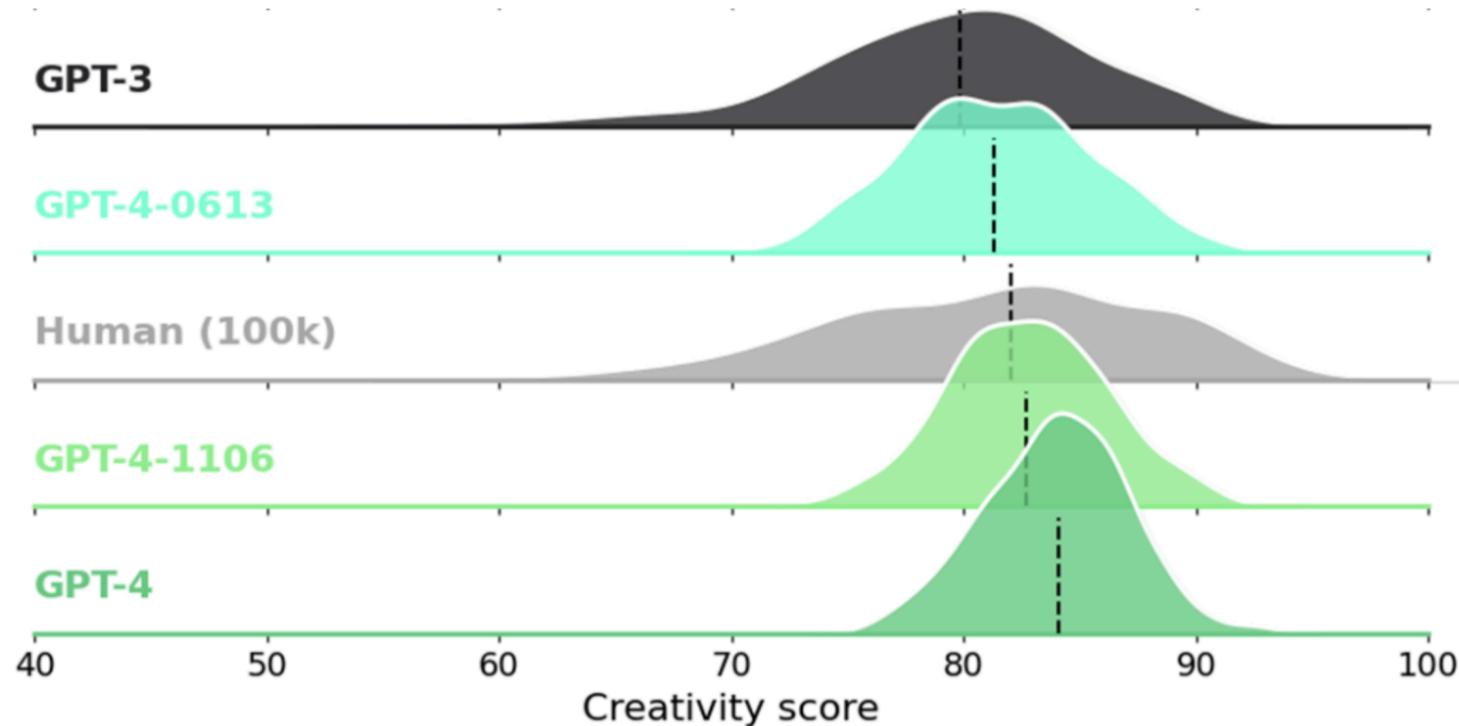
1. Review Literature

2. Locate data sources

3. Analyse data

# AI for forming research questions/hypothesis

- Stochastic parrots?
- Creativity = cognitive task!
  - Their creativity is ~empirically testable
- Even if flawed / limited still useful (instant, cheap)
- AI + humans



Bellemare-Pepin et al 2025

# Creativity Benchmark

A collective industry benchmark for creativity.

<https://creativitybenchmark.ai>

## Task

1. Review Literature

2. Locate data sources

3. Analyse data

# AI for forming research questions/hypothesis + testing them

## Theoretical Physics with Generative AI

Stephen D.H. Hsu

Dec 2025

publication in Physics Letters B after peer review. *Remarkably, the main idea in the paper originated de novo from GPT-5.* GPT-5, Gemini, and Qwen-Max were used extensively to perform calculations, find errors, and generate the finished paper.



Cell

A Cell Press journal

## AI mirrors experimental science to uncover a mechanism of gene transfer crucial to bacterial evolution

[José R. Penadés](#) <sup>1,2,3,4,7,8</sup> · [Juraj Gottweis](#) <sup>5,7</sup> · [Lingchen He](#) <sup>1,2</sup> · ... · [Vivek Natarajan](#) <sup>5</sup> · [Alan Karthikesalingam](#) <sup>5</sup> · [Tiago R.D. Costa](#) <sup>2,3,6</sup> ... [Show more](#)

### Highlights

- AI co-scientist predicted a complex gene transfer mechanism before its publication
- Top AI-generated hypotheses opened new research directions
- **AI bypassed human bias to propose overlooked biological possibilities**
- Benchmarking showed AI co-scientist outperformed other LLMs on this task

Nov 2025

## The AI Scientist: Towards Fully Automated Open-Ended Scientific Discovery

[Chris Lu](#), [Cong Lu](#), [Robert Tjarko Lange](#), [Jakob Foerster](#), [Jeff Clune](#), [David Ha](#)

## Robin: A multi-agent system for automating scientific discovery

[Ali Essam Ghareeb](#), [Benjamin Chang](#), [Ludovico Mitchener](#), [Angela Yiu](#), [Caralyn J. Szostkiewicz](#), [Jon M. Laurent](#), [Muhammed T. Razzak](#), [Andrew D. White](#), [Michaela M. Hinks](#), [Samuel G. Rodrigues](#)

Scientific discovery is driven by the iterative process of background research, hypothesis generation, experimentation, and data analysis. Despite recent advancements in applying artificial intelligence to scientific discovery, no system has yet automated all of these stages in a single workflow. Here, we introduce Robin, the first multi-agent system capable of fully automating the key intellectual steps of the scientific process. By

## Early science acceleration experiments with GPT-5

[Sébastien Bubeck](#)<sup>1</sup>, [Christian Coester](#)<sup>2</sup>, [Ronen Eldan](#)<sup>1</sup>, [Timothy Gowers](#)<sup>3</sup>, [Yin Tat Lee](#)<sup>1</sup>, [Alexandru Lupsasca](#)<sup>1,4</sup>, [Mehtaab Sawhney](#)<sup>5</sup>, [Robert Scherrer](#)<sup>4</sup>, [Mark Sellke](#)<sup>1,6</sup>, [Brian K. Spears](#)<sup>7</sup>, [Derya Unutmaz](#)<sup>8</sup>, [Kevin Weil](#)<sup>1</sup>, [Steven Yin](#)<sup>1</sup>, [Nikita Zhivotovskiy](#)<sup>9</sup>

## Task

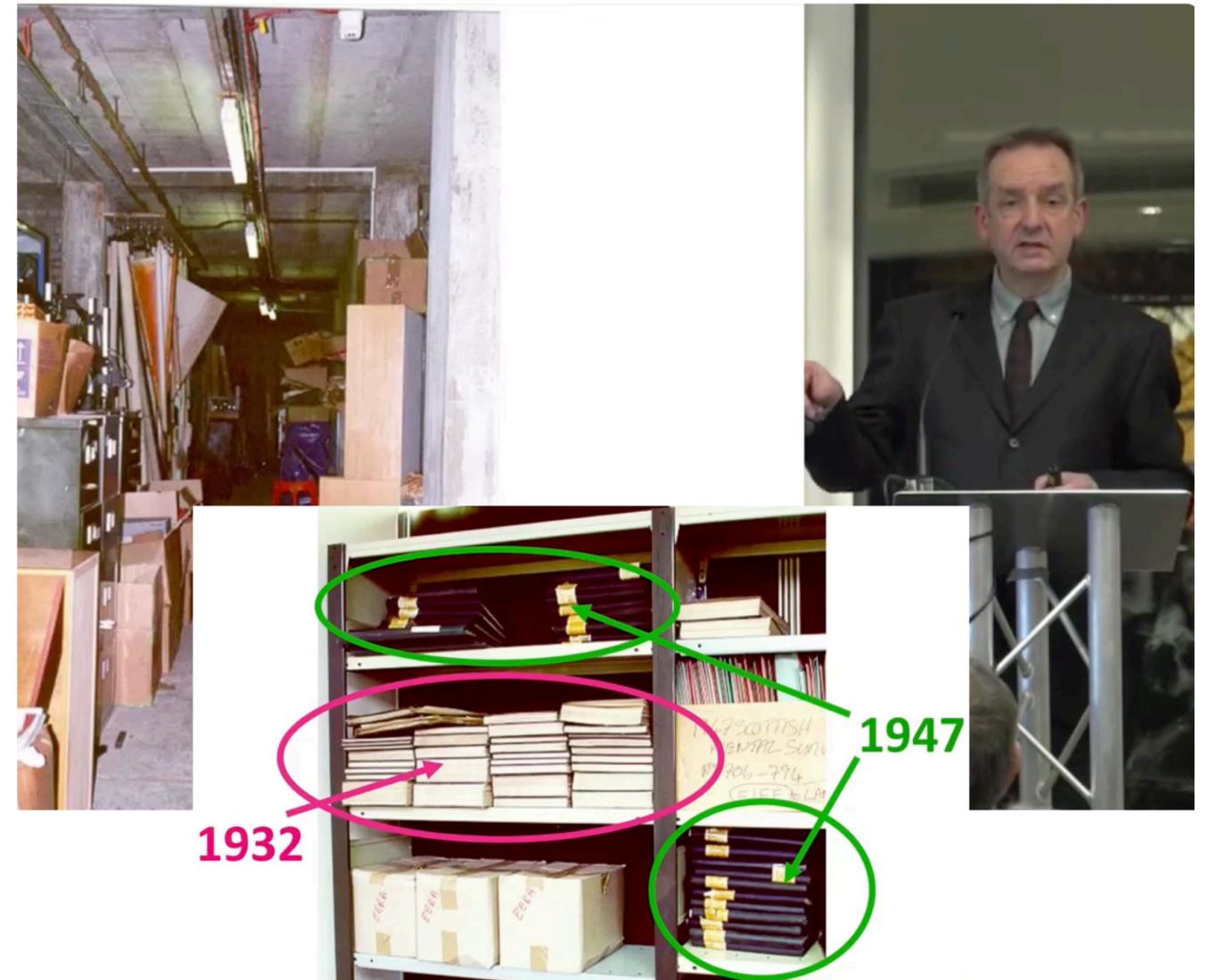
1. Review Literature
- 2. Locate data sources**
3. Analyse data

