



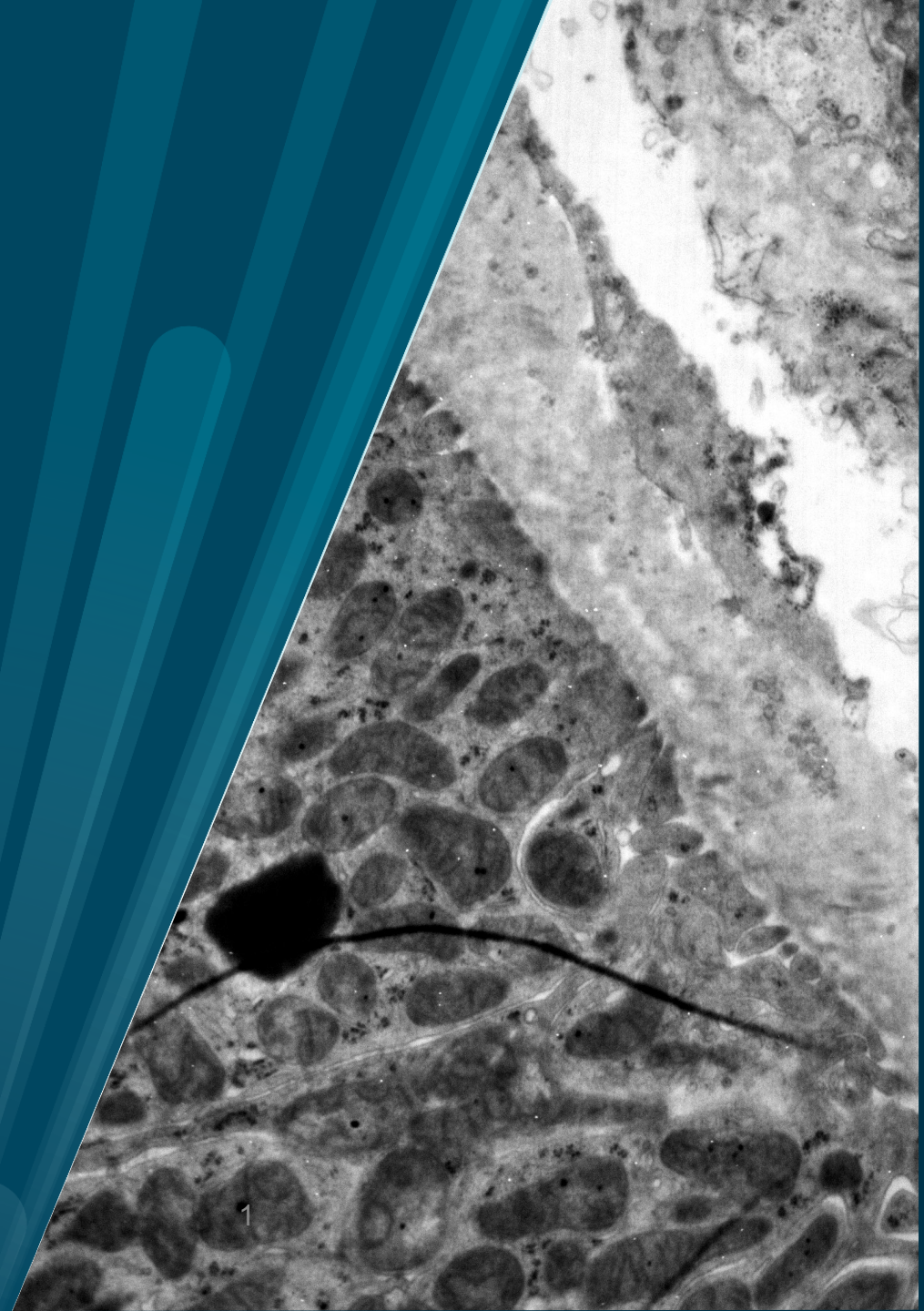
UiT The Arctic University of Norway

Code, Chaos, and Collaboration

The Quest for Reproducible AI Research

*Elisabeth Wetzer, Associate Professor in Machine Learning
Machine Learning Group*

elisabeth.wetzer@uit.no



Machine Learning Group



- machine-learning.uit.no
- visual-intelligence.no
- integreat.no

What is Reproducibility?

Reproducibility

External researchers can validate an experiment by following its documentation— either using the original data/ code or reimplementing it — to confirm the results and implementation correctness beyond mere repeatability.

Replicability

Independent researchers confirm the same hypothesis with intentionally varied implementations or using fundamentally different experimental designs— showing the findings hold across altered conditions.

Repeatability

Repeatability is the ability of the original researchers to rerun an experiment with the same setup, procedures, and analysis and obtain the same results and conclusions.

The Reproducibility Crisis

Is the widespread failure to reliably replicate published machine learning results.

Erodes trust and credibility

If results can't be reliably reproduced, researchers, funders, and the public lose confidence in published findings and scientific claims

Wastes time and resources and slows progress

Irreproducible work forces others to chase false leads, repeat experiments, or build on unstable foundations

