# Big Tech-Funded AI Papers Have Higher Citation Impact, Greater Insularity, and Larger Recency Bias

Max Martin Gnewuch
*University of Göttingen*
Göttingen, Lower Saxony, Germany
maxmartin.gnewuch@stud.uni-goettingen.de

Jan Philip Wahle
*University of Göttingen*
Göttingen, Lower Saxony, Germany
wahle@uni-goettingen.de

Terry Ruas
*University of Göttingen*
Göttingen, Lower Saxony, Germany
ruas@uni-goettingen.de

Bela Gipp
*University of Göttingen*
Göttingen, Lower Saxony, Germany
gipp@uni-goettingen.de

arXiv:2512.05714v1 [cs.DL] 5 Dec 2025

*Abstract*—Over the past four decades, artificial intelligence (AI) research has flourished at the nexus of academia and industry. However, Big Tech companies have increasingly acquired the edge in computational resources, big data, and talent. So far, it has been largely unclear how many papers the industry funds, how their citation impact compares to non-funded papers, and what drives industry interest. This study fills that gap and quantifies the number of industry-funded papers at $10$ top AI conferences (e.g., ICLR, CVPR, AAAI, ACL) and their citation influence by analyzing $\approx49.8\,K$ papers, $\approx1.8\,M$ citations from AI papers to other papers, and $\approx2.3\,M$ citations from other papers to AI papers from $1998$–$2022$ in Scopus. We investigate the volume and evolution of industry funding in AI research, the citation impact of the papers, the diversity and temporal range of their citations, and the subfields in which industry predominantly acts through $7$ research questions. Our findings reveal that the industry present has grown markedly since **2015**, from less than $2\,\%$ to more than $11\,\%$ in **2020**. Between $2018$ and $2022$, $12\,\%$ of industry-funded papers achieved high citation rates as measured by h5-index, compared to $4\,\%$ of non-industry-funded papers and $2\,\%$ of non-funded papers. Top AI conferences engage more with industry-funded research than non-funded research, as measured by our newly proposed metric, the *Citation Preference Ratio* ($CPR$). We show that industry-funded research is increasingly insular — citing predominantly other industry-funded papers while referencing fewer non-funded papers. Furthermore, industry-funded work cites more recent work and fewer older papers than non-funded works. These findings reveal new trends in AI research funding, including a shift towards (1) more industry-funded papers and their growing citation impact, (2) greater insularity of industry-funded works than non-funded works, and (3) preference of industry-funded research to cite recent work. While industry funding contributes markedly to AI research, these new trends also raise questions about Big Tech's allocation of resources and potential control over research topics. All data and code are publicly available: https://github.com/Peerzival/impact-big-tech-funding.

*Keywords*—Scientometrics, Big Tech, Funding, Power Dynamics, Monopolization, Echo Chambers, Ethical AI, Bias.
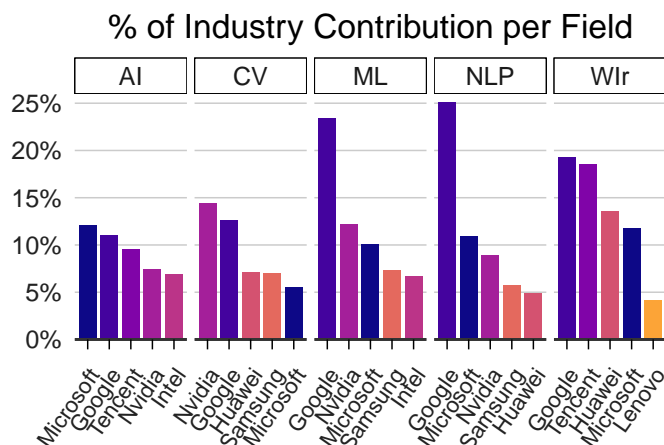
Fig. 1: The share of the top 5 companies regarding all industry funding across key domains between 2018 to 2022.

## I. INTRODUCTION

Artificial Intelligence (AI) has evolved into a "general purpose technology" [1], [2] comparable to other historic economic transformations such as the steam engine, electrification, and the Internet [2], offering new opportunities across industries [3]. This evolution can reshape society by affecting job markets [4], healthcare [5], and addressing global challenges like climate change [6]. However, integrating AI applications can lead to unexpected consequences and introduce new types of risks. For example, developing advanced weapons and autonomous military systems could trigger an arms race, escalating global tensions and undermining security [7]. If policymakers do not support displaced workers adequately, this trend may lead to economic instability [4]. Recognizing these risks alongside opportunities, the development of AI systems must be inclusive, ensuring that their benefits are widely accessible rather than

limited to a privileged few [1], [8], [9]. Research teams that are diverse in terms of age, gender, race, and ethnicity and that engage with diverse literature across time and fields are better equipped to create more democratic, fair, and unbiased AI [10], [11]. However, early evidence suggests that AI development is becoming increasingly insular [12] with few powerful actors leading innovations [1], [13], [14].

**The role of Big Tech.** The industrial landscape of AI experiences strong interest by a small number of large companies, mainly in the US and China (Figure 1), often referred to as *Big Tech* (e.g., Meta, Google, Amazon, Nvidia, Microsoft, Baidu, Tencent) [1], [2], [15]–[17]. These companies have access to 3 key resources essential for modern AI research: big data, computational power, and access to a highly skilled AI workforce [2], [14], [15], [17]. This oligopolistic/monopolistic control over critical resources could grant a disproportionate amount of power to a small number of corporations [14], [15], who contribute to shaping what is examined (and not examined) in AI research [13]. As a result, there is growing concern about scientific independence, influence, and power concentration in AI research [13], [14], [18].

**Tensions between commercial and public interests.** The goals of industry and those of the public or society can differ markedly. Connections between industry and the AI research community mean that the questions and incentives shaping the field are not always within the control of individual researchers [13]. The tech industry's incentives can influence the trajectory of a field, including which questions are deemed worth pursuing and which answers lead to grants, awards, and tenure [13]. The increasing industry investment in AI research does not diminish the potential for substantial societal benefits. However, commercial motivations often drive companies to prioritize profit-oriented goals and topics. Such incentives can but do not have to align with the public interest and result in beneficial tools, hardware, or software. Without public alternatives, AI may follow a similar pattern observed in the pharmaceutical industry, where investments can overlook the needs of those that deserve particular attention, such as lower-income groups [14], [19].

**Research Gap.** Prior work has examined Big Tech presence, including analyses of Big Tech funding for AI researchers [18], the rising share of Big Tech–affiliated publications at top conferences [1], the migration of researchers from academia to private sector and talent shifts, research focus, value differences, and impact disparities between academia and industry [12], [20]–[22]. These works rely on author affiliations instead of direct funding attribution and use a subset of publicly traded companies. Further prior work lacks reference from Big Tech funding to other funding types, such as public funding and papers without funding. By doing so, we aim to foster further discourse on the interplay between academic and industry research and advance a more nuanced understanding of their relationship.