1. Unlocking historic (new) data

2. With collected data

- Manual variable work -> AI assisted (e.g. Harmony for harmonisation; pdf annotation, metadata curation etc)

3. One research question -> all studies



“Found end of corridor in 1997”

The **Lothian Birth Cohort 1936**: a study to examine influences on cognitive ageing from age 11 to age 70 and beyond

[IJ Deary, AJ Gow, MD Taylor, J Corley, C Brett...](#) - BMC geriatrics, 2007 - Springer

... A total of 1091 participants make up the **Lothian Birth Cohort 1936**. They undertook: a medical interview and examination; physical fitness testing; extensive cognitive testing (reasoning, ...

☆ Save 📄 Cite Cited by 563 Related articles All 27 versions

## Task

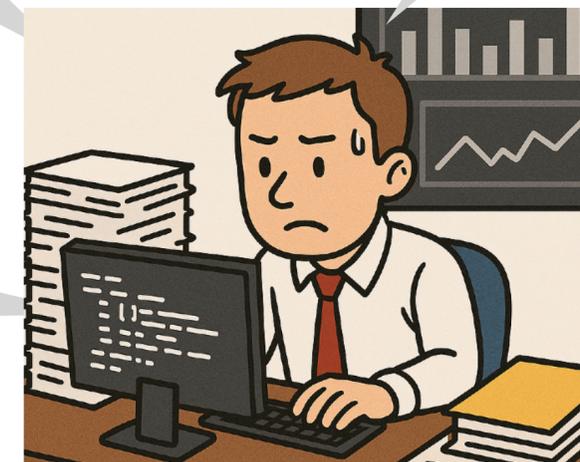
1. Review Literature
2. Locate data sources
- 3. Analyse data**

# Coding: humans only

Error!?

I've spent years typing  
“( % etc

\*check documentation\*



\*check forum\*

\*2 hours later\*  
I used ` instead of ‘



It works! But across a scientific lifetime, how much is wasted?

```
#load the libraries
library(gamlss)
library(sitar)
library(tidyverse)
library(haven)
library(psych) |

#load all data
dat <- read_dta('cleaned_datafor_r.dta')

#bmi
bmi46<- subset(dat, select = c("bmi46", "sex", "fscb" , "pa42b"))

#bmi
describeBy(bmi46, group="sex" ,digits=2)
describeBy(bmi46, group="fscb" ,digits=2)
describeBy(bmi46, group="pa42b" ,digits=2)

#output estimates
bmi46sex_no <- gamlss(bmi46 ~ sex ,
                    sigma.formula = ~sex ,
                    family = NO (mu.link = log), data = na.omit(bmi46) )
```

## Task

1. Review Literature
2. Locate data sources
- 3. Analyse data**

# Coding: humans + AI autocomplete

```
TS Calculator.ts > Calculator
1 class Calculator {
2
3   public add(num1: number, num2: number): number {
4     return num1 + num2;
5   }
6
7   public subtract(num1: number, num2: number): number {
8     return num1 - num2;
9   }
10 }
```

Independent report

# AI coding assistant trial: UK public sector findings report

Published 12 September 2025

GDS ran a trial of AI coding assistants (AICAs) across government from November 2024 to February 2025. A total of 2,500 licences were made available across central government organisations.

- 67% of respondents reporting a reduction in time spent searching for information or examples
- 65% reporting faster task completion
- 56% reporting more efficient problem solving

On average, users reported time savings of 56 minutes per working day, which equates to approximately 28 working days saved per user annually, based on a standard working calendar. Of this, 24 minutes per day were associated with creation of code or analysis

User sentiment was positive, with 58% expressing they would not want to return to their pre-AICA working conditions.

## Task

1. Review Literature
2. Locate data sources
- 3. Analyse data**

# Coding: humans + AI: plain language -> code



Use any!

## Task

1. Review Literature
2. Locate data sources
- 3. Analyse data**

# Coding: humans + AI: plain language -> code

CHAT

+

⌵

🕒

⚙️

⋮

|

🗄️

✕



## Ask about your code

AI responses may be inaccurate.

[Generate Agent Instructions](#) to onboard AI onto your codebase.

```
import matplotlib.pyplot as plt
import matplotlib.cm as cm
import numpy as np
from pathlib import Path

# --- Extract the country rows from your table1 ---
# (they start after the 'Countries' row)
country_start = row_names.index('Countries') + 1
country_rows = row_names[country_start:]
```



 Main\_analyses\_supplemental.ipynb • Cell 14 +

Explore and understand your code

Ask ⌵ GPT-4.1 ⌵

```
import matplotlib.pyplot as plt
import matplotlib.cm as cm
import numpy as np
from pathlib import Path

# --- Extract the country rows from your table1 ---
# (they start after the 'Countries' row)
country_start = row_names.index('Countries') + 1
country_rows = row_names[country_start:]
```

Agent 

Plan

✓ Ask

Edit

---

AIAgentExpert

[Configure Custom Agents...](#)

Auto	10% discount
✓ GPT-4.1	0x
GPT-4o	0x
GPT-5 mini	0x
Grok Code Fast 1	0x
Raptor mini (Preview)	0x
Claude Haiku 4.5	0.33x
Claude Opus 4.5 (Preview)	3x
Claude Sonnet 4	1x
Claude Sonnet 4.5	1x
Gemini 2.5 Pro	1x
Gemini 3 Pro (Preview)	1x
GPT-5	1x
<a href="#">Manage Models...</a>	

## Task

1. Review Literature
2. Locate data sources
- 3. Analyse data**

# Coding: humans + AI agents

<https://julius.ai/>

D

 Mental health D... spreadsheet 4.66 MB

I uploaded the following files: Mental health Depression disorder Data.csv



J

Thank you for uploading the file "Mental health Depression disorder Data.csv". I'll take a look at the data and provide you with an overview of its contents. Let's start by examining the structure and the first few rows of the dataset.

`</>` Show Code ▾

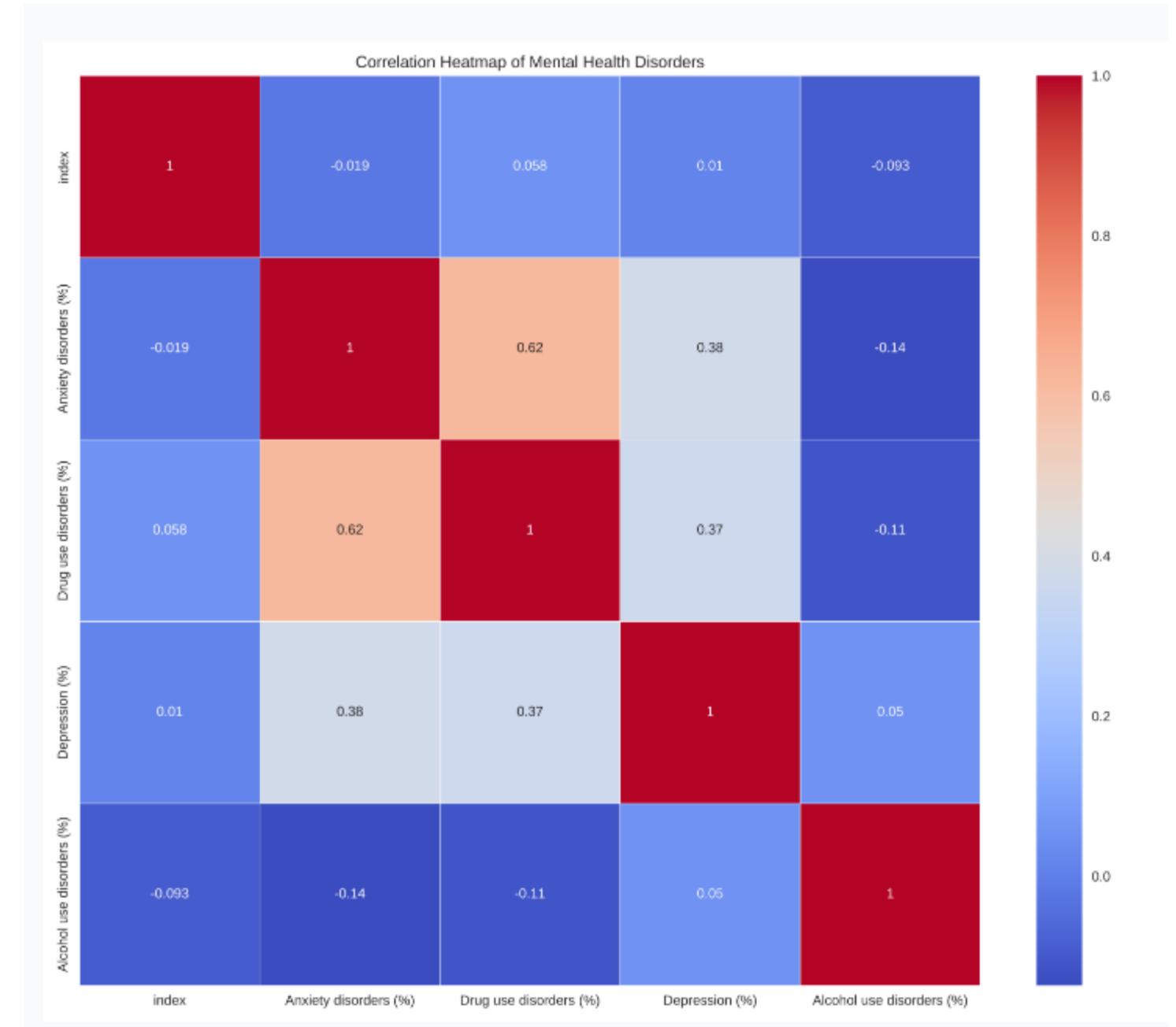
I've analyzed the Mental health Depression disorder Data.csv file you uploaded. Let me provide you with an overview of the dataset and some initial insights.