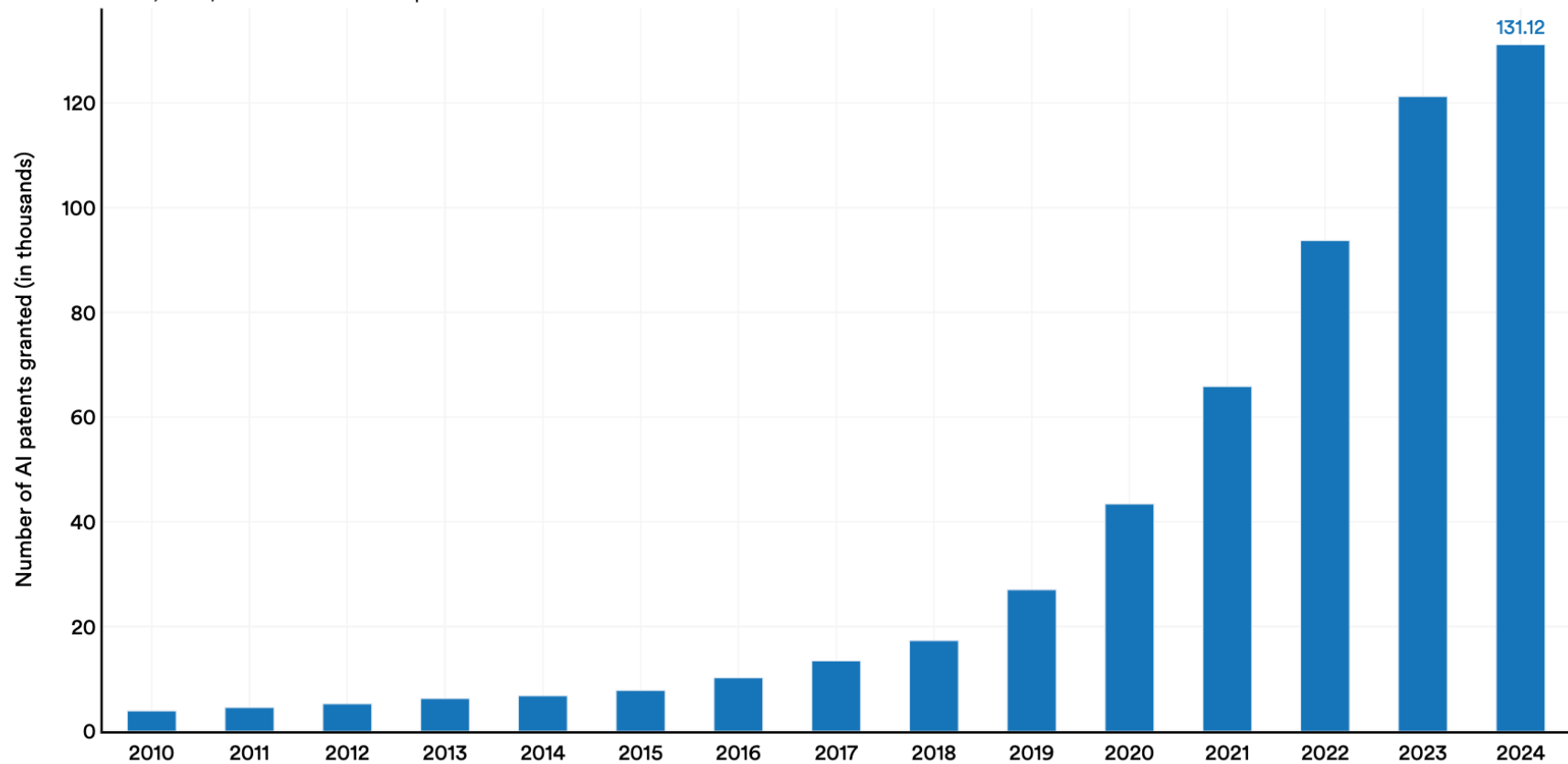
Misguides priorities and harms real-world impact