**This study.** In this work, we systematically assess the influence of funded papers by Big Tech and public funding types on the citation practices of AI papers. Using the Scopus database, we compiled a new dataset of metadata associated with $\approx 49.8$ K AI papers published at 10 top conferences[1] from 1998 to 2022, along with $\approx 1.8$ M citations from AI papers and $\approx 2.3$ M citations to AI papers. The dataset contains both ties of authors to industry through author affiliations and extracted funding statements through acknowledgement sections of papers. In this analysis, we mainly focus on the acknowledgement sections of direct funding to a paper as they are a high-precision way of tracing direct funding. In our dataset, each citation includes details about the year of publication, the funding agency identified in the paper for both industry and non-industry funders, and the field of study for both the papers *cited* by AI research and those *citing* AI research. We use this new dataset to address different questions of industry in AI:

1) **Evolution of industry presence in AI**: How has the volume of industry-funded research evolved over time? Which companies have high publication output? What specific fields within AI attract the greatest attention from industry?

2) **Engagement of the AI community**: How much do papers from top AI venues cite industry-funded research? How has this citation behaviour changed over time? How does the citation impact of industry-funded papers compare to other AI papers? Which companies have the highest citation impact?

3) **Insularity of industry-funded research**: Does industry-funded research predominantly cite other industry-funded research instead of research from other funding sources? How much does industry-funded research cite different academic fields compared to other AI papers?

4) **Recency bias of industry-funded research**: How far back in time do industry-funded papers cite on average compared to other AI papers? Do industry-funded papers cite more recent or older work than other AI papers?

We show that the industry present has grown markedly since 2015 from less than $2\%$ to more than $11\%$ in 2020. Since 2020, the percentage of industry-funded papers has slightly declined by $14\%$, while the AI community's engagement in industry-funded research has increased. $12\%$ of industry-funded papers have a high citation impact as measured by h5-index, compared to $4\%$ of non-industry-funded and $2\%$ of non-funded papers. We provide evidence that industry-funded research is increasingly insular, citing other industry-funded research $2\%$ more often than other funding types. Furthermore, our experiments show that industry-funded research references more recent papers and cites fewer older ones than non-funded research.

These results have potential implications for the academic community. The decreasing presence of industry-funded papers could signal Big Tech increasingly using their own channels for publication while moving out of these venues, which also implies reduced funding for venues — a key funding source that many venues still rely on. Our results suggest that publishing in top AI conferences has become less vital for Big Tech to achieve high citation rates and visibility for industry-funded

---

[1]According to csranking.org

research. The disproportionate engagement with research that Big Tech is funding and the limited diversity of companies [10], [11] may affect long-term innovation. Our results suggest that a few AI companies are increasingly impacting modern AI research as measured by citations. Evidence exists that these companies may not reflect the interests of the broader population [10], [11], and our results support that papers from Big Tech prioritize recent literature disproportionately and have insular citation cycles compared to other AI papers.

## II. RELATED WORK

**Big tech influence on AI research.** Several studies have attempted to quantify Big Tech's influence in AI research using various methods and data sources. [18] revealed that $59\%$ of papers published in top AI journals addressing ethical and societal implications feature at least one author with financial ties to Big Tech. Similarly, [1] analyzed participation trends in AI conferences after the rise of deep learning in 2012, finding a marked increase in representation from major technology companies.

[20] used bibliometric data to trace researcher migration from academia to industry. Their findings indicate that $25\%$ of AI researchers at top 5 Nature Index institutions transition to industry, highlighting Big Tech's competitive offers to attract top talent.

The study settings of [23] allowed the authors to examine the link between private-sector affiliations and research impact, measured through citations and attention scores. Their analysis shows that the private sector is dominant in shaping AI research.

Work by [22] compared citation networks and memetic propagation between Big Tech and academia. Their results suggest that Big Tech-affiliated papers are disproportionately cited, with the most impactful research stemming from collaborations between academia and industry.

Several studies have tried to analyze the content of AI papers and study the influence of Big Tech. [12] study subject co-occurrence in arXiv papers, while [21] investigated funding sources, reporting an increasing presence of Big Tech.

Yet, there is a key gap in our understanding of Big Tech's influence. Most studies described above equate influence with institutional affiliation shares (with the exception of [22] and [23]). However, this approach does not consider the number of the paper's citations, which is a key proxy for the impact of a paper [22], [24]–[26]. Furthermore, the reliance on publicly traded companies excludes key non-profit research organizations (e.g., OpenAI), particularly in the fast-moving field of AI, where startups emerge fast and have a marked impact.

**Citation patterns in scientific works.** Related work has attracted significant attention to citation patterns exploring various aspects such as self-citation [27], publication venue [28], [29], paper quality [30], publication language [31], geographic location [25], gender [32]–[34], and field of study [35].

Interdisciplinary research field engagement represents a core component of responsible research [36], [37]. Numerous breakthroughs have emerged from such interactions, including Einstein's photoelectric effect [38] and Bohr's atomic model [39]. Similarly, medicine has significantly benefited from integrating neuroscience [40] and ecology [41].

An area of recent particular interest is the temporal aspect of citations. [42] identified an increasing tendency to cite older papers between 1990 and 2013. [24] reported a recency bias in Natural Language Processing (NLP), showing that post-2015 publications increasingly favour recent work. [43] analyzed citation trends in NLP and another academic field, highlighting the field's strong focus on recent research. So far, it has been unclear whether differences exist in citation patterns in industry research and other AI work. Understanding these patterns is crucial, given Big Tech's potential to shape AI research trajectories and their broader implications for society and science.

Our study addresses gaps from related work by exploring 7 novel research questions, including the volume of papers with acknowledgment of industry funding (Q1), citation behaviours related to funding types (Q2, Q3, Q4), the topics which influence-funded research cites (Q5), and how far back in time industry-funded papers cite (Q6, Q7).

## III. METHODOLOGY

We derive a new dataset from Scopus[2], a database that includes papers published between 1902–2024, totalling $\approx 69.5\,\mathrm{M}$ papers, $\approx 2.2\,\mathrm{B}$ citations, and providing funding information for $\approx 21.0\,\mathrm{M}$ papers. Scopus extracts funding information primarily through the acknowledgment sections of papers, which we use as a lower bound. We provide an overview of the key dataset statistics in Table I. Scopus has a broad coverage of journals and conferences, high quality control and moderation standards, and many peer-reviewed publications with inclusive content coverage [44].

To trace citations from Big Tech to other papers, we construct a citation graph that extends 2 levels deep from the papers published at top AI conferences as part of our **data collection**. We focus our analysis on the period from 2018 to 2022 to capture the marked impact of transformer-based models on AI research. We chose 2022 as the end date because some conference data only extend to 2022 in the Scopus version available. Next, in **data processing**, we determine whether the extracted names belong to Big Tech funders using manual and automatic methods. The final dataset contains $49\,811$ papers.