First, let's look at the structure of the data:

(Optional) click on columns or rows to focus on specific data

 Export 

	index	Entity	Code	Year	Schizophrenia (%)	Bipolar disorder (%)	Eating disorders (%)	Anxiety disorders (%)	Drug use disorders (%)	Depression (%)	Alcohol use disorders (%)
0	0	Afghanistan	AFG	1990	0.16056	0.697779	0.101855	4.82883	1.677082	4.071831	0.672404
1	1	Afghanistan	AFG	1991	0.160312	0.697961	0.099313	4.82974	1.684746	4.079531	0.671768



## Task

1. Review Literature
2. Locate data sources
- 3. Analyse data**

# Coding: humans + AI agents

## Associations between Birth Weight and BMI

Notebook generated using API access to `gpt-4.1-2025-04-14` via the `data-analysis-crow` framework on 23 June 2025.

### Prompt:

The attached dataset contains information about participants in a longitudinal study. Investigate whether there is an association between birth weight and BMI using linear regression analyses adjusted for sex.

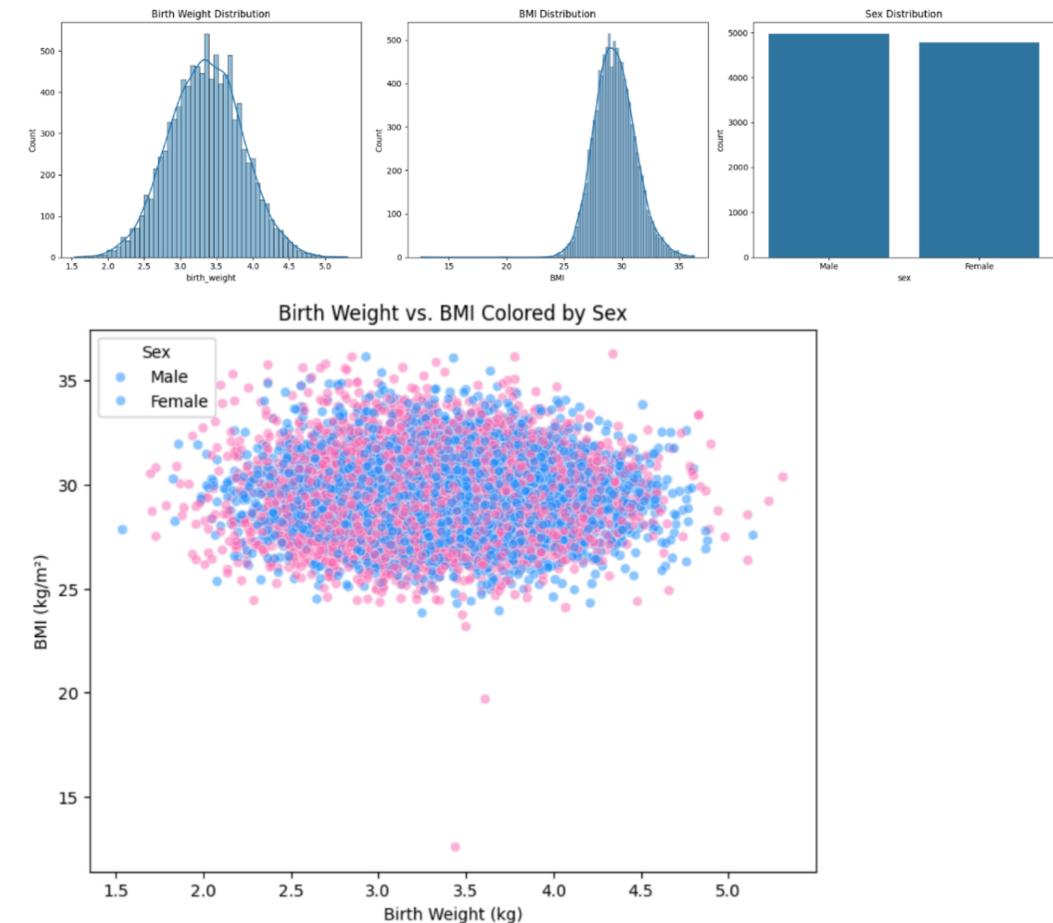
```
# Step 1: Load Data and Perform Descriptive Statistics
#
# Analysis Planning:
#
# - Relevant file: 'cohort-data.csv' looks like the main dataset. We will load this file with pandas.
# - We will inspect the file's shape and column names with .head() and .info() to understand the structure and
# - Use .describe() for numerical summaries (mean, std, min, max, etc.) for birth weight and BMI.
# - Check for missing data with .isnull().sum().
# - If there are missing/strange values (e.g., negative birth weights or BMI), we will note, and consider handling
# - Plan to present summaries as printed outputs and initial visual inspection if columns allow.
# - Key statistics: distribution and range of birth weight, BMI; balance of sexes.
# - Data quality issues: unexpected values, missing data, unexpected data types; will address as found.

import pandas as pd

# Load the data file
cohort = pd.read_csv('cohort-data.csv')

# Examine the structure and first few rows
display(cohort.head())
print(cohort.info())
print(cohort.describe(include='all'))
# Check for missing data
print('Missing values per column:')
print(cohort.isnull().sum())
```

	ID	sex	education	birth_weight	adult_body_weight	height	salary	depression_symptoms
0	1	0	4	3.29	91.882406	1.791942	70603.442728	20.433706
1	2	1	3	3.48	81.210334	1.619915	53981.525954	13.147941



1	2	US	1558	PHX	CLT	3	15	222	1
2	3	AA	2400	LAX	DFW	3	20	165	1
3	4	AA	2466	SFO	DFW	3	20	195	1
4	5	AS	108	ANC	SEA	3	30	202	0

### # Dataset Features Discussion

#### ## Overview

This airlines dataset contains **539,383 flight records** with **9 features** that describe various aspects of flight information.

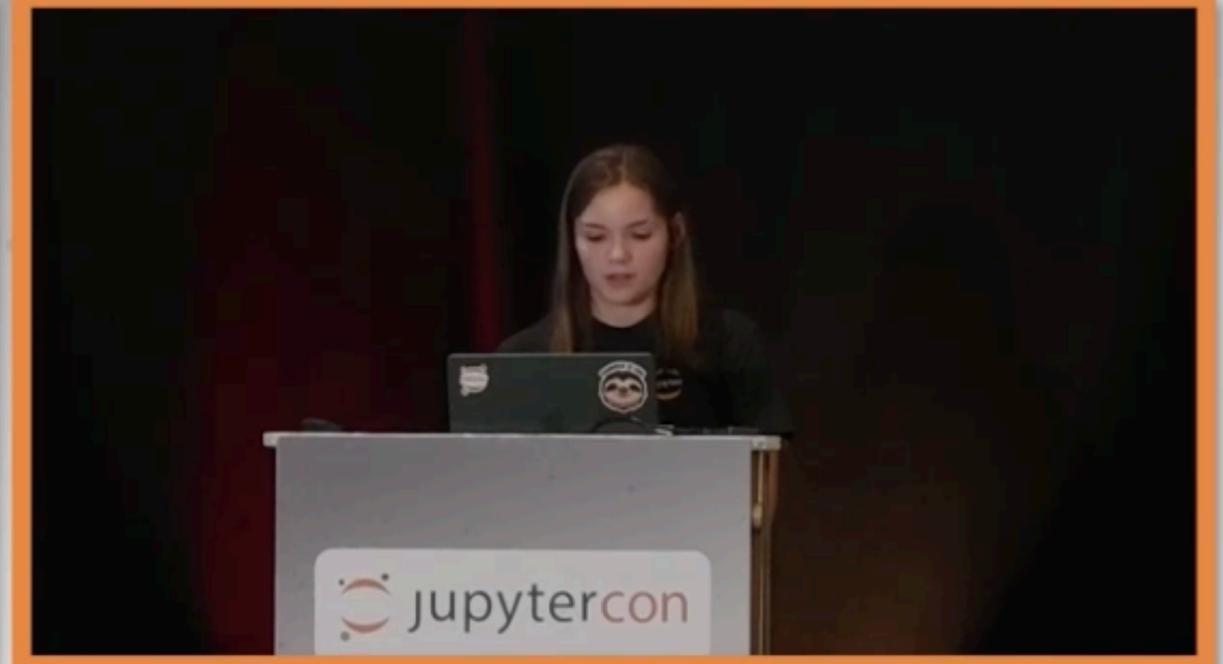
#### ## Feature Descriptions

- id**: Unique identifier for each flight record
- Airline**: Airline carrier code (e.g., CO, US, AA, AS)
- Flight**: Flight number assigned by the airline
- AirportFrom**: Departure airport code (e.g., SFO, PHX, LAX)
- AirportTo**: Arrival/destination airport code (e.g., IAH, CLT, DFW)
- DayOfWeek**: Day of the week (numeric representation, likely 1-7)
- Time**: Scheduled departure time
- Length**: Flight duration in minutes
- Delay**: Binary indicator (0 = on-time, 1 = delayed) - **Target variable for prediction**

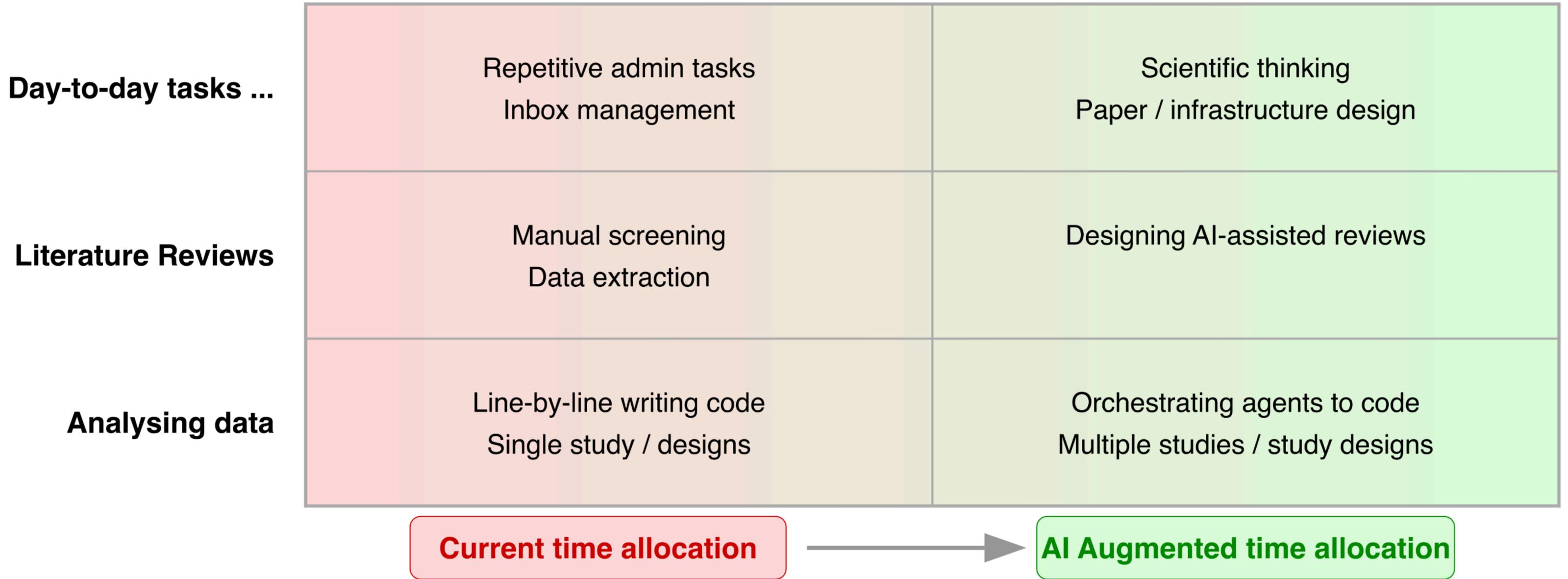
#### ## Key Observations

- The dataset appears suitable for classification tasks (predicting flight delays)
- Features include both categorical (Airline, airports) and numerical (Time, Length) variables
- The target variable 'Delay' is binary, making this a binary classification problem

```
[ ]: # todo: try to use statistics to compute the probability that a flight leaving SAn Diego On THursday will
# be delayed |
```



 **Jupytercon**





## Centre for Longitudinal Studies (CLS)

Code repositories from CLS

11 followers

United Kingdom

<https://cls.ucl.ac.uk/>

@clscohorts

[company/ucl-centre-for-longitudin...](https://www.linkedin.com/company/ucl-centre-for-longitudinal-studies/)

### Harmonisation / Derived variable code

Core datasets (useful variables across domains e.g., health, socioeconomic, demographic)

- [MCS Core](#)
- [Next Steps Core](#)
- BCS70 + NCDS forthcoming in 2026

Cross-cohort harmonisation:

- [Fertility Histories](#)
- [Smoking](#)
- [Asthma](#)

Code for single cohorts:

- [Time use data in MCS](#)

### Genomic Projects

- [CLS Polygenic Index Repository](#) - Tools and scripts for calculating polygenic scores in longitudinal studies

### Data Management Code

Tools and resources for longitudinal data management:

- See [CLS Data Management Resources](#) - code for common data management tasks

<https://github.com/cls-data>

# CLS PGI Pipeline (v1.0)

## Contributors

**Tim Morris** - Designed the pipeline and led development and implementation.

**Gemma Shireby** - Designed the pipeline and led curation of genotype data.

**Georg Otto** - Contributed to pipeline design including feature suggestions.

**David Bann** - Provided code annotations and documentation support.

**Liam Wright** - Checked PGIs and designed PGI visualisations.

[https://github.com/CLS-Data/CLS\\_PGI\\_repository](https://github.com/CLS-Data/CLS_PGI_repository)

## CLS Data Handling Guide

🔍 Search CLS Data Handling Guide

[Introduction](#)

[Downloading Documentation](#)

[NCDS](#) ▾

[BCS70](#) ▾

[Next Steps](#)

**MCS** ▲

[Creating a Simple Folder Structure](#)

[Data Discovery](#)

[Data Structures](#)

[Working with the Household Grid](#)

[Combining Data Across Sweeps](#)

[Combining Data Within A Sweep](#)

[Reshaping Data from Long to Wide \(or Wide to Long\)](#)

This section presents code to clean and handle data from the Millennium Cohort Study (MCS). The MCS has the most complex data structures of CLS's cohorts, so we provide an [introduction to the data](#) as necessary background.

### TABLE OF CONTENTS

- [Creating a Simple Folder Structure](#)
- [Data Discovery](#)
- [Data Structures](#)
- [Working with the Household Grid](#)
- [Combining Data Across Sweeps](#)
- [Combining Data Within A Sweep](#)
- [Reshaping Data from Long to Wide \(or Wide to Long\)](#)

<https://cls-data.github.io/>

# NCRM hackathon Nov 2025

**Survey Variable Selector & Merge Script Generator**  
Upload your lookup CSV to browse, filter, and select variables. Generate merge scripts in R, Stata, or Python.

Choose lookup.csv | lookup.csv (30,229 variables)

Select All Visible | Deselect All | Load Preset | Save Selection | Toggle View

VARIABLE NAMES: e.g., income\*, educ?, health  
LABELS: e.g., employment, disease  
WAVE/FOLDER: All waves  
DATASET: All datasets

Variable	Label	Folder	Dataset	Position
BCSID	research case identifier	0y	bcs1derived.dta	1
BD1AGEFB	1970: Age of mother at first birth	0y	bcs1derived.dta	5
BD1CNTRY	1970: Country of Interview	0y	bcs1derived.dta	2
BD1FAGE	1970: Age of father at CM's birth (from s2 var e009)	0y	bcs1derived.dta	8
BD1FAGM	1970: Age of father at present marriage	0y	bcs1derived.dta	10
BD1MAGE	1970: Age of mother at CM's birth (from s1 var a0005a/s2 var e008)	0y	bcs1derived.dta	7
BD1MAGM	1970: Age of mother at present marriage	0y	bcs1derived.dta	9
BD1PSOC	1970: Social class at birth: fathers occup or mothers (vars a0014 + a0018)	0y	bcs1derived.dta	4
BD1REGN	1970: Standard Region of residence	0y	bcs1derived.dta	3
BD1TEENM	1970: Ever a teenage mother (BD1AGEFB grouped)	0y	bcs1derived.dta	6
a0002	Multiplicity Code	0y	bcs7072a.dta	2
a0005a	Mothers age at Delivery	0y	bcs7072a.dta	3
a0006a	Region of Birth of Mother	0y	bcs7072a.dta	4
a0007a	Region of Birth of Father	0y	bcs7072a.dta	5
a0008a	Reaion of Birth of Mother's Mother	0v	bcs7072a.dta	6

Showing 1 to 50 of 30,229 variables | Show 50 entries | Previous 1 2 3 4 5 ... 605 Next

**Configuration**

Ollama Model Name ?  
gemma3:4b

Ensure Ollama is running locally: ollama serve

Session ID: 29b32a7b-3d6f-49f0-af60-1945659a391b

## Is your industry in this list?

Select an option:

- Select...
- Yes
- No

AI: Okay, let's try this: "Could you tell me what types of plants you primarily sell or grow – for example, are you focused on h