Incentives that reward flashy but fragile results skew research agendas

Machine Learning Explosion

Number of AI patents granted worldwide, 2010–24

Source: AI Index, 2026 | Chart: 2026 AI Index report

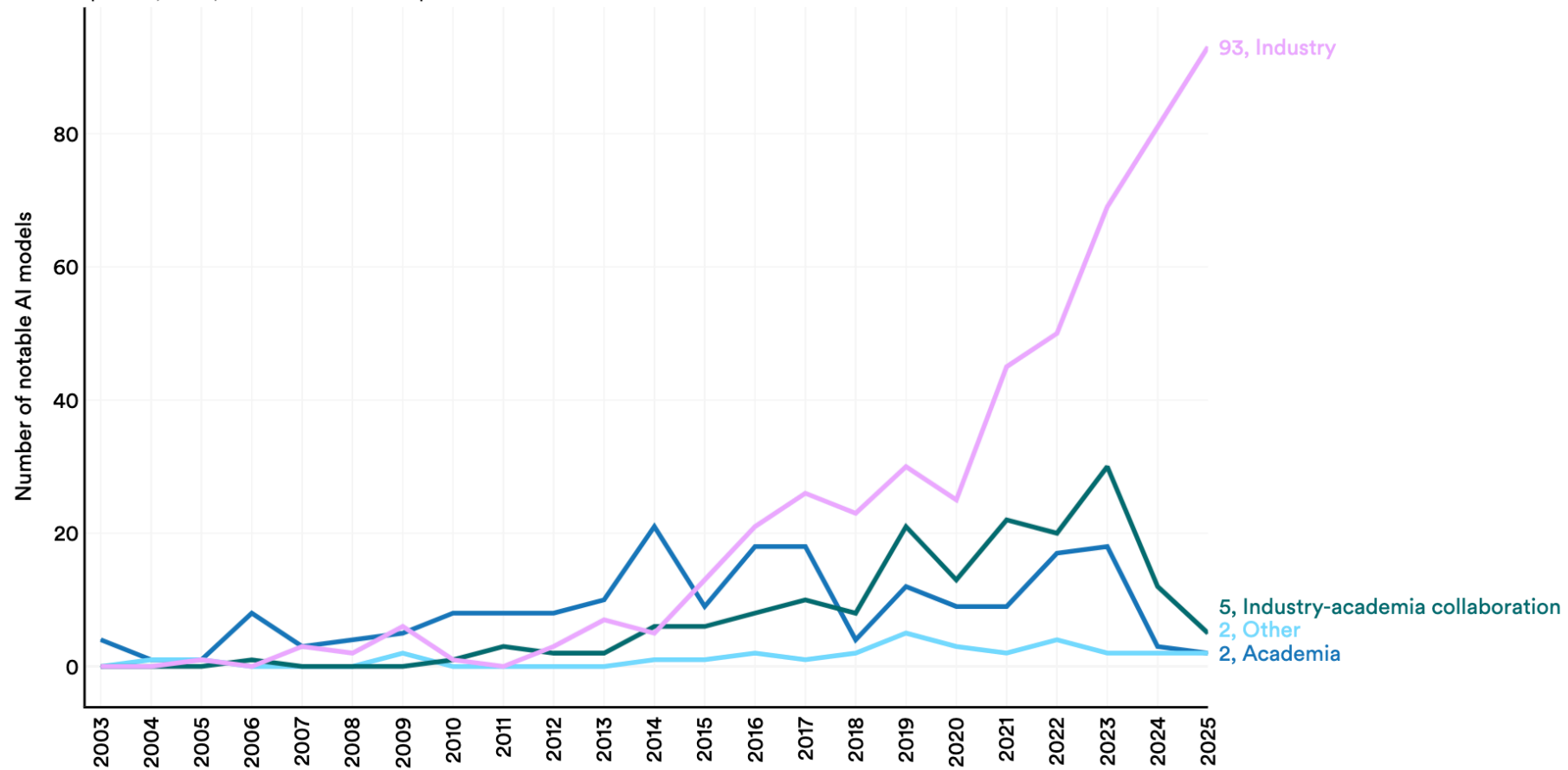


AI Index Report
2026,
Human-Centered
AI, Stanford
University

Machine Learning Explosion

Number of notable AI models by sector, 2003–25

Source: Epoch AI, 2026 | Chart: 2026 AI Index report



AI Index Report
2026,
Human-Centered
AI, Stanford
University

Machine Learning Explosion

Submission Trend

+136.5%

Last 5 events (2021-2025): 9,122 → 21,575 submissions (+136.5%).



Yearly Statistics

YEAR	LOCATION	ACCEPTED	SUBMITTED	ACCEPTANCE RATE	SECOND TRACK	NOTE
2025	San Diego, California, USA	5,290	21,575	24.5%		
2024	Vancouver, Canada	4,043	15,671	25.8%		
2023	New Orleans, USA	3,218	12,343	26.1%		
2022	New Orleans, USA	2,672	10,411	25.7%		
2021	Online, Online	2,334	9,122	25.6%		
2020	Online, Online	1,898	9,467	20.0%		
2019	Vancouver, Canada	1,428	6,743	21.2%		
2018	Montreal, Canada	1,011	4,856	20.8%		
2017	Long Beach, California, USA	678	3,240	20.9%		
2016	Barcelona, Spain	568	2,406	23.6%		

Machine Learning Explosion



In 2026, the number of submissions likely hit over 40 000

Machine Learning Explosion

Submission Trend

+136.5%

Last 5 events (2021-2025): 9,122 → 21,575 submissions (+136.5%).



NEURAL INFORMATION
PROCESSING SYSTEMS

Yearly Statistics

YEAR	LOCATION	ACCEPTED	SUBMITTED	ACCEPTANCE RATE	SECOND TRACK	NOTE
2025	San Diego, California, USA	5,290	21,575	24.5%		
2024	Vancouver, Canada	4,043	15,671	25.8%		
2023	New Orleans, USA	3,218	12,343	26.1%		
2022	New Orleans, USA	2,672	10,411	25.7%		
2021	Online, Online	2,334	9,122	25.6%		
2020	Online, Online	1,898	9,467	20.0%		
2019	Vancouver, Canada	1,428	6,743	21.2%		
2018	Montreal, Canada	1,011	4,856	20.8%		
2017	Long Beach, California, USA	678	3,240	20.9%		
2016	Barcelona, Spain	568	2,406	23.6%		