### A. Data Collection

For the analysis of AI subfields, we select 5 key domains: *Artificial Intelligence* (AI), *Computer Vision* (CV), *Machine Learning* (ML), *Natural Language Processing* (NLP), and *Web & Information Retrieval* (WIr), based on csranking.org[3]. We include 2 top conferences per domain using the same rank and their h5-index. Of the 10 conferences, 7 are in the Scopus database, while the 3 unmatched conferences (i.e., ECCV, NeurIPS, and WWW) get replaced by those with the

[3]https://csrankings.org/#/index?ai&vision&mlmining&nlp&inforet&world

| Time range | 1902–2024 |
|---|---|
| #papers | 69 491 766 |
| #funded papers[†] | 21 047 938 |
| #citations | 2 199 264 185 |
| #papers AI* | 114 090 |
| #funded papers AI[†] | 45 893 |
| #out-citations from AI | 3 308 618 |
| #in-citations to AI | 6 012 570 |

[†]Lower bound. *Sum of articles in Table II.

TABLE I: Overall dataset statistics.

| Field | Conference | Number of articles | h5-index ($\downarrow$) |
|---|---|---|---|
| AI | AAAI | 6793 | 212 |
| | IJCAI | 4203 | 133 |
| CV | CVPR | 9923 | 422 |
| | ICCV* | 4035 | 291 |
| ML | ICLR | 3778 | 303 |
| | ICML* | 5733 | 288 |
| NLP | ACL | 8674 | 192 |
| | EMNLP | 5513 | 176 |
| WIr | SIGIR | 1874 | 90 |
| | WSDM* | 707 | 77 |

TABLE II: The selected top AI conferences ordered by field and increasing h5-index.

third-highest h5-index in their respective fields (i.e., ICCV, ICML, and WSDM). The final list of top AI conferences contains AAAI, IJCAI, CVPR, ICCV, ICML, ICLR, ACL, EMNLP, WSDM, and SIGIR. Table II provides details of the selected top AI conferences.

**Citation Graph.** To define the search scope of Big Tech, we construct a citation graph. The graph originates from papers published at the selected top conferences (root level). It expands through the outgoing citations of these papers (level 1) as well as papers cited by those in level 1 (level 2). We do the same for analogous incoming citations. We extract the names of the funding agencies for each paper in the root papers and those in levels 1 and 2 and identify industry funders (IFs) contributing to AI research. This two-level expansion addresses the limitations of previous approaches focusing only on Big Tech contributions in conferences, which extracted companies by market capitalization. We allow papers outside of conferences because some key AI companies do not publish directly at AI conferences. For example, organizations such as OpenAI, Hugging Face, and Mistral play a fundamental role in advancing the field, but they mostly publish through their content delivery networks.

### B. Data Processing

The process of generating our list of IFs involved five distinct steps:

1) **Extraction of Funding Agencies**: We extracted the names of funding agencies and the number of funded papers from

TABLE III: Overview of funding agency dataset statistics.

| Attribute | Amount |
|---|---|
| Total funding agencies | 78 333 |
| Funding agencies $\geq$ 10 occurrences | 4206 |
| Total corporate funding Agencies* | 3136 |
| Corporate funding agencies (manual analysis) | 382 |
| Corporate funding agencies after standardization | 216 |
| Corporate funding agencies (automatic analysis) | 2754 |

*Sum of manual and automatic analysis.

Scopus.
2) **Manual Analysis**: For agencies with more than ten funded papers, we manually reviewed their names to determine their classification as IF. A funding agency was designated as IF if it was neither public nor non-profit.
3) **Standardization of IF Names**: We standardized IF names to address variations in company nomenclature and alternative descriptions for the same organization (e.g., Amazon, Amazon Research).
4) **Automatic Analysis**: We used fuzzy matching to examine the remaining funding agencies with fewer than ten occurrences. A funding agency was classified as IF if its name included one of the 216 standardized IF names identified in the 3. step.
5) **Integration of Results**: Finally, we combined the results of the manual and automatic analyses into a unified list of IF names. The full list of extracted funding agencies and number of funded papers can be found at https://github.com/Peerzival/impact-big-tech-funding.

In total, we processed 78 333 funding agencies and identified 3136 as IF. Table III summarizes key statistics for the dataset.

### IV. EXPERIMENTS

We use the dataset described above to answer a series of questions about industry funding in AI.

**Q1.** *How many papers have received industry funding in top AI conferences? How does this number vary by research field? Has this number stayed roughly the same, or has it changed markedly over the years?*

**Ans.** We calculate the *percentage of industry-funded papers* ($PIFP$) in a given year $y$ as: $PIFP(y) = \sum_{\forall f_i \in F} \frac{IF(f_i, y)}{P(f_i, y)}$ where $IF(f_i, y)$ represents the number of industry-funded papers in field $f_i$ in year $y$, and $P(f_i, y)$ is the total number of papers in that field. $F$ is the set of all subfields in AI. To isolate industry funding trends within individual subfields, we also compute the *field-specific industry funding percentage* ($FIFP$) for each year, where $FIFP(x, y) = \frac{IF(x, y)}{\sum_{\forall f_i \in F} P(f_i, y)} \cdot 100$

**Results.** Figure 2 tracks the evolution of industry-funded papers in AI research. Figure 2a shows the long-term perspective on industry-funded AI research, while Figure 2b narrows the focus to the selected time frame of 2018–2022. The proportion of industry-funded AI research increases markedly, from $0.6\%$ in 1998 to $9\%$ in 2022 (Figure 2a). A sharp growth occurs between 2016 ($4\%$) and 2020 ($11\%$), reflecting a $180\%$ increase. However, our results show a decline post-2020, with

the share of industry-funded papers dropping from $11\%$ in 2020 to $9\%$ in 2022. Between 2018 and 2022, $10\%$ of all papers published at top AI conferences were funded by Big Tech (Figure 3a). Figure 3b shows the percentage of industry-funded papers in each subfield out of all industry-funded papers. NLP ($31\%$) and CV ($27\%$) have the largest share of industry-funded publications in Scopus. Observe that despite its growth, WIr ($3\%$) lags behind. Overall, $61\%$ of papers are funded, with Big Tech providing funding for $10\%$ of papers on average (Figure 3a). Non-profit or public organizations averagely fund $51\%$ of papers. In ML, $36\%$ of papers are funded by non-profit or public organizations, and $53\%$ of are not funded, while all other fields (i.e., industry-funded and non-industry-funded) have a funding percentage above $50\%$. The highest concentration of Big Tech-funded research occurs in NLP ($11\%$) and ML ($10\%$).

**Discussion.** The results align with [14]'s observations of increasing industry presence in AI between 2016 and 2020. However, our analysis reveals lower percentage values due to our high-precision method of focusing on acknowledging direct funding rather than author institution affiliations. Funding shapes research directions directly [45], indicating industry influence more precisely than affiliations alone.