Your Answer

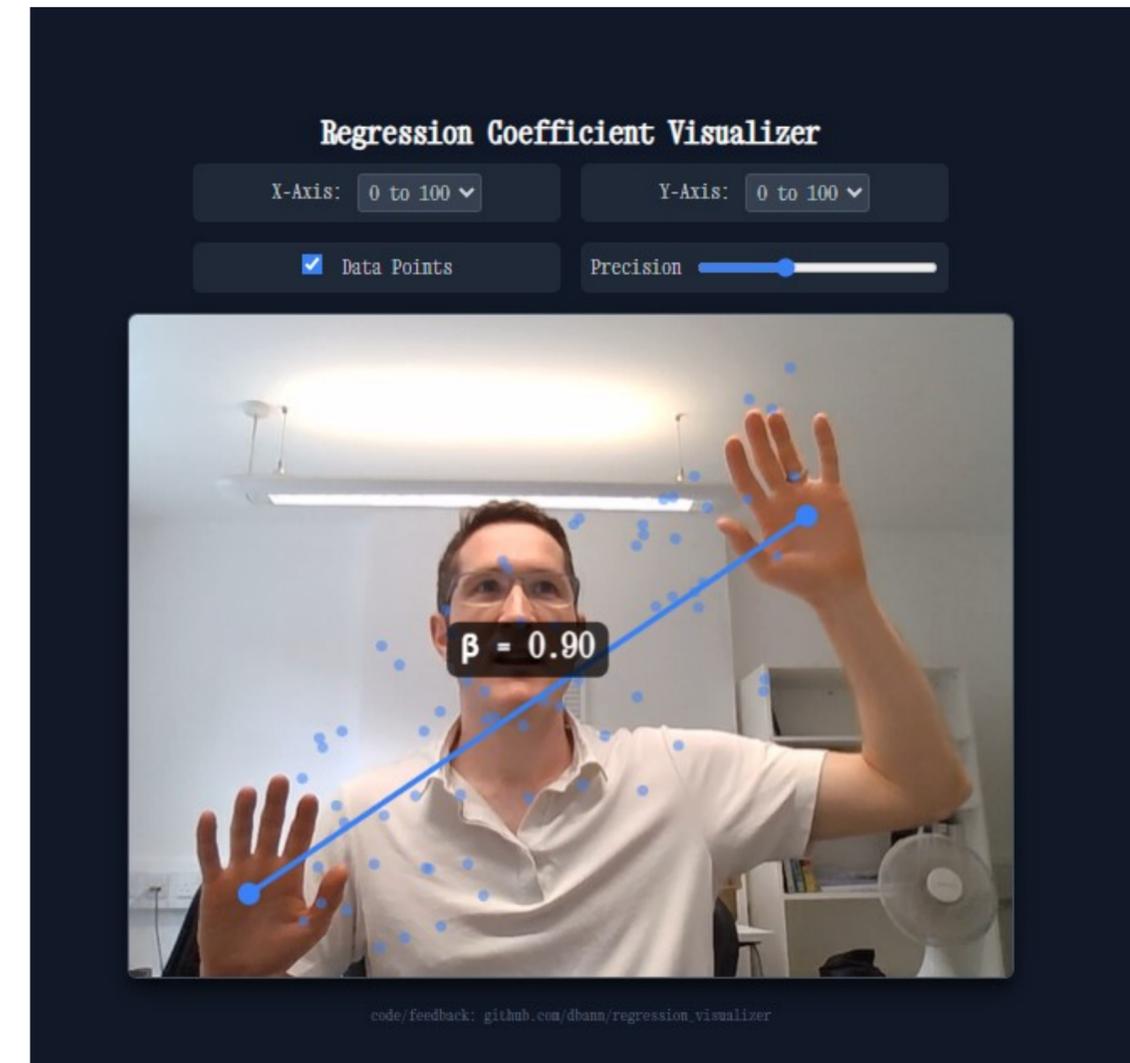
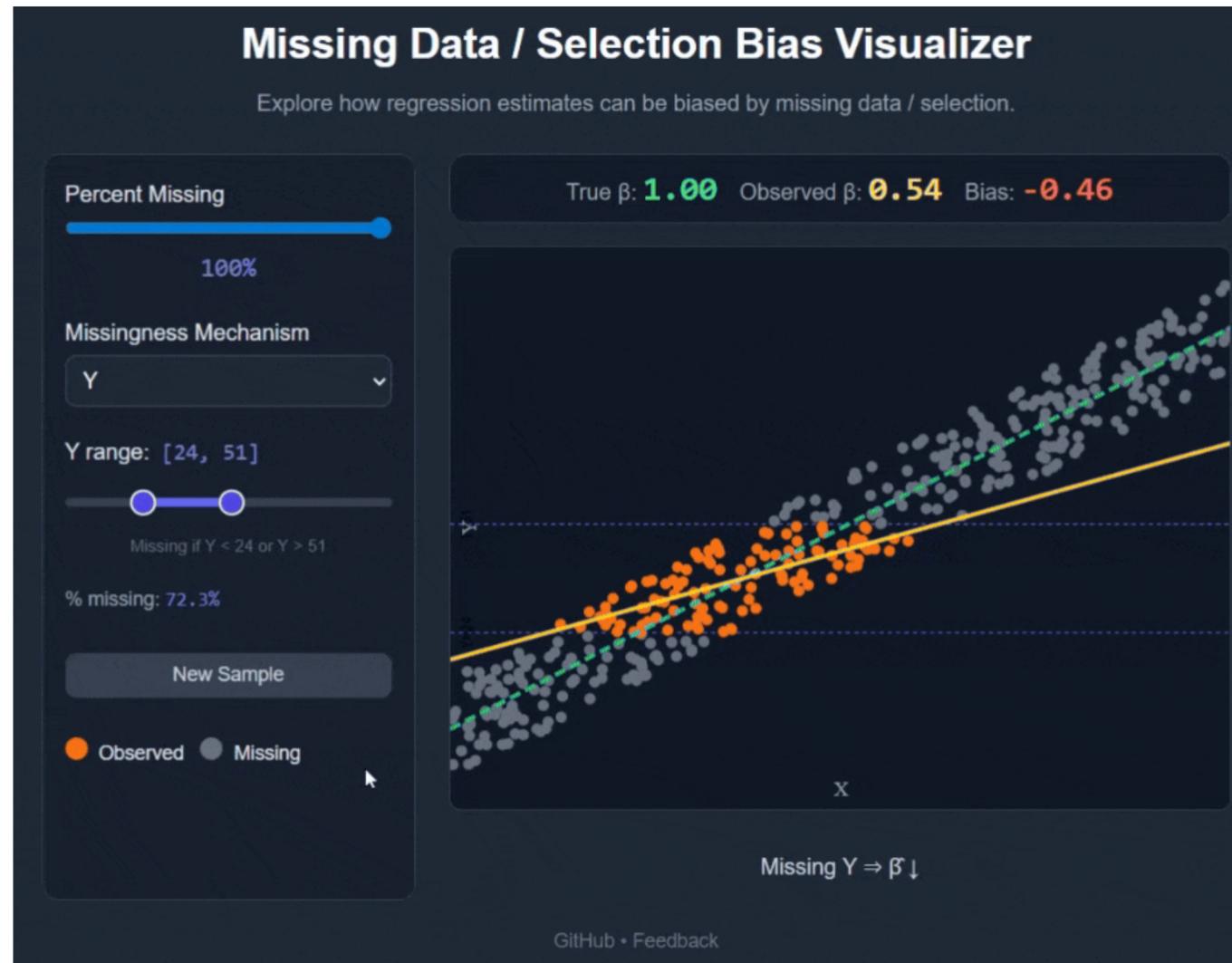
Refine Search

**Liam Wright** liam.wright@ucl.ac.uk

<https://github.com/dbann/surveymerge>



# Teaching, making visualisers



<https://github.com/dbann/selection>

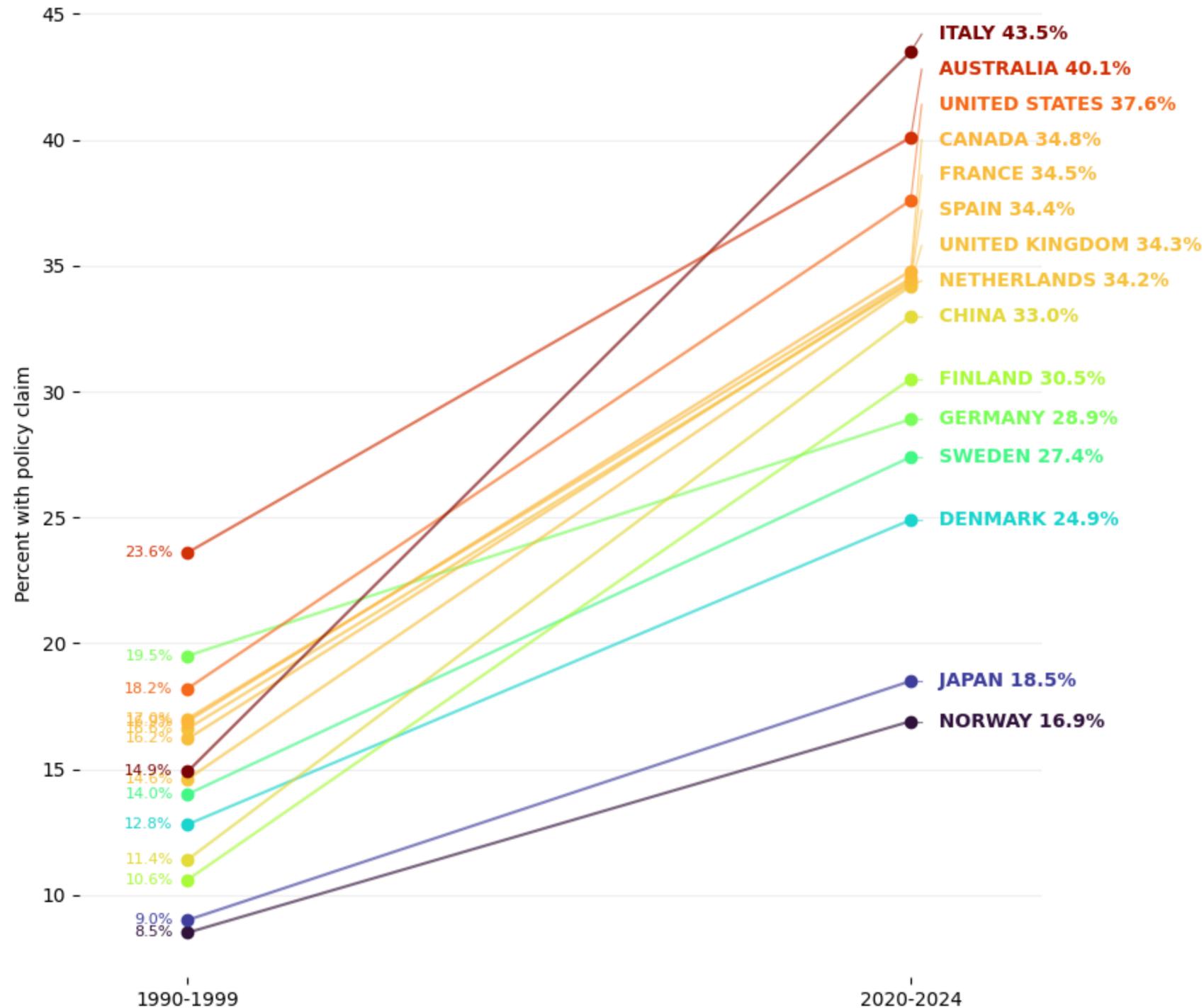
[https://github.com/dbann/regression\\_visualizer](https://github.com/dbann/regression_visualizer)

## Task

1. Review Literature
2. Locate data sources
3. Analyse data

# Example of a human + AI paper

Slopegraph: % with policy claim (1990-1999 → 2020-2024)



- n=44k abstracts, \$3
- Concordance in classification (kappa)
  - Human-human = 0.83
  - LLM (Deepseek 3,1) -human = 0.73-0.84
- LLMs used to learn
  - Python, APIs, Git (Vs Code)
- Humans + LLMs
  - New research, new field (metascience)
  - New data / syntax resource
  - Help early career researcher (Mengyao Wang)

# Science Funding Cost Calculator

Explore the hidden costs of the grant application process. How much time and money is spent on unfunded proposals?

## Parameters

Words requested by funder



Applications funded (%)



2nd Stage Selection Rate (%)



Scientists' Hourly Rate (£)

Total Funding Applications

Writing Speed

Currency: £



**£3,560,000**

Total Wasted Cost



**71,200**

Hours Wasted

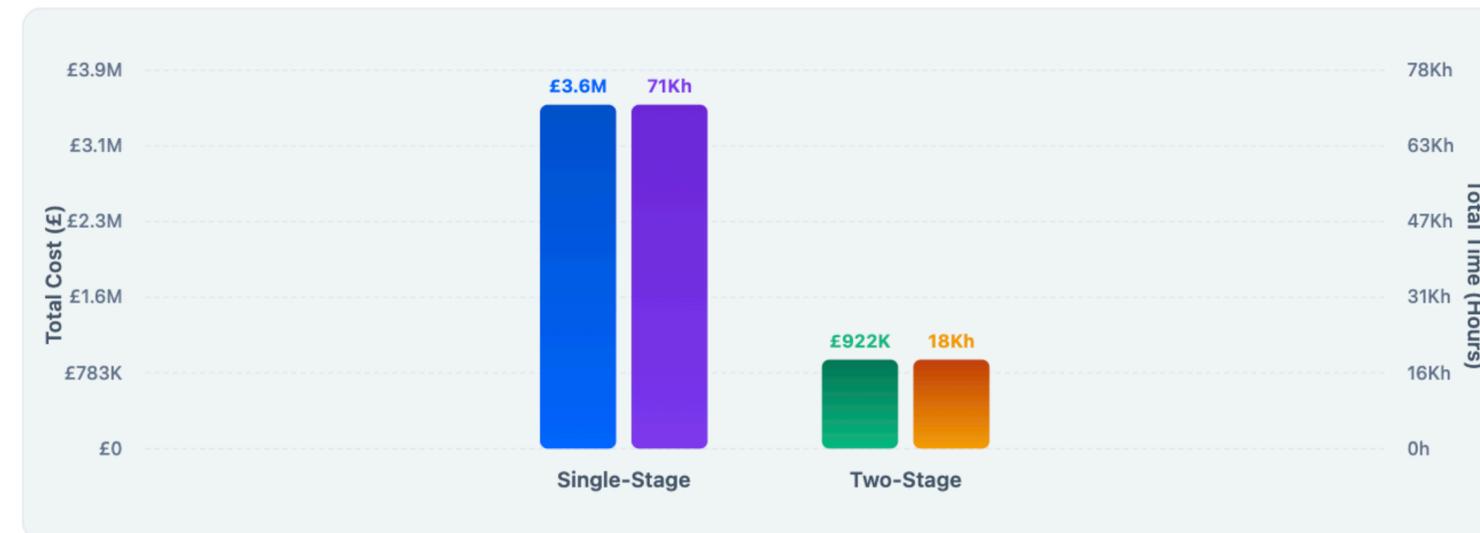


**11%**

Success Rate

## Cost Breakdown and Comparison with 2-Stage System

Comparing single-stage vs. two-stage review processes



## Potential Cost Savings

With 2-stage system (1000 words first stage, 30% proceeding):

**£2,638,000**

## Potential Time Savings

Reduced application writing burden:

**52,760h**

## Task

1. Review Literature
2. Locate data sources
3. Analyse data

# Does good > bad? 1) Our data are private

Use open-weight LLMs!

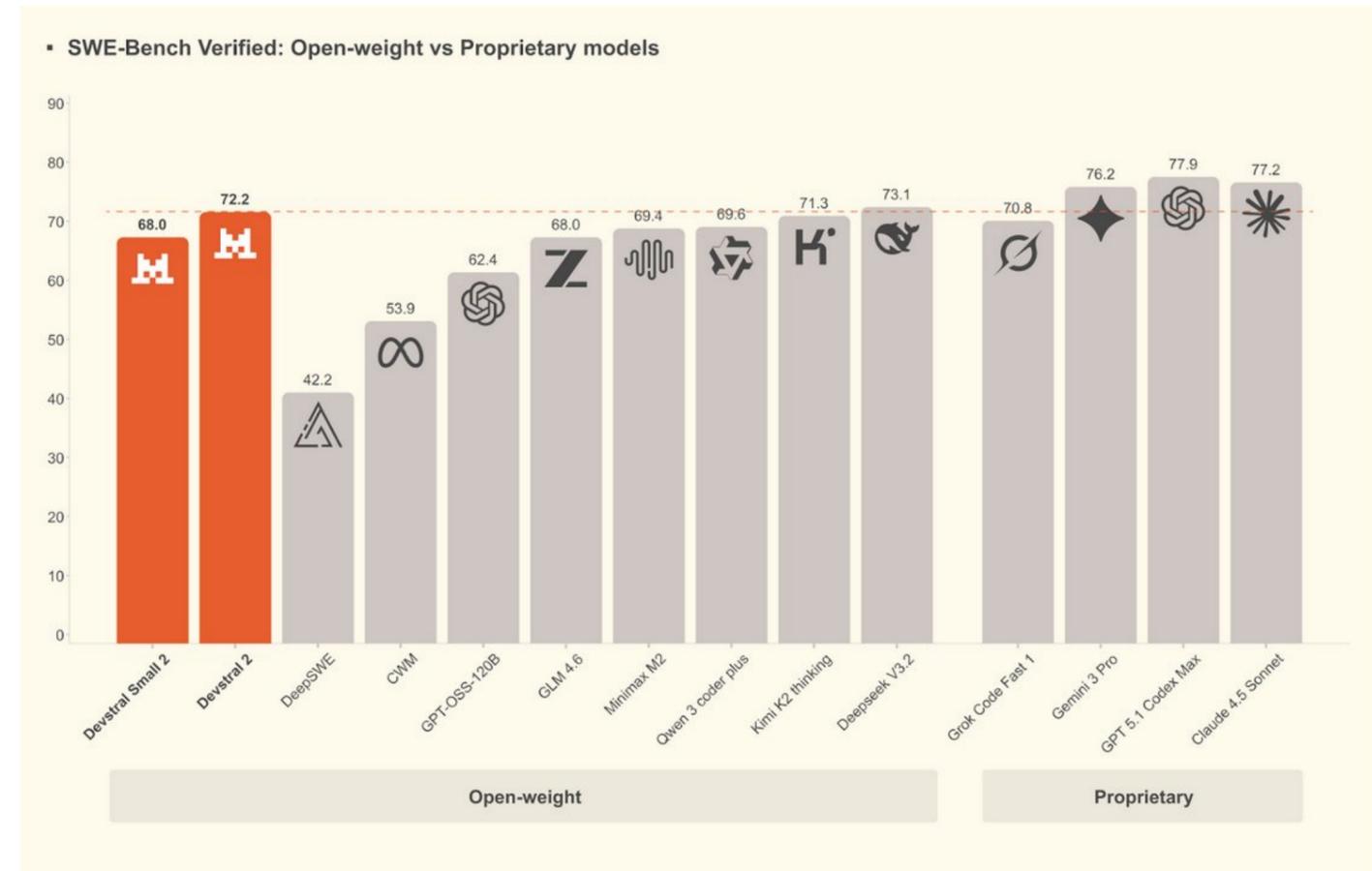
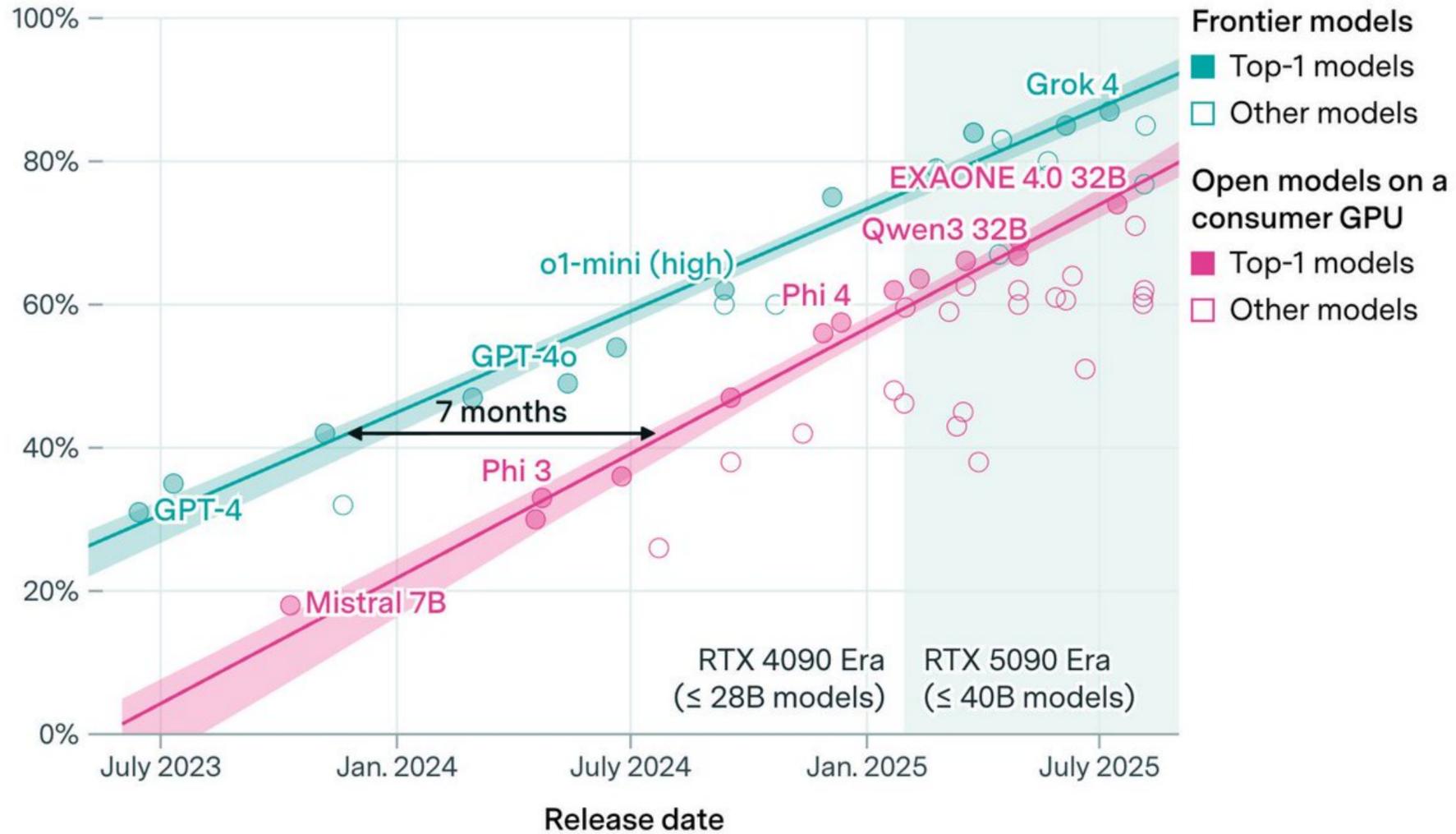
 <b>Closed weight LLMs</b>	<b>Open weight variant</b>	 <b>Open weight LLMs</b>
<a href="#">ChatGPT</a> (OpenAI  )	<a href="#">ChatGPT-OSS</a> (OpenAI  )	<a href="#">Deepseek</a> (Deepseek  )
<a href="#">Gemini</a> (Google  )	<a href="#">Gemma</a> (Google  )	<a href="#">Kimi k2</a> (Moonshot AI  )
<a href="#">Claude</a> (Anthropic  )	n/a	<a href="#">Mistral</a> (Mistral AI  )
<a href="#">Grok</a> (xAI  )	<a href="#">Earlier Grok releases</a> (xAI  )	<a href="#">Qwen</a> (Alibaba Qwen  )
		<a href="#">Llama</a> (Meta  )
		<a href="#">Ernie</a> (Baidu  )
		<a href="#">GLM</a> (Zhipu AI  )

Task
1. Review Literature
2. Locate data sources
3. Analyse data

# Does good > bad? 1) Our data are private

EPOCH AI

GPQA-Diamond accuracy



<https://mistral.ai>

Here's R code to create a ridgeline plot using the `ggridges` package:

```
r Copy  
  
# Install and load required packages  
install.packages(c("gggridges", "ggplot2", "dplyr"))  
library(gggridges)  
library(ggplot2)  
library(dplyr)  
  
# Example 1: Using built-in dataset  
# Create a sample dataset  
set.seed(123)  
data <- data.frame(  
  value = c(rnorm(100, 20, 5), rnorm(100, 30, 7), rnorm(100, 40, 6)),  
  group = rep(c("Group A", "Group B", "Group C"), each = 100),  
  time = rep(1:100, 3)  
)  
  
# Basic ridgeline plot  
ggplot(data, aes(x = value, y = group, fill = group)) +  
  geom_density_ridges(alpha = 0.7) +  
  theme_minimal() +  
  labs(title = "Ridgeline Plot Example",  
       x = "Value",  
       y = "Group")  
  
# Example 2: Using mtcars dataset  
# Prepare data for ridgeline plot  
mtcars_grouped <- mtcars %>%  
  mutate(cyl = factor(cyl)) %>%  
  select(mpg, cyl)
```

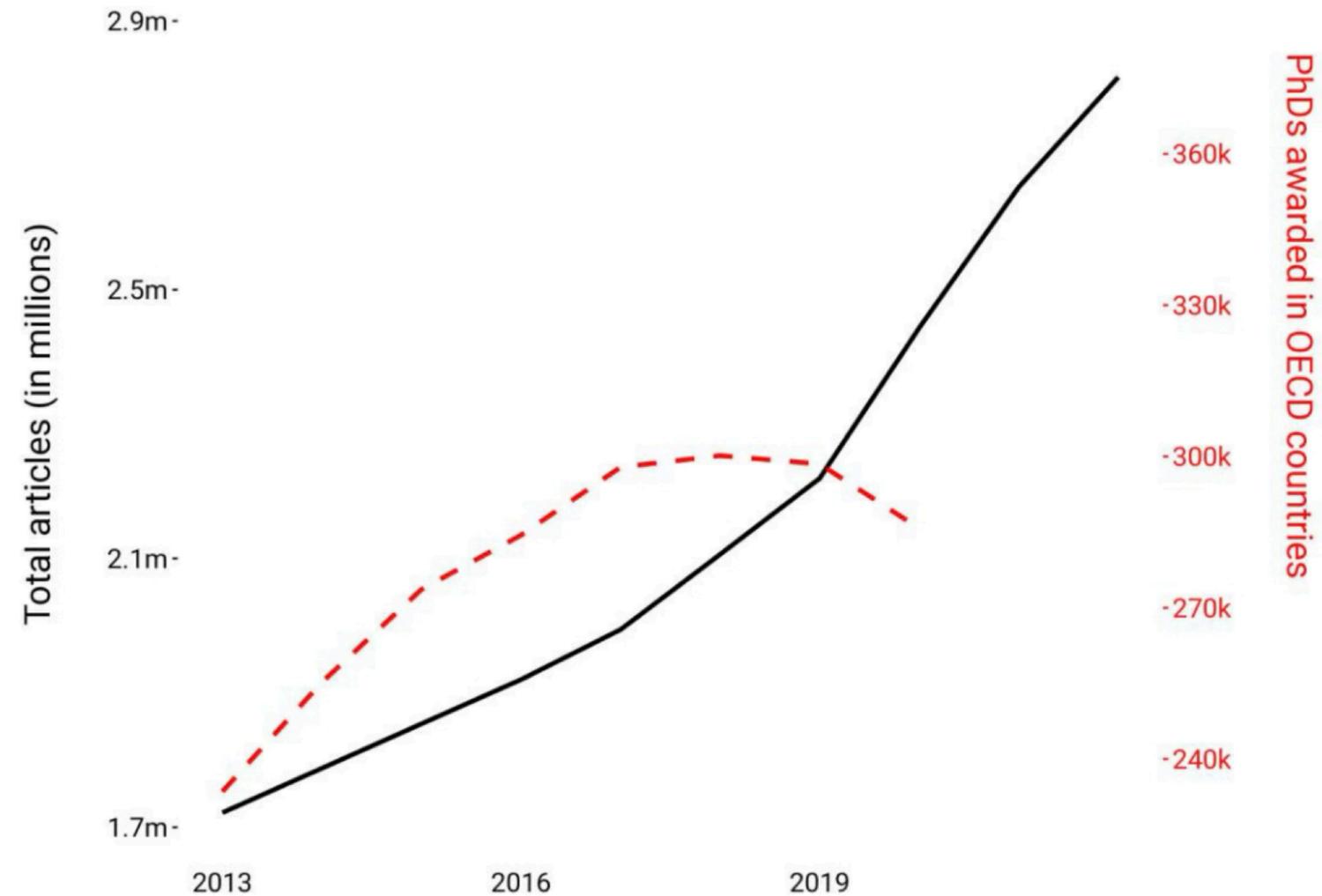
Send a message to the model...



## Task

1. Review Literature
2. Locate data sources
3. Analyse data

# Does good > bad? 2) AI slop



[Hanson et al Quan Sci Stud. 2024](#)

Human problem; how to incentive quality > quantity, globally?

## Task

1. Review Literature
2. Locate data sources
3. Analyse data

# Does good > bad? 3) Environmental impact

1) Do we have accurate (full lifecycle) estimates? 2) what's the comparison?

Measuring the environmental impact of delivering AI at Google Scale



Cooper Elsworth, Keguo Huang, David Patterson, Ian Schneider, Robert Sedivy, Savannah Goodman, Ben Townsend, Parthasarathy Ranganathan, Jeff Dean, Amin Vahdat, Ben Gomes, James Manyika

for the median Gemini Apps text prompt over one year. We identify that the median Gemini Apps text prompt uses less energy than watching nine seconds of television (0.24 Wh) and consumes the equivalent of five drops of water (0.26 mL). While these impacts are low compared to other daily activities, reducing the environmental

3) Is it worth it? Balance impact vs benefit (AGI, ASI -> science, health, education)

Duty to advance research....  
Threats e.g., to funding, response rates, stalling discovery...

# Conclusion

- 1. Quantitative research is mostly cognitive work**
- 2. AI is becoming highly capable and useful**

**AI is the future - even if progress plateaus.**

- Overdetermined given capability, influence on science, education, society
- Profound change, at infancy: we should adapt + contribute
- Understand impact, capabilities, responsible use + use AI to help learn
- Accelerate research: more efficient, reproducible, ambitious, enjoyable

## Thanks (feedback):



- **Tom Berman**
- **Liam Wright**
- **Snehal Pinto Pereira**
- **Ben Weidmann**
  
- Centre for Longitudinal Studies, Social Research Institute, IOE / UCL
- ARC (Yevgeniya Kovalchuk, Ed Lowther, Mack Nixon, Niklas Loynes)
- Patrick Sturgis for slide 6 idea
- NCRM + this event (Gabriele Durrant, Penny White; Peter Tennant)



**david.bann@ucl.ac.uk**

Funders: ESRC, UKRI Digital Research Infrastructure Programme

**Conflicts of interest:** none

# In case of interest

## Generative AI Tools for Quantitative Research: A Practical Guide

Presenter(s): [David Bann and Liam Wright](#)



Generative Artificial Intelligence (GenAI)—systems that produce text, code, images and more—has advanced rapidly in recent years. Large language models (LLMs), like ChatGPT, have surged in capability.

Below we outline practical guidance for incorporating GenAI models into quantitative research workflows, focusing on coding. This should be useful to quantitative social scientists, but also to health researchers or data scientists more broadly. 'Supporting Materials' section includes links to further reading and guides on creating web visualisation tools and using local LLMs.

You may have already accessed services like ChatGPT through a web browser, copying and pasting syntax, but these can now be directly integrated into coding tools (e.g., Interactive Development Environments [IDEs] like [VS Code](#), [Positron](#), or [RStudio](#)).

<https://www.ncrm.ac.uk/resources/online/all/?id=20859>

## Why can't Epidemiology be automated (yet)?

[David Bann](#), [Ed Lowther](#), [Liam Wright](#), [Yevgeniya Kovalchuk](#)

<https://github.com/edlowther/automated-epidemiology>

International Journal of  
**Epidemiology**

## Methods Futures Briefing #006

Series edited by Robert Meckin and Mark Elliot

## Artificial Intelligence

By [David Bann](#)<sup>1</sup> and [Liam Wright](#)<sup>1</sup>

<https://www.ncrm.ac.uk/resources/futures/>



# Appendix slides

Autonomy Level	Example Systems	Unlocking AGI Level(s)	Example Risks Introduced
<b>Autonomy Level 0: No AI</b> <i>human does everything</i>	Analogue approaches (e.g., sketching with pencil on paper)  Non-AI digital workflows (e.g., typing in a text editor; drawing in a paint program)	No AI	n/a (status quo risks)
<b>Autonomy Level 1: AI as a Tool</b> <i>human fully controls task and uses AI to automate mundane sub-tasks</i>	Information-seeking with the aid of a search engine  Revising writing with the aid of a grammar-checking program  Reading a sign with a machine translation app	Possible: Emerging Narrow AI  Likely: Competent Narrow AI	de-skilling (e.g., over-reliance)  disruption of established industries
<b>Autonomy Level 2: AI as a Consultant</b> <i>AI takes on a substantive role, but only when invoked by a human</i>	Relying on a language model to summarize a set of documents  Accelerating computer programming with a code-generating model  Consuming most entertainment via a sophisticated recommender system	Possible: Competent Narrow AI  Likely: Expert Narrow AI; Emerging AGI	over-trust  radicalization  targeted manipulation
<b>Autonomy Level 3: AI as a Collaborator</b> <i>co-equal human-AI collaboration; interactive coordination of goals &amp; tasks</i>	Training as a chess player through interactions with and analysis of a chess-playing AI  Entertainment via social interactions with AI-generated personalities	Possible: Emerging AGI  Likely: Expert Narrow AI; Competent AGI	anthropomorphization (e.g., parasocial relationships)  rapid societal change
<b>Autonomy Level 4: AI as an Expert</b> <i>AI drives interaction; human provides guidance &amp; feedback or performs sub-tasks</i>	Using an AI system to advance scientific discovery (e.g., protein-folding)	Possible: Exceptional Narrow AI  Likely: Expert AGI	societal-scale ennui  mass labor displacement  decline of human exceptionalism
<b>Autonomy Level 5: AI as an Agent</b> <i>fully autonomous AI</i>	Autonomous AI-powered personal assistants <i>(not yet unlocked)</i>	Likely: Exceptional AGI; ASI	misalignment  concentration of power

# The misalignment of incentives in academic publishing and implications for journal reform

Jennifer S. Trueblood, David B. Allison, [...] Andrei R. Teodorescu

e.g.,

**Take publishing back.** Non-profit, scientist-run journals.

**Change evaluations.** Move beyond impact factor: reward rigor, openness, contribution

**Evaluate changes + experiment** e.g., game theory

# Association of Birth Weight with Body Mass Index at Age 42 in a Longitudinal Cohort

## Introduction

Birth weight has long been studied as an early-life predictor of adult health. The “developmental origins of health and disease” hypothesis (Barker hypothesis) proposes that an adverse intrauterine environment, often reflected by low birth weight (LBW), can predispose individuals to chronic conditions in later life <sup>1</sup>. Indeed, some studies report that LBW is associated with higher risks of obesity, insulin resistance, and type 2 diabetes in adulthood <sup>2</sup>. Conversely, high birth weight (e.g. >4 kg) has been linked to greater body mass index (BMI) and higher odds of overweight or obesity in childhood and adult life <sup>3</sup>. For example, a large pooled analysis of twin cohorts found that each 1 kg increase in birth weight was associated with roughly a 0.5–0.9 kg/m<sup>2</sup> higher BMI in later life ( $p < 0.001$ ) <sup>4</sup>. These findings suggest a positive linear tracking of body size from birth to adulthood.

However, the relationship is not entirely straightforward. Some evidence indicates a **J-shaped** or **U-shaped** association, where *both* low and high birth weights confer higher risk of adverse outcomes. A systematic review noted that individuals at the extremes of birth weight tend to have elevated odds of adult overweight/obesity <sup>5</sup>. In particular, while high birth weight infants are more likely to become obese, the impact of low birth weight on adult adiposity has been less clear and might depend on other factors <sup>3</sup>. It has been suggested that the association between low birth weight and adult obesity could be **obscured by sex differences** <sup>6</sup>. Males and females have different fetal growth patterns and metabolic responses; indeed, famine and cohort studies have found that in utero undernutrition affects later adiposity in a sex-specific manner <sup>6</sup>. For instance, one study observed that low birth weight was associated with *lower* adult BMI in men but a tendency toward *higher* adiposity in women <sup>7</sup>. This underscores the importance of analyzing sex-specific effects when examining birth weight and obesity outcomes.

STROBE Statement—Checklist of items that should be included in reports of *cohort studies*

	Item No	Recommendation	Page No
Title and abstract	1	(a) Indicate the study’s design with a commonly used term in the title or the abstract	1
		(b) Provide in the abstract an informative and balanced summary of what was done and what was found	n/a
<b>Introduction</b>			
Background/rationale	2	Explain the scientific background and rationale for the investigation being reported	1
Objectives	3	State specific objectives, including any prespecified hypotheses	1