<https://openaccept.org/c/ai/neurips/>

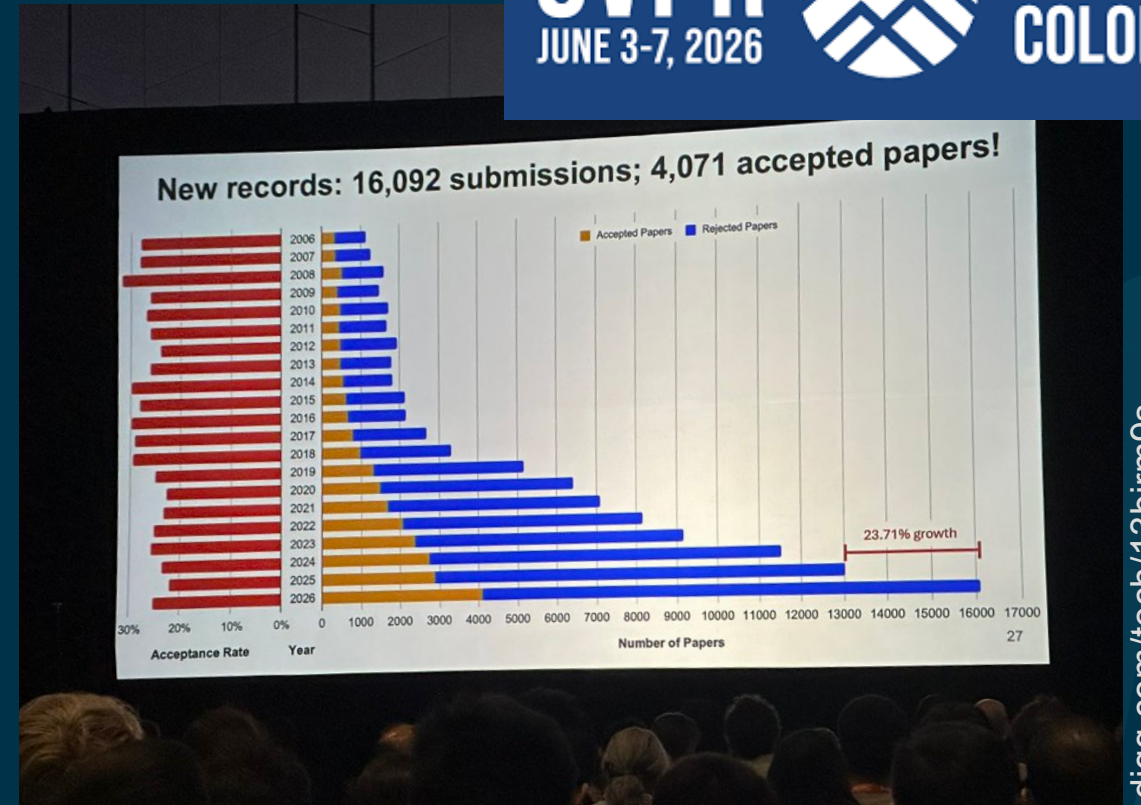
CVPR

JUNE 3-7, 2026



DENVER

COLORADO



digg.com/tech/13bjrm9s

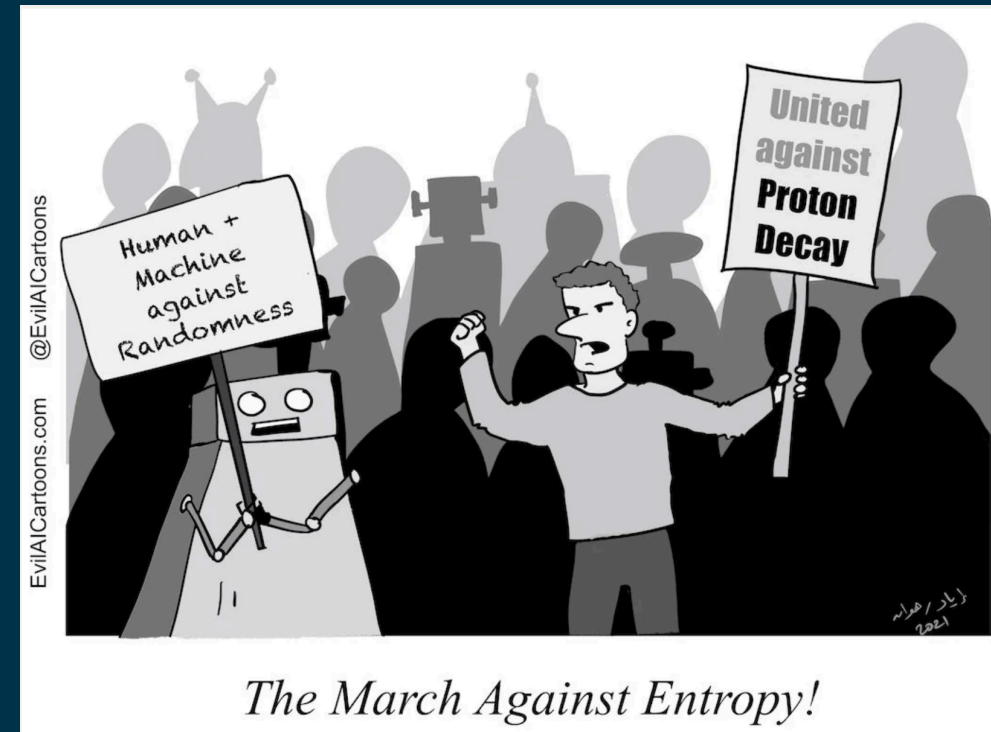
Machine Learning Explosion

- It is the same trend across all the big ML conferences:
ICLR, ICML, NeurIPS, AAAI, ECCV, ICCV, CVPR,...
- The number of available reviewers does not scale with this trend

Machine Learning Explosion

Result

- ▶ Lots of noise in reviewing
- ▶ We need quality insurance beyond the review process



Machine Learning Explosion

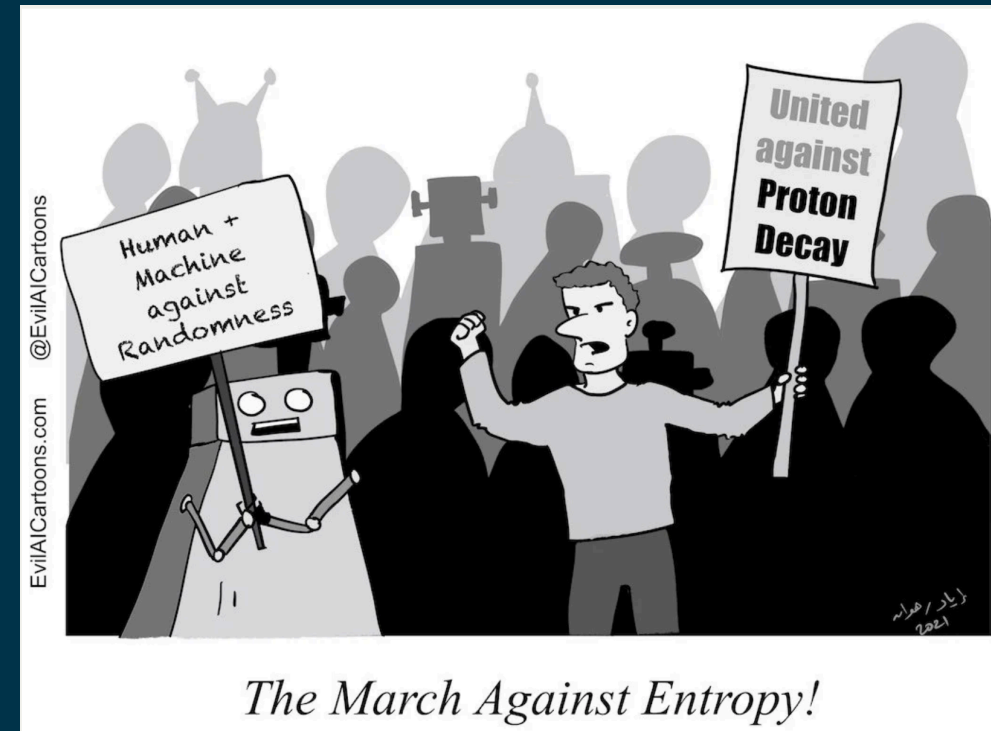
Result

- ▶ Lots of noise in reviewing
- ▶ We need quality insurance beyond the review process

NeurIPS Consistency Experiment:

23% of submissions received different decisions

50.6% of accepted papers would have been rejected



Machine Learning Explosion

- We have a problem far beyond reproducibility
- Nonetheless, if scientific fields like medicine, psychology or other experimental fields can operate under uncontrollable factors like patient diversity and massive experimental costs, we should be able to perform sound research



What makes for a good ML paper?



What makes for a good ML paper?

**New method addressing
knowledge gap**



What makes for a good ML paper?

New method addressing knowledge gap



Apply method to a handful of benchmark datasets



What makes for a good ML paper?

New method addressing knowledge gap



Apply method to a handful of benchmark datasets



Compare to 3-4 competing methods



What makes for a good ML paper?

New method addressing knowledge gap



Apply method to a handful of benchmark datasets



Compare to 3-4 competing methods



What makes for a good ML paper?