The decline in industry funding post-2020 does not necessarily imply a reduced industry presence in AI research. Preprint platforms such as arXiv can be confounding factors and show a non-linear increase in AI-related papers [46], [47]. Yet our observations suggest that the industry is shifting away from traditional conference publications in favour of less costly and time-consuming alternatives. For example, the technical report for LLaMA [48] and ChatGPT[4] was released only on arXiv and corporate website, respectively. The industry's ability to conduct exclusive research — owing to their access to critical resources — explains why these papers remain highly relevant despite avoiding peer-reviewed venues. The COVID-19 pandemic further accelerated the use of preprint servers as a rapid publication medium [49], while macroeconomic pressures such as supply chain disruptions, plummeting revenues, and rampant inflation have driven companies to cut costs [50], possibly influencing their reduced investment in academic conference publications. The potential trend of industry migration from publishing at top AI conferences brings several challenges. This trend may lead to less correction of errors through errata and retractions [49]. It also may result in diminished conference sponsorship, as Big Tech may now be independent of conference publications to gain attention.

The high proportion of funded papers shown in Figure 3a — particularly in AI ($73\%$) and NLP ($64\%$) — highlights the financial demands of cutting-edge AI research, exceeding what unfunded individuals or groups can typically pay for. Public and non-profit institutions are the primary supporters of AI research, although the rising share of industry funding in specific fields reflects a growing corporate interest in the field's commercial applications. However, since Scopus bases funding information

on paper acknowledgments, the true industry presence is likely underreported, as disclosures are sometimes omitted or forgotten [51]. Furthermore, funding for AI research extends beyond the flow of money. Funding can include other benefits such as access to models, datasets, computational power, and expertise [2].

**Q2.** *Papers from which funding sources — industry (Big Tech), public, or non-profit —a re most prominently cited in AI research papers? How has this citation trend evolved over time?*

**Ans.** As we know from Q1, $9\%$ of all papers published between 2018 and 2022 were funded by Big Tech. Another key question is whether industry-funded research attracts more citations than non-industry-funded and non-funded papers. Industry ownership of key resources may increase the visibility of these studies, even if their quantity is low. For a fair comparison, we consider the volume of funded papers by funding types (i.e., *industry*-funded, *non-industry*-funded, and *non-funded*) and normalize the citation data according to each funding type's size. We introduce a new metric, called the *Citation Preference Ratio* ($CPR$), which measures whether a paper with a specific funding type is cited more or less frequently than expected, based on its availability. A higher CPR indicates that a paper with a specific funding type is cited more often than the volume of papers would suggest. The CPR from AI to a funding type $f$ is defined as follows:

$$CPR_{AI}(f) = \frac{C(f)}{E(f)} \tag{1}$$

$$\text{where } C(f) = \sum_{\forall f_i \in F} C^{f_i \to f}, \tag{2}$$

$$\text{and } E(f) = \left( \sum_{\forall f_i \in F} \sum_{\forall f_j \in F} C^{f_i \to f_j} \right) \cdot \frac{N_f}{N} \tag{3}$$

where $C(f)$ is the number of citations from all funding types to funding type $f$, $E(f)$ is the expected number of citations to $f$ proportional to its share of total papers, $C^{f_i \to f_j}$ is the number of citations from funding type $f_i$ to funding type $f_j$, $F$ is the set of all subfields, and $N$ is number of papers across all funding types. A $CPR > 1$ shows that the funding type $f$ is more often cited by AI than expected; $CPR < 1$ shows that AI less often cites the funding type $f$ than expected; and $CPR = 1$ shows that the citations are proportional to availability (no citation preference).

**Results.** Figure 4a reveals a consistent upward trend in citation preference towards funded papers since 2019. By 2021, the AI community started citing more industry-funded papers than expected by the number of papers. Despite this growing trend, non-funded research was cited more frequently than industry-funded until 2022. However, non-funded research's CPR has declined since 2019, demonstrating that the AI community is citing fewer non-funded papers than expected by the growth. Notably, the CPR of non-industry-funded papers remains low. Although there are growing numbers of papers annually, top AI papers cite non-funded papers disproportionately less. In
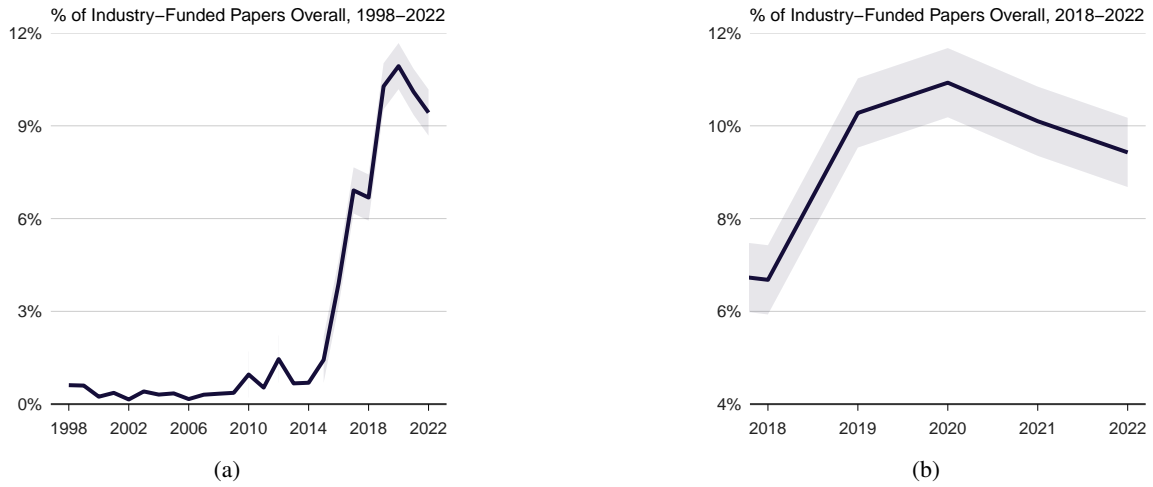
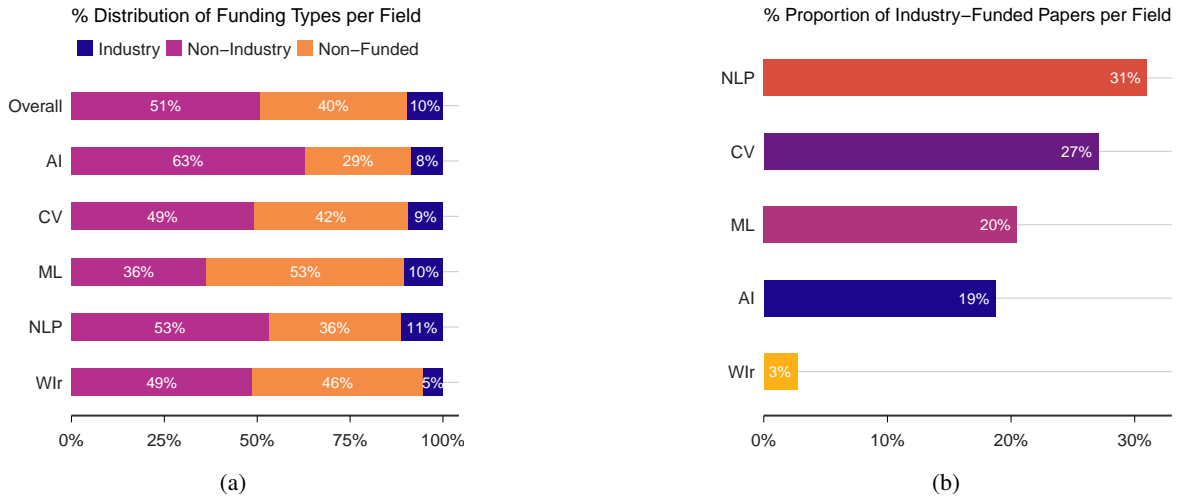Fig. 2: The *PIFP* (a) from 1998 to 2022 and (b) from 2018 to 2022.



Fig. 3: The percentage distribution of (a) all funding types and (b) industry-funded papers for AI subfields.

contrast, they cite industry-funded papers disproportionately more than the growth in papers would suggest.

**Discussion.** The increasing CPR for industry-funded and non-industry-funded papers demonstrates a link between funding and academic visibility in AI research. Voices in the community suggest that industry-funded research of Big Tech labs often serves as a label for excellent methodology and implementation. The marked increase in citation preference of industry-funded research can also be attributed to the industry's contribution of tools and resources, such as frameworks, datasets, and models, which become foundational in AI research and development [1], [2], [14]. The publications of these tools are frequently cited in academic and industry research alike. An example is the paper introducing PyTorch, a widely used machine learning library created by Meta [52].

In response to the growing relevance of AI, government agencies and non-profit organizations have increased their funding [53]. For example, the top 3 funding agencies in our citation graph based on occurrences were the EU Horizon 2020 Program ($3\%$ of occurrences), the Engineering and Physical

Sciences Research Council ($1.8\%$ of occurrences), and the German Research Foundation ($1.8\%$ of occurrences). This shift may have increased the quantity and quality of non-industry-funded publications, leading to rising citation counts. Our results show that research without external funding lags behind in citation metrics.

**Q3.** *To what extent do industry-funded papers cite other industry-funded papers as opposed to other funding types?*

**Ans.** The growing presence of industry funding in AI research raises questions about potential self-reinforcing cycles, where industry-funded research may disproportionately cite other industry-funded work, potentially creating echo chambers [54] that could shape research narratives and discourse even without explicit intention. We calculate the difference in industry-funded outgoing citation percentage to a funding type $f$ versus the average outgoing citations from various funding types to $f$. We rely on the *Outgoing Relative Citational Prominence* ($ORCP$) metric by [55] with one key modification: we adjust the notion of a paper being in specific research fields to a paper being funded by specific funding types (details in Section D). If
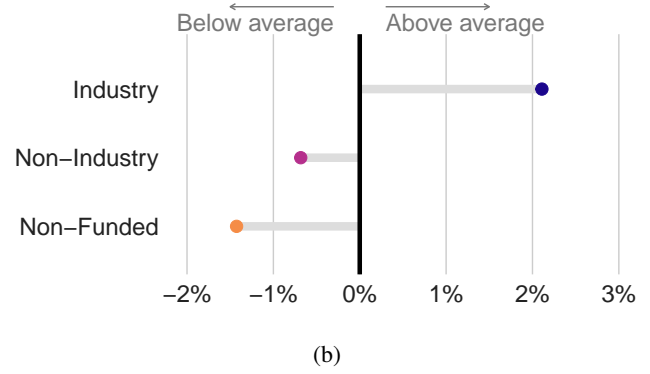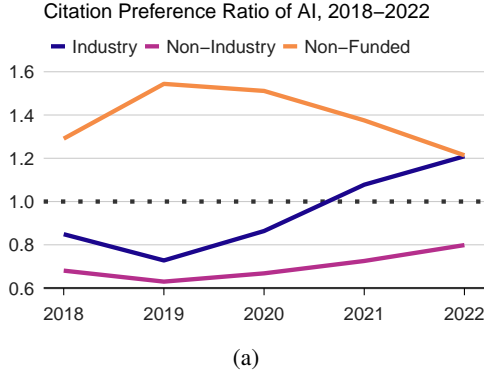
Fig. 4: (a) The *CPR* of AI. (b) Industry-funded research's *ORCP* towards industry-funded, non-industry-funded, and non-funded papers. In both cases, higher values indicate more insularity.

industry-funded research ($IF$) has an ORCP greater than $0$ for $f$, then $IF$ cites $f$ more often than other funding types cite $f$ on average.

**Results.** Figure 4b shows the ORCP scores of industry-funded papers across funding types, with industry-funded research citing other industry-funded research more than average ($ORCP = 2\%$). Notably, despite the presence of extensive non-industry-funded and non-funded research of comparable quantity, both of these funding types have an $ORCP < 0$, implying that industry-funded research cites non-industry-funded and non-funded work markedly less than how much the other funding types cite non-industry-funded and non-funded research. Among all funding types, the highest ORCP to a funding type occurs within the same funding type, indicating that citations to papers of the same funding type are more common than cross-type citations. Non-industry-funded and non-funded research has higher ORCP scores to itself than industry-funded work has to itself. Additionally, industry-funded research has the lowest ORCP among the three funding types when cited by non-industry-funded and non-funded papers.

**Discussion.** The findings show that papers from one funding type prefer to cite papers from the same funding type than others. Despite this pattern, the degree of insularity remains moderate, reflected by an ORCP of $2\%$. Citation insularity emerges from a complex interplay of factors beyond mere research specialization. While industry and public funding sources encompass diverse organizations with varied research priorities, several structural factors may contribute to self-referential citation practices. These can include shared methodological frameworks and established collaborative networks for industry-funded research. Similarly, despite their diverse missions, public funding agencies often create interconnected research communities through targeted programs and review processes, potentially reinforcing specific citation patterns.

**Q4.** *How well are industry-funded papers cited? How does the citation impact vary between different funding types?*
**Ans.** We examine the citation impact of industry-funded papers as a measure of influence on other researchers. We analyze median citations, mean citations, and the h5-index [56] across

TABLE IV: The total number of AI papers published in the last five years, the median and mean number of citations, and the h5-index for different funding types (by decreasing h5-index).

| Funding Type | Count | Median | Mean | h5-index ($\downarrow$) |
|---|---|---|---|---|
| Non-Ind.-Funded | 25 190 | 44 | 170.68 | 893 |
| Ind.-funded | 4801 | 86 | 304.07 | 570 |
| Non-Funded | 19 820 | 10 | 48.86 | 380 |

different funding types. The h5-index is a proxy for impact and influence despite its known limitations in capturing all research dimensions [57], [58].

**Results.** Table IV shows the mean and median number of citations and the h5-index for all papers of a funding type. Observe how funded papers have the highest mean citations (industry-funded: 304.07; non-industry-funded: 170.68), median citations (industry-funded: 86; non-industry-funded: 44), and h5-index (industry-funded: 570; non-industry-funded: 893). Non-industry-funded research has the highest h5-index (893) and the largest volume of papers (25 190), demonstrating a substantial number of highly cited papers and total contributions. Conversely, industry-funded research achieves a disproportionately high h5-index (570) relative to the number of papers, reflecting a focus on high-yield research outputs. By comparison, unfunded research ranks lowest in citation impact among all funding types. Industry-funded research stands out, with $12\%$ of papers having high citation impact, compared to only $4\%$ of non-industry-funded and $2\%$ of non-funded papers. However, industry-funded show consistent citation patterns, suggesting a steady impact compared to other funding types, which tend to have more one-hit successes.

We conducted Levene's test to assess the relationship between funding type and citation counts. The results indicate a weak positive correlation ($\rho = 0.133$), which is highly statistically significant ($p < 0.0001$). While this finding confirms a correlation between funding type and citation counts, it also highlights that funding type accounts for only $13.3\%$ of the observed variance in citation counts. This observation emphasizes the importance of additional factors, such as research quality, topical relevance, and author reputation.

**Discussion.** The number of citations to funded papers shows a

connection between research funding and citations. Big Tech-funded research has a disproportionately high citation impact relative to its publication volume compared to other types, a phenomenon attributable to multiple underlying mechanisms. A possible reason is that industry often recruits students and faculties, which brings established citation networks, increasing visibility and engagement. However, a high citation count is not necessarily an indicator of research excellence; citations are only one factor. A balanced public and private funding approach is essential for fostering a robust and equitable AI research ecosystem. Metrics beyond citations (e.g., transparency) help evaluate research impact, emphasize reproducibility, societal relevance, and ethical considerations. Public institutions can encourage industry participation in open science while providing researchers with the resources necessary to make decisions about engaging in or opting out of industry collaborations.

**Q5.** *Which fields do industry-funded papers cite? How diverse are the outgoing citations in these papers?*

**Ans.** We analyze outgoing citations by funding type to determine whether non-industry-funded and non-funded research concentrates on industry-favored topics. We calculate each funding type's share of citations directed to various fields, defined as the percentage of citations to a given field from a given funding type over all citations from a given funding type to any field. For papers associated with multiple fields, each field receives a citation. We use the fields Scopus pre-classification system[5], which categorizes a paper based on the aims and scope of the title and its content.

**Results.** Figure 5 shows the distribution of citations to the top 10 cited fields for industry-funded papers. Across all funding types, 8 out of 10 most cited fields belong to computer science, highlighting a strong focus on this field and a low outgoing citation field diversity. The top 10 fields account for over 70 % of outgoing citations in each funding type. Industry-funded research shows the highest concentration, with 77 % of citations directed to these top 10 fields, compared to 75 % and 72 % for non-industry-funded and non-funded papers, respectively. The 4 main fields cited remain consistent across funding types, constituting more than 50 % of citations within the top 10 fields, indicating a common primary interest across funding types. However, industry-funded papers show an increased interest in linguistics, while non-industry-funded and non-funded papers emphasize AI more prominently. Additionally, non-industry-funded research shows a stronger orientation toward theoretical work, contrasting with the industry and non-funded paper's emphasis on networks and signal processing.

**Discussion.** The convergence of research fields across funding types, alongside the growing engagement with industry-funded work, demonstrates Big Tech influence. Building on the findings of [12], we show that industry-funded research exhibits lower thematic diversity compared to non-industry-funded and non-funded research, demonstrated by the high citation density in
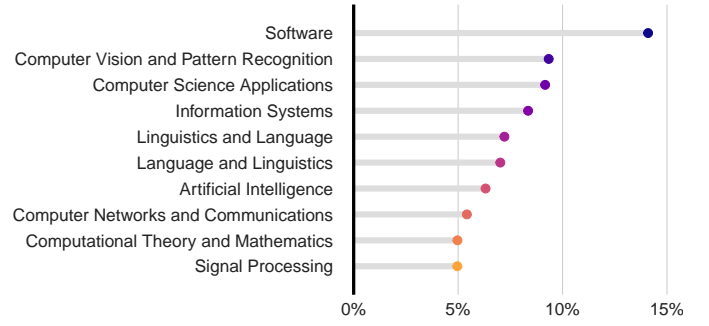
Fig. 5: Percentage of outgoing citations from industry-funded papers to the top 10 most cited fields.

the 10 ten most-cited fields. Our results reflect a concentration of industry-funded work in data-hungry and computationally intensive fields, such as computer vision and information retrieval (information systems). In contrast, research not funded by Big Tech shows a relatively strong focus on AI and theoretical fields such as mathematics. This emphasis does not preclude engagement with data-driven or computationally intensive fields. Instead, it indicates a broader distribution of research interests compared to the narrower focus of industry-funded work.

**Q6.** *What is the average age of cited papers for different funding types?*

**Ans.** Scholars frequently engage with related work across disciplines to validate or challenge earlier findings, situate their contributions, and extend the boundaries of knowledge. However, the tendency to not cite enough relevant good work from the past (more than a few years old) — referred to as "citation amnesia" [24] — remains a pressing issue. Forgetting some older works can be helpful, as it makes space for new ideas. However, too much forgetting can lead to unnecessary reinvention of methods. Building on insights from [24], [43], we investigate temporal citation patterns to quantify this phenomenon. We analyze the citations for each funding type to other papers and calculate how far back in time the *cited* papers are published. When a paper $x$ cites a paper $y_i$, then the age of the citation ($AoC$) is the difference between the year of publication of $x$ and $y_i$, as in $AoC(x, y_i) = YoP(x) - YoP(y_i)$. We then calculate the $AoC$ for each citation in a paper and average them (*mAoC*).

**Results.** Figure 6 shows the distribution of *mAoCs* for all papers of each funding type and overall across the years after the publication of the *cited* paper. The y-axis represents the average percentage of citations papers received for the years (x-axis). Most citations occur for papers published 2 years prior ($AoC = 2$). The citation patterns for all funding types show a similar increasing trend, followed by an abrupt decline in the following years. Non-funded papers have a lower peak value but maintain higher citation rates than other funding types. Additionally, industry-funded papers have the highest percentage of citations in the publication year, while non-funded papers have the lowest. The *mAoC* for papers published between 2018 and 2022 shows industry-funded papers have the

Fig. 6: Distribution of citation age as measured by *AoC* for papers in AI (overall and by funding type).



Fig. 7: Citation ages as measured by *mAoC* for industry-funded papers between 2018 and 2022. The standard deviation for each year is displayed below the respective violin plot. The median (white diamond) and the mean *mAoC* (dark line) are shown for each year.

lowest mean *mAoC* of $4.79$, followed closely by non-industry-funded papers with $4.92$, and non-funded papers at $5.03$ (details in Section E in Section E).

**Discussion.** Our results show that papers typically receive the highest number of citations 2 years after publication, followed by a strong decline, which is consistent with related work [24], [43]. At their peak, non-funded papers receive fewer citations than funded research, but their citation decline is more gradual. This dynamic, coupled with the higher *mAoC* value (Section E), suggests that non-funded research accumulates citations over a longer time and cites older work compared to funded research. Non-funded work may focus on citing foundational or theoretical contributions that continue attracting citations over time. In contrast, funded research may prioritize emerging topics and cutting-edge innovations. This encourages new developments to build upon recent advances, often neglecting foundational work.

**Q7.** *What is the distribution of mAoC in industry-funded papers? How does this distribution vary over the years?*

**Ans.** We compute the *mAoC* for each industry-funded paper and the median and mean *mAoC* for all such papers over time. Section E provides more details on the *mAoC*.

**Results.** Figure 7 shows the violin plots for distributions of *mAoC* in industry-funded papers across years. Each plot highlights the median (white diamond within the grey rectangle), representing the recency of citations for that year. Over time, the median *mAoC* for industry-funded papers shows a consistent decline, indicating an increasing focus on recent work. The shrinking size of the second and third quartiles (halves of the grey rectangle) indicates that citations are increasingly concentrated around the median, reflecting a narrowed citation age. From 2018 to 2022, the mean *mAoC* closely follows the median trend, but it remains consistently higher, revealing a right skew in the data. This skew is due to several papers citing much older papers. The decreasing standard deviation also suggests diminishing citation age diversity, possibly reinforcing the trend toward citing newer literature. By 2022, the violin's
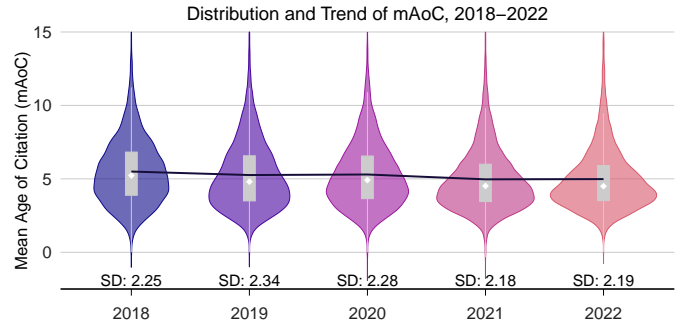
density transforms into a spinning tractroid top[6], reflecting a concentration of *mAoC* values near five years.

**Discussion.** The results show a decline in citations to older works within industry-funded papers and a reduction in the temporal diversity of citations. Although the exact causes of this trend remain uncertain, multiple factors contribute to the evolving citation dynamics. The substantial impact of transformers on NLP and ML can increasingly favour more recent publications. The pressures of the "publish or perish" paradigm further exacerbate this trend, encouraging researchers to divide their work into smaller, publishable units [60]. Our results add to and are consistent with the mean-citation age results in [43].Our analysis focuses on AI publications between 2018 and 2022, situating these results in the overall trajectory of how temporal citation patterns in AI have evolved recently.

## V. CONCLUDING REMARKS

This work examined the impact of Big Tech funding on AI research. To enable this analysis, we compiled a dataset that includes $\approx 49.8\,\mathrm{K}$ AI papers, their funding agencies, citations from other papers to AI papers, and citations from AI papers to other papers. We introduced the *Citation Preference Ratio*, a novel metric that shows a growing trend within the AI research community to reference Big Tech-funded work. We also employed established metrics such as the *Relative Citational Prominence*, which highlights an increasing insularity in industry-funded research, and the *Mean Age of Citation* metric, which shows the tendency of industry-funded research to cite recent literature. While the presence of Big Tech-funded research in top AI conferences is declining, its citational influence continues to grow. Industry-funded papers cite fewer contributions from non-industry-funded and non-funded research despite a more extensive growth in paper volume than industry-funded papers (growing insularity).

The manual and automated methods for identifying funding agencies in our work introduce risks of interpretation and matching errors. Relying solely on citation-based metrics

---

[6]Form of the iconic spinning top in the movie Inception.

captures influence but cannot fully explain the qualitative drivers behind these trends. A broader analysis of industry impact, encompassing resource allocation and the interplay between public and private funding, are still necessary. From an ethical perspective, this work underscores the importance of interpreting these citation-based findings responsibly, remembering that metrics alone cannot dictate value, nor should they justify prioritizing one funding source over another.

As researchers, we are not only observers of these trends and fast paced developments. Instead, we have agency in this process. One could further say we, as a field, have a responsibility to reflect on these trends, to discuss and vote for appropriate actions to shape and direct the field's future. Public institutions play a critical role in this effort by improving their funding policies and providing researchers with the tools and support they need to receive diverse project funding. Bridges between industry and public funding — to learn from one another and benefit from infrastructure — are vital for a healthy research environment. Government efforts to provide competitive infrastructure, such as computing, data, and compensation, are also key to attracting talent for open research.

REFERENCES

[1] N. Ahmed and M. Wahed, "The de-democratization of ai: Deep learning and the compute divide in artificial intelligence research," 2020.

[2] P. Verdegem, "Dismantling ai capitalism: the commons as an alternative to the power concentration of big tech," *AI & society*, vol. 39, no. 2, pp. 727–737, 2024.

[3] I. M. Cockburn, R. Henderson, and S. Stern, "The impact of artificial intelligence on innovation," National Bureau of Economic Research, Working Paper 24449, March 2018. [Online]. Available: http://www.nber.org/papers/w24449

[4] C. B. Frey and M. A. Osborne, "The future of employment: How susceptible are jobs to computerisation?" *Technological forecasting and social change*, vol. 114, pp. 254–280, 2017.

[5] F. Jiang, Y. Jiang, H. Zhi, Y. Dong, H. Li, S. Ma, Y. Wang, Q. Dong, H. Shen, and Y. Wang, "Artificial intelligence in healthcare: past, present and future," *Stroke and vascular neurology*, vol. 2, no. 4, 2017.

[6] R. Dobbe and M. Whittaker, "Ai and climate change: how they're connected, and what we can do about it," *AI Now Institute*, vol. 17, 2019.

[7] M. Brundage, S. Avin, J. Clark, H. Toner, P. Eckersley, B. Garfinkel, A. Dafoe, P. Scharre, T. Zeitzoff, B. Filar *et al.*, "The malicious use of artificial intelligence: Forecasting, prevention, and mitigation," *arXiv preprint arXiv:1802.07228*, 2018.

[8] S. Plan, "The national artificial intelligence research and development strategic plan," *National Science and Technology Council, Networking and Information Technology Research and Development Subcommittee*, 2016.

[9] A. Koenecke, A. Nam, E. Lake, J. Nudell, M. Quartey, Z. Mengesha, C. Toups, J. R. Rickford, D. Jurafsky, and S. Goel, "Racial disparities in automated speech recognition," *Proceedings of the national academy of sciences*, vol. 117, no. 14, pp. 7684–7689, 2020.

[10] C. Kuhlman, L. Jackson, and R. Chunara, "No computation without representation: Avoiding data and algorithm biases through diversity," 2020.

[11] S. M. West, M. Whittaker, and K. Crawford, "Discriminating systems," *AI Now*, pp. 1–33, 2019.

[12] J. Klinger, J. Mateos-Garcia, and K. Stathoulopoulos, "A narrowing of ai research?" 2022. [Online]. Available: https://arxiv.org/abs/2009.10385

[13] M. Whittaker, "The steep cost of capture," *Interactions*, vol. 28, no. 6, pp. 50–55, 2021.

[14] N. Ahmed, M. Wahed, and N. C. Thompson, "The growing influence of industry in ai research," *Science*, vol. 379, no. 6635, pp. 884–886, 2023.

[15] G. A. Montes and B. Goertzel, "Distributed, decentralized, and democratized artificial intelligence," *Technological Forecasting and Social Change*, vol. 141, pp. 354–358, 2019.

[16] M. Abdalla, J. P. Wahle, T. Ruas, A. Névéol, F. Ducel, S. Mohammad, and K. Fort, "The elephant in the room: Analyzing the presence of big tech in natural language processing research," in *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, A. Rogers, J. Boyd-Graber, and N. Okazaki, Eds. Toronto, Canada: Association for Computational Linguistics, Jul. 2023, pp. 13 141–13 160. [Online]. Available: https://aclanthology.org/2023.acl-long.734

[17] M. Riedl, "Ai democratization in the era of gpt-3," *The Gradient*, vol. 25, 2020.

[18] M. Abdalla and M. Abdalla, "The grey hoodie project: Big tobacco, big tech, and the threat on academic integrity," in *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, ser. AIES '21. New York, NY, USA: Association for Computing Machinery, 2021, p. 287–297. [Online]. Available: https://doi.org/10.1145/3461702.3462563

[19] P. Trouiller, P. Olliaro, E. Torreele, J. Orbinski, R. Laing, and N. Ford, "Drug development for neglected diseases: a deficient market and a public-health policy failure," *The Lancet*, vol. 359, no. 9324, pp. 2188–2194, 2002.

[20] R. Jurowetzki, D. Hain, J. Mateos-Garcia, and K. Stathoulopoulos, "The privatization of ai research(-ers): Causes and potential consequences – from university-industry interaction to public research brain-drain?" 2021. [Online]. Available: https://arxiv.org/abs/2102.01648

[21] A. Birhane, P. Kalluri, D. Card, W. Agnew, R. Dotan, and M. Bao, "The values encoded in machine learning research," 2022. [Online]. Available: https://arxiv.org/abs/2106.15590

[22] S. Giziński, P. Kaczyńska, H. Ruczyński, E. Wiśnios, B. Pieliński, P. Biecek, and J. Sienkiewicz, "Big tech influence over ai research revisited: Memetic analysis of attribution of ideas to affiliation," *Journal of Informetrics*, vol. 18, no. 4, p. 101572, 2024.

[23] M. Färber and L. Tampakis, "Analyzing the impact of companies on ai research based on publications," *Scientometrics*, vol. 129, no. 1, pp. 31–63, 2024.

[24] J. Singh, M. Rungta, D. Yang, and S. M. Mohammad, "Forgotten knowledge: Examining the citational amnesia in nlp," *arXiv preprint arXiv:2305.18554*, 2023.

[25] M. Rungta, J. Singh, S. M. Mohammad, and D. Yang, "Geographic citation gaps in nlp research," *arXiv preprint arXiv:2210.14424*, 2022.

[26] S. M. Mohammad, "Examining citations of natural language processing literature," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, D. Jurafsky, J. Chai, N. Schluter, and J. Tetreault, Eds. Online: Association for Computational Linguistics, Jul. 2020, pp. 5199–5209. [Online]. Available: https://aclanthology.org/2020.acl-main.464

[27] S. Della Sala and J. Brooks, "Multi-authors' self-citation: A further impact factor bias?" *Cortex; a journal devoted to the study of the nervous system and behavior*, vol. 44, no. 9, pp. 1139–1145, 2008.

[28] M. Callaham, R. L. Wears, and E. Weber, "Journal prestige, publication bias, and other characteristics associated with citation of published studies in peer-reviewed journals," *Jama*, vol. 287, no. 21, pp. 2847–2850, 2002.

[29] J. P. Wahle, T. Ruas, S. M. Mohammad, and B. Gipp, "D3: A massive dataset of scholarly metadata for analyzing the state of computer science research," *arXiv preprint arXiv:2204.13384*, 2022.

[30] G. Buela-Casal and I. Zych, "Analysis of the relationship between the number of citations and the quality evaluated by experts in psychology journals," *Psicothema*, pp. 270–276, 2010.

[31] R. P. C. Lira, R. M. C. Vieira, F. A. Gonçalves, M. C. A. Ferreira, D. Maziero, T. H. M. Passos, and C. E. L. Arieta, "Influence of english language in the number of citations of articles published in brazilian journals of ophthalmology," *Arquivos Brasileiros de Oftalmologia*, vol. 76, pp. 26–28, 2013.

[32] S. M. Mohammad, "Gender gap in natural language processing research: Disparities in authorship and citations," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, D. Jurafsky, J. Chai, N. Schluter, and J. Tetreault, Eds. Online: Association for Computational Linguistics, Jul. 2020, pp. 7860–7870. [Online]. Available: https://aclanthology.org/2020.acl-main.702/

[33] P. Chatterjee and R. M. Werner, "Gender disparity in citations in high-impact journal articles," *JAMA Network Open*, vol. 4, no. 7, pp. e2114509–e2114509, 2021.

[34] S. Abdalla, M. Abdalla, M. Saad, D. Jones, S. Podolsky, and M. Abdalla, "Ethnicity and gender trends of uk authors in the british medical journal and the lancet over the past two decades: a comprehensive longitudinal analysis," *EClinicalMedicine*, vol. 64, 2023.

[35] R. Costas, M. Bordons, T. N. Van Leeuwen, and A. F. Van Raan, "Scaling rules in the science system: Influence of field-specific citation characteristics on the impact of individual researchers," *Journal of the American Society for Information Science and Technology*, vol. 60, no. 4, pp. 740–753, 2009.

[36] A. L. Porter, D. J. Roessner, and A. E. Heberger, "How interdisciplinary is a given body of research?" *Research evaluation*, vol. 17, no. 4, pp. 273–282, 2008.

[37] L. Leydesdorff, C. S. Wagner, and L. Bornmann, "Interdisciplinarity as diversity in citation patterns among journals: Rao-stirling diversity, relative variety, and the gini coefficient," *Journal of informetrics*, vol. 13, no. 1, pp. 255–269, 2019.

[38] A. Einstein, "Über einen die erzeugung und verwandlung des lichtes betreffenden heuristischen gesichtspunkt," 1905.

[39] N. Bohr, "I. on the constitution of atoms and molecules," *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science*, vol. 26, no. 151, pp. 1–25, 1913.

[40] S. Postel, K. Bawa, L. Kaufman, C. H. Peterson, S. Carpenter, D. Tillman, P. Dayton, S. Alexander, K. Lagerquist, L. Goulder *et al.*, *Nature's services: Societal dependence on natural ecosystems*. Island Press, 2012.

[41] D. W. Pearce and R. K. Turner, *Economics of natural resources and the environment*. Johns Hopkins University Press, 1989.

[42] A. Verstak, A. Acharya, H. Suzuki, S. Henderson, M. Iakhiaev, C. C. Y. Lin, and N. Shetty, "On the shoulders of giants: The growing impact of older articles," 2014. [Online]. Available: https://arxiv.org/abs/1411.0275

[43] J. P. