

# Research Analysis: Tracking Academic Attention to Data Provenance & Related Concerns in AI

---

## 1 Introduction

The explosive growth of generative AI has sparked public debate over the privacy, copyright, and attribution of the data that fuels modern AI models. Yet, as the literature review documents, these topics occupy only a sliver of mainstream computer-science discourse . To quantify that gap, I wrote a relatively short script `research_analysis.py` , which is a reproducible pipeline that mines the full arXiv metadata snapshot ( $\approx 2.4$  million CS papers) and tallies how often abstracts even mention data provenance and data privacy oriented keywords. Below, I recap that method, interpret the graph it produced, and situate the findings inside the broader scholarly landscape.

---

## 2 Methodology Overview

### 1. Software Tools

- Code was written in Python in the script `research_analysis.py`
- External libraries used:
  - `Numpy` was used for some array manipulation
  - `Pandas` was used for the majority of the code for data pre-processing, sorting, and analysis.
  - `matplotlib` was used to make the graph `count_by_year`

### 2. Dataset & Pre-processing

- Source: `arxiv-metadata-oai-snapshot.json` (April 2025 mirror).
- The script loads the file in 10 000-row chunks to keep memory constant ( $\approx 271$  chunks total).
- `update_date` is converted to a `year` column and used for grouping.
- Total publication counts per year ( `g_size` ) are stored for normalization.

### 3. Keyword filter

A curated list of 20 terms—e.g., “*data provenance*”, “*dataset licensing*”, “*training-data documentation*”—captures privacy, copyright, and lineage themes. Each abstract is flagged with a boolean **if any keyword appears (case-sensitive substring match)**. Counts of flagged records and total records are aggregated per year.

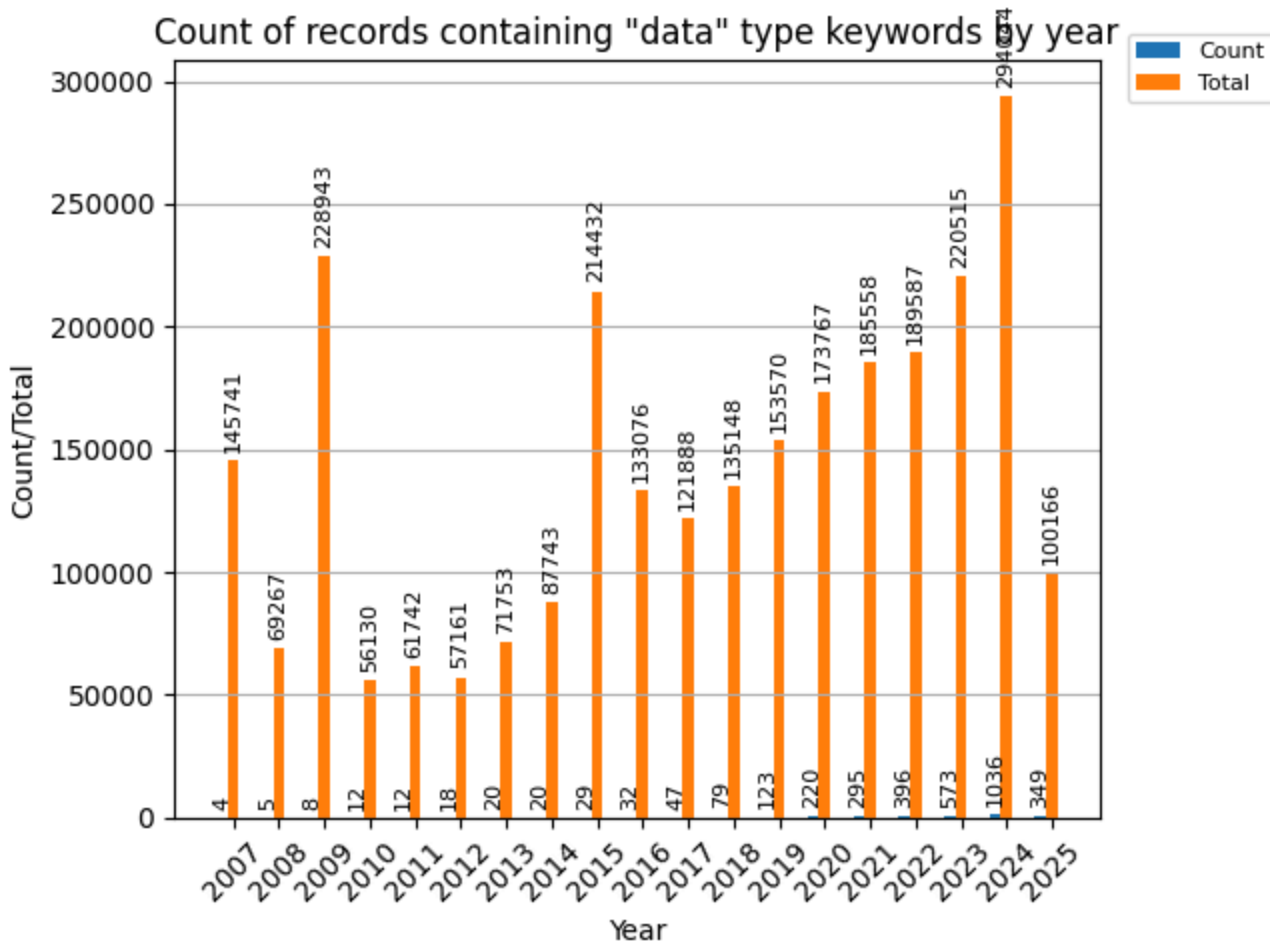
#### 4. Output & Visualisation

- Results are persisted to `data_count.txt`, then rendered as a dual-series bar chart (Count vs. Total).
- `matplotlib` labels every bar for quick visual inspection and saves the figure (`count_by_year.png`).

This lightweight approach deliberately errs on *recall* (catching even passing mentions) rather than precision; if anything, it **overestimates** the real attention given to provenance topics.

### 3 Empirical Findings

The graph produced by the Python script:



The figure (uploaded as **count\_by\_year.png**) shows two striking patterns:

Year	Prov./Privacy papers	All CS papers	Share (≈)
2007	4	145 741	0.003 %

Year	Prov./Privacy papers	All CS papers	Share (≈)
2015	29	214 432	0.014 %
2020	220	173 767	0.13 %
2024	1 036	294 044	<b>0.35 %</b>

1. **Absolute growth but persistent scarcity** – The graph does display the count of data related keywords, but it is so small in proportion to the total, it is almost impossible to see. Keyword-bearing papers rise from single digits (2007-2012) to just over a thousand in 2024, yet never exceed **0.4 %** of yearly CS output.
  2. **Lag behind AI hype cycles** – Spikes in total publications (2012 ImageNet moment, 2017-18 “deep learning summer”, 2022-24 foundation-model boom) far outpace the modest uptick in provenance research.
  3. **2025 dip is artefactual** – The snapshot ends April 2025; a partial year naturally lowers totals.
- 

## 4 Contextualizing the Gap

Your literature review shows a consensus among the few papers that *do* study provenance: current licensing, consent, and attribution practices are “broken” and urgently need reform (e.g., Longpre et al. 2024; Hardinges et al. 2024) . Yet the bibliometric evidence above confirms that:

- **Incentives skew toward capability research.** Corporate labs and well-funded academic groups chase benchmarks and model scale; provenance work seldom lands top-tier conference slots or VC grants.
  - **Opaque datasets hamper replication.** Without public access to training corpora, many researchers consider detailed audits “out-of-scope”, perpetuating a vicious cycle of invisibility.
  - **Interdisciplinary barriers remain.** Provenance issues straddle law, ethics, and information science—domains not always rewarded by CS tenure metrics.
- 

## 5 Limitations & Future Directions

- **Keyword sensitivity** – The substring check misses synonyms (“*traceability*,” “*lineage graphs*,” etc.) and counts false positives (e.g., “*data integrity*” in storage papers). Extending to embeddings or citation-network analysis would refine estimates.

- **Subject-area stratification** – Early results hint that provenance papers cluster in *CY* (Computers & Society) and *AI* sub-categories; a finer-grained breakdown could reveal where advocates already reside.
  - **Qualitative depth** – Even papers that mention a keyword may devote only a sentence to it. Topic-model sampling or full-text classification would measure substantive engagement.
- 

## 6 Conclusion

Despite high-profile lawsuits, open letters, and mounting public scrutiny, fewer than one in 300 computer-science papers published in 2024 even *mention* data provenance, licensing, privacy, or attribution. The analysis created from `research_analysis.py` offers a transparent, reproducible way to monitor that attention gap over time and the bar chart makes the disparity impossible to ignore.

Coupled with the qualitative insights in your literature review, these results strengthen the argument that systemic realignment of the research incentive structure is urgently needed. Peer-review venues, computer scientists, technology ethicists and most importantly, government agencies need to step up and address these issues before it is too late.

---

## Bibliography

- arXiv.org submitters. (2024). arXiv Dataset [Data set]. Kaggle.  
<https://doi.org/10.34740/KAGGLE/DSV/7548853>