New method addressing knowledge gap



Apply method to a handful of benchmark datasets



Compare to 3-4 competing methods



Implementation Challenge

- Alternative I:
No code available; method is described in paper


Implementation Challenge

- Alternative I:
No code available; method is described in paper
 - Page limits imposed by conferences and journals
 - Default settings often not reported, but may change over time
 - A lot of additional time used for implementation, leaves no resources for fine-tuning

Implementation Challenge

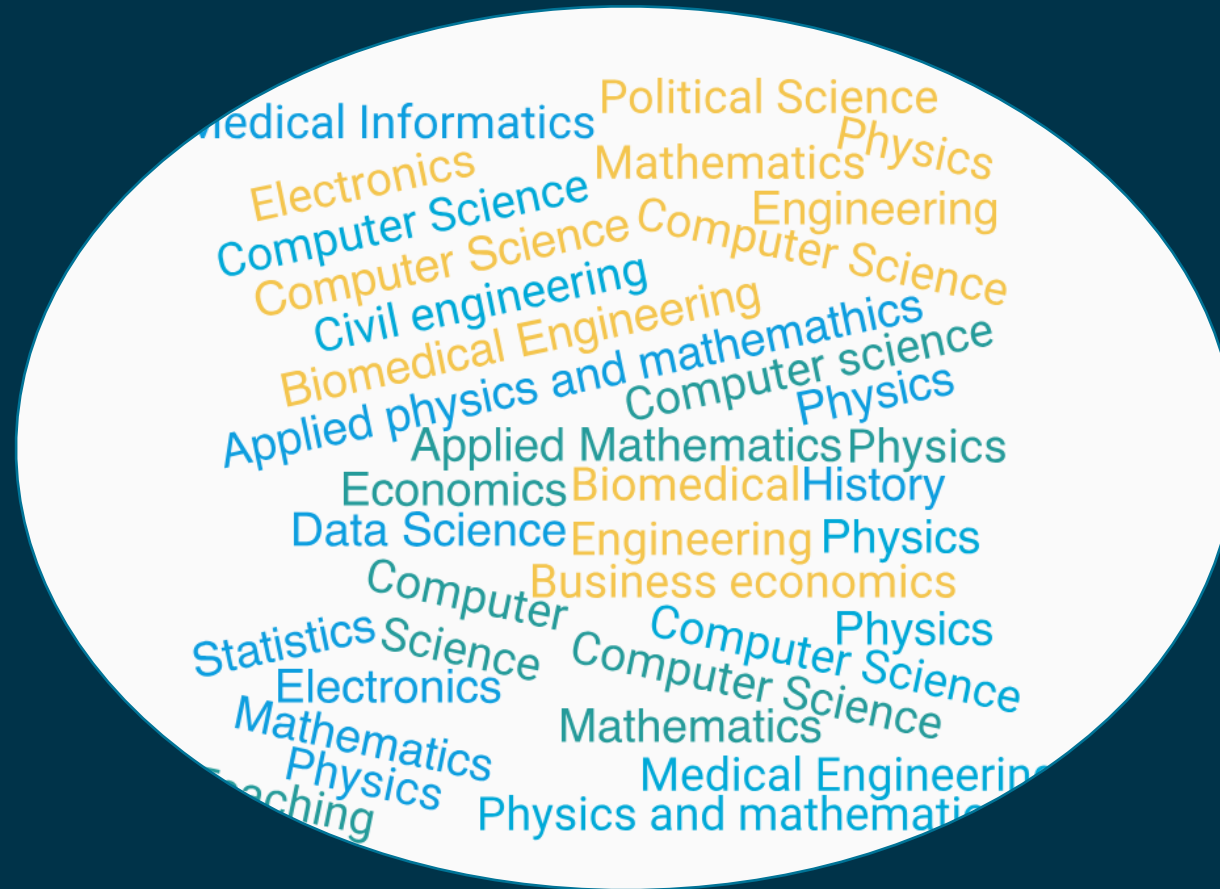
- Alternative I:
No code available; method is described in paper
- Alternative II:
Code repo is available; but takes significant work to make it work

Implementation Challenge

- Alternative I:
No code available; method is described in paper
- Alternative II:
Code repo is available; but takes significant work to make it work
- Alternative III:
Code repo is available and actually usable 

Why is that?

ML as a diverse field



ML as a diverse field

- Many come from adjacent fields such as signal processing, mathematics or physics
- About 30% had more than two codings courses in their life
- Most reported being self-taught programmers or having had one course in programming, but no training in best coding practices

ML as a diverse field

- Some seniors who ended up in ML may not be passionate programmers
- Underestimate the time and effort that goes into code clean-up
- Don't value research time spent on software publishing

```
Course name: j
# A Tic Tac Toe code Python to help a user and
import os
Score_File = "scores.txt"

# Function for handling score
def load_scores():
    if not os.path.exists(Score_File):
        return {"X": 0, "O": 0, "Draw": 0}

    Scores = {"X": 0, "O": 0, "Draw": 0}
    with open(Score_File, "r") as f:
        for line in f:
            key, value = line.strip().split("=")
            scores[key] = int(value)
    return scores

def save_scores(scores):
    with open(Score_File, "w") as f:
        for key in scores:
            f.write(f"{key}={scores[key]}\n")

def show_scores(scores):
    print("\n == Score Board ==")
    print(f"Player X wins: {scores['X']}")
    print(f"Player O wins: {scores['O']}")
    print(f"Draws: {scores['Draw']}")
    print(" == == == == == \n")
```

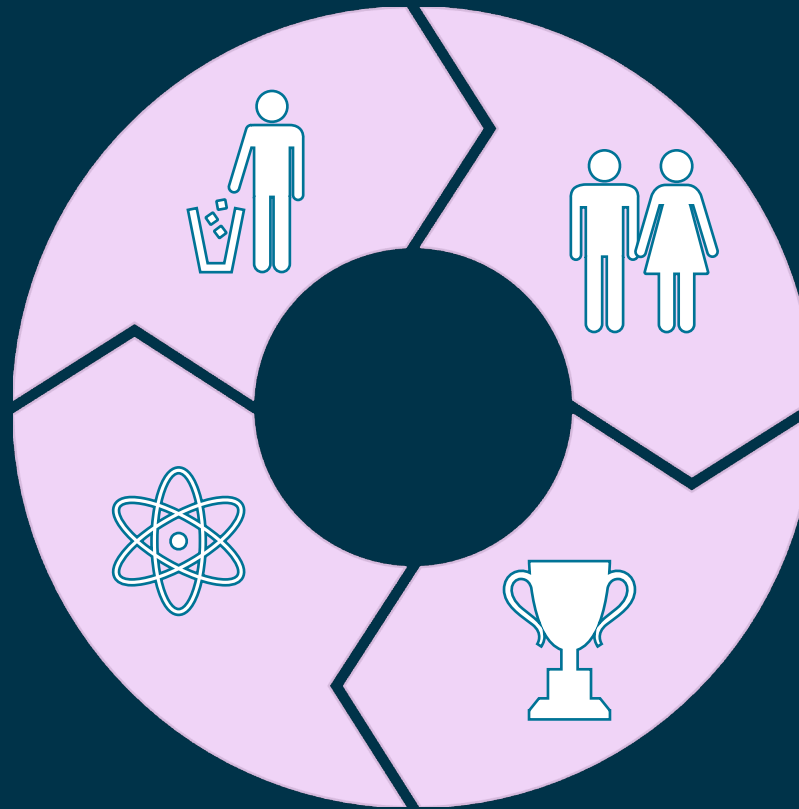
Reasons you should care about your research code

Efficiency

We have a responsibility to not waste resources

Scientific Integrity

It should be in your interest to have your findings challenged if your arguments or experiments were flawed



Collaborations

Good code can spark community contributions and collaborations

Impact

You want to be the competing method everyone picks because they do not have to spend weeks getting to run

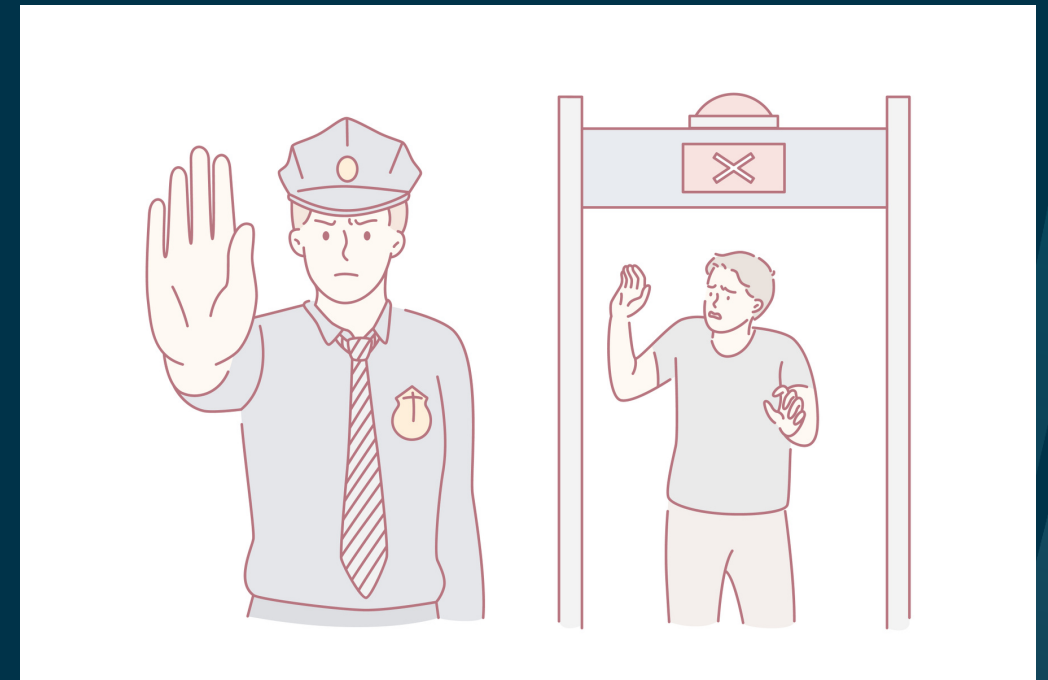
The purpose is not to keep tabs on you

```
README
```

All the libraries and the python version installed are listed below. All of them are not necessary, the core ones are e2cnn and pytorch.

Package Version

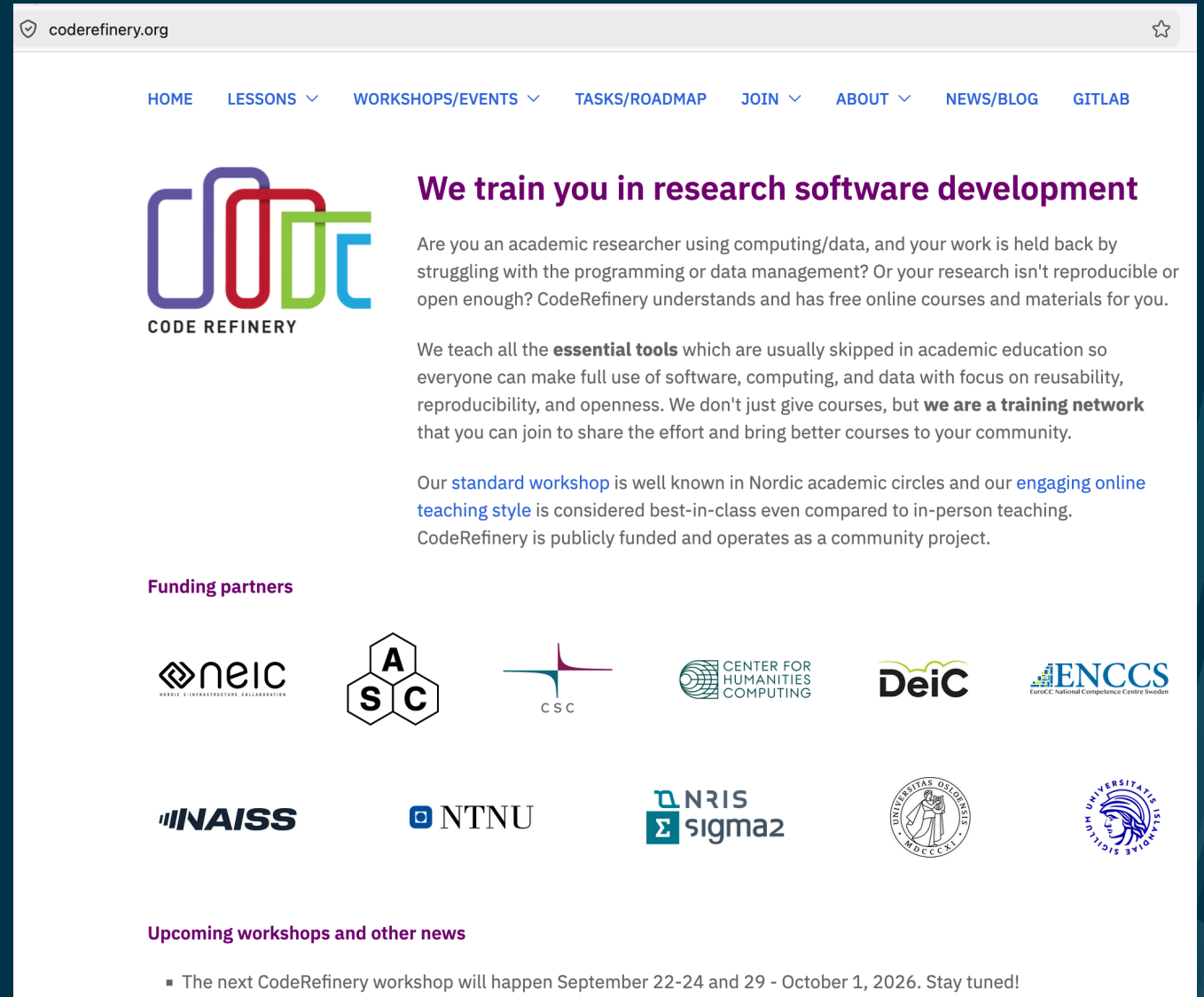
alabaster 0.7.12
albumintations 0.5.2
appdirs 1.4.4
astroid 2.3.3
atomicwrites 1.3.0
attrs 19.3.0
autopep8 1.5
Babel 2.8.0
backcall 0.1.0
bleach 3.1.0
cachey 0.2.1
certifi 2019.11.28 cffi 1.14.0
chardet 3.0.4
cloudpickle 1.3.0
cryptography 2.9
cyclcr 0.10.0
dask 2021.2.0
decorator 4.4.1
defusedxml 0.6.0
diff-match-patch 20181111
docstring-parser 0.7.3
docutils 0.16
e2cnn 0.1
entrypoints 0.3
flake8 3.7.9
freetype-py 2.1.0
HeapDict 1.0.1



So how do we get better in this?

Increase RSE literacy

- Code Refinery offers plenty of learning material, courses and workshops



The screenshot shows the homepage of coderefinery.org. The navigation menu includes: HOME, LESSONS, WORKSHOPS/EVENTS, TASKS/ROADMAP, JOIN, ABOUT, NEWS/BLOG, and GITLAB. The main heading is "We train you in research software development". Below this, there is a paragraph explaining the organization's mission: "Are you an academic researcher using computing/data, and your work is held back by struggling with the programming or data management? Or your research isn't reproducible or open enough? CodeRefinery understands and has free online courses and materials for you." This is followed by another paragraph: "We teach all the **essential tools** which are usually skipped in academic education so everyone can make full use of software, computing, and data with focus on reusability, reproducibility, and openness. We don't just give courses, but **we are a training network** that you can join to share the effort and bring better courses to your community." A third paragraph states: "Our [standard workshop](#) is well known in Nordic academic circles and our [engaging online teaching style](#) is considered best-in-class even compared to in-person teaching. CodeRefinery is publicly funded and operates as a community project." The "Funding partners" section displays logos for: neic (Academic Infrastructure Collaboration), ASC (Academic Software Collaboration), CSC (Center for Scientific Computing), Center for Humanities Computing, DeIC (Digital Education Innovation Centre), and ENCCS (EuroCC National Competence Centre Sweden). The second row of logos includes: NAISS (National Academic Infrastructure Support Service), NTNU (Norwegian University of Science and Technology), NRIS sigma2 (National Research Infrastructure for Scientific Computing), and logos for the University of Oslo and the University of Århus. The "Upcoming workshops and other news" section contains a bullet point: "The next CodeRefinery workshop will happen September 22-24 and 29 - October 1, 2026. Stay tuned!"

Increase RSE literacy

- RSE Group at UiT offers drop-in help desk every week



UiT Norges arktiske universitet

Research Software Engineering

Contact us

Email

Send an email to rse@uit.no and we will reply as soon as we can.

Where to find us on the campus during office hours

We are available almost every [Wednesday from 14:00-16:00](#) at our Help Desk sessions. Click [here](#) for location and scheduling information.

Increase RSE literacy

- Increase Reward and Appreciation in the Field



The screenshot shows the homepage of The Journal of Open Source Software (JOSS). At the top left is the JOSS logo, a blue gear with a white center, followed by the text "The Journal of Open Source Software". A hamburger menu icon is in the top right. Below the header is a light blue notification box with a close button (X) in the top right corner. The notification text reads: "Important Update: JOSS has updated its submission scope requirements, affecting what is eligible for submission and what information is required in your paper. **Read the announcement** → • **View updated requirements** →". Below the notification, the main text states: "The Journal of Open Source Software is a **developer friendly**, open access journal for research software packages." This is followed by a subtext: "Committed to publishing quality research software with zero article processing charges or subscription fees." At the bottom, there are two dark blue buttons: "Submit a paper to JOSS" and "Volunteer to review" with a small icon of a person and a gear.

Who are our repos for?

- If we all use Claude & co to implement the competing method for our setting
 - ... should we consider our paper and code well enough documented and reproducible if Claude can replicate our results?

Who are our repos for?

- If we all use Claude & co to implement the competing method for our setting
 - ... should we consider our paper and code well enough documented and reproducible if Claude can replicate our results?
 - ... should this be a service journals provide to test min. requirement?

Trends

- Agentic Reproducibility

Preprint. Under review.

Read the Paper, Write the Code: Agentic Reproduction of Social-Science Results

Benjamin Kohler
ETH Zurich

David Zollikofer
ETH Zurich

Johanna Einsiedler
University of Basel

Alexander Hoyle*
ETH Zurich

Elliott Ash*
ETH Zurich

Abstract

Recent work has used LLM agents to reproduce empirical social science results with access to both the data and code. We broaden this scope by asking: Can they reproduce results given only a paper’s methods description and original data? We develop an agentic reproduction system that extracts structured methods descriptions from papers, runs reimplementations under strict information isolation—agents never see the original code, results, or paper—and enables deterministic, cell-level comparison of reproduced outputs to the original results. An error attribution step traces discrepancies through the system chain to identify root causes. Evaluating four agent scaffolds and four LLMs on 48 papers with human-verified reproducibility, we find that agents can largely recover published results, but performance varies substantially between models, scaffolds, and papers. Root cause analysis reveals that failures stem both from agent errors and from underspecification in the papers themselves.

Trends

- Agentic Reproducibility

Preprint. Under review.

Read the Paper, Write the Code: Agentic Reproduction of Social-Science Results

Benjamin Kohler
ETH Zurich

David Zollikofer
ETH Zurich

Johanna Einsiedler
University of Basel

Alexander Hoyle*
ETH Zurich

Elliott Ash*
ETH Zurich

Abstract

Recent work has used LLM agents to reproduce empirical social science results with access to both the data and code. We broaden this scope by asking: Can they reproduce results given only a paper's methods descrip-

tion and code? We build a reproduction system that takes a paper, runs reimplementations, and performs a cell-level comparison of results. Our attribution step traces results back to root causes. Evaluating agents with human-verified reproductions of published results, but without scaffolds, and papers. We find that agents struggle with scaffolds, and papers. We find that agents struggle with scaffolds, and papers.

REPRO-Bench: Can AI agents Automate Research Reproducibility Assessments?



Daniel Kang

Follow

5 min read · Jul 28, 2025



Trends

- Agentic Reproducibility

ARA: Agentic Reproducibility Assessment For Scalable Support Of Scientific Peer-Review

Kevin Riehl¹

ETH Zürich, IVT & Agentic Systems Lab (ASL) kriehl@ethz.ch

&Andres L. Marin¹

European Commission, Joint Research Centre & University of Konstanz

E-Mail: andres.laverde-marin@ec.europa.eu

&Nikiforos Zacharof Ideas Forward

nikiforos.zacharof@ideasforward.com

&Fan Wu ETH Zürich, ASL E-Mail: fanwu@ethz.ch

&Patrick Langer ETH Zürich, ASL E-Mail: planger@ethz.ch

&Robert Jakob ETH Zürich, ASL E-Mail: rjakob@ethz.ch

&Anastasios Kouvelas ETH Zürich, IVT kouvelas@ethz.ch

&Georgios Fontaras European Commission, Joint Research Centre

georgios.fontaras@ec.europa.eu

&Michail A. Makridis ETH Zürich, IVT mmakridis@ethz.ch

Preprint. Under review.

Read the Paper, Write the Code: Agentic Reproduction of Social-Science Results

Benjamin Kohler
ETH Zurich

David Zollikofer
ETH Zurich

Johanna Einsiedler
University of Basel

Alexander Hoyle*
ETH Zurich

Elliott Ash*
ETH Zurich

Abstract

Recent work has used LLM agents to reproduce empirical social science results with access to both the data and code. We broaden this scope by asking: Can they reproduce results given only a paper's methods descrip-

REPRO-Bench: Can AI agents Automate Research Reproducibility Assessments?



Daniel Kang

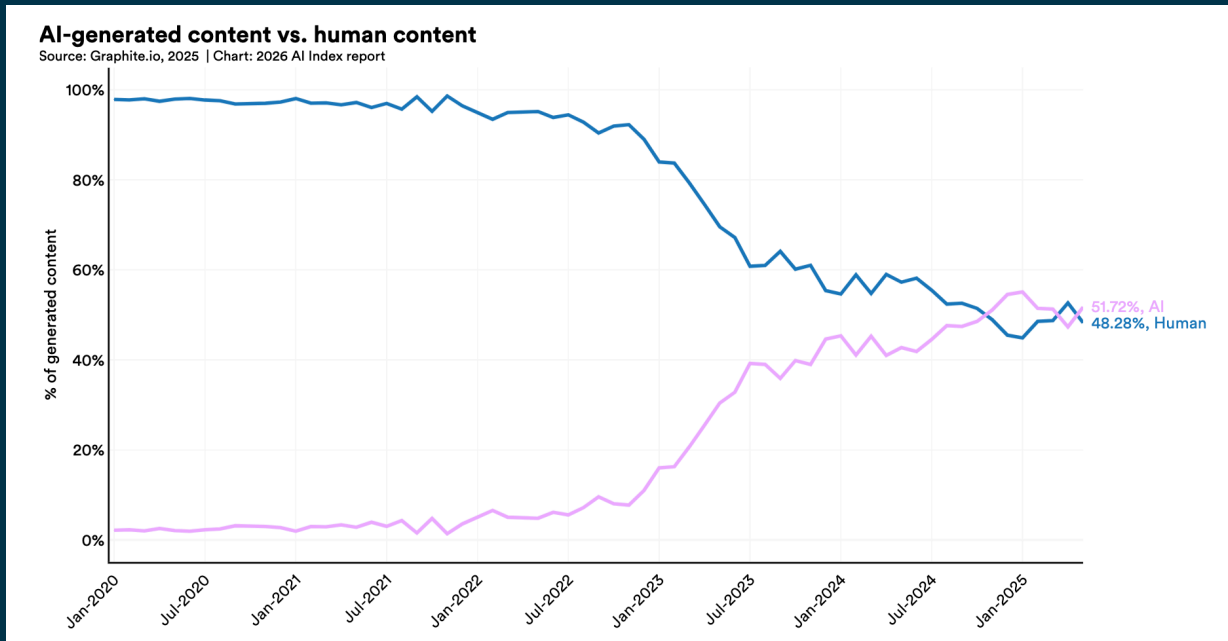
Follow

5 min read · Jul 28, 2025

tion, runs a reproduction system that papers, runs reimplementations that never see the original data, cell-level comparison of results, or attribution step traces back to root causes. Evaluating agents with human-verified references, but not published results, but also scaffolds, and papers. This is from agent errors and is.



What will this mean for coding practices?



AI Index Report 2026,
Human-Centered AI, Stanford University

Key Take-Away

Value Reproducible Code and Open Data

- ◆ Teach students best coding practices, not only CS students
- ◆ Reward time spent on code clean up and software publishing
- ◆ Introduce minimal requirements on code usability

Thank you!

Contact:

elisabeth.wetzer@uit.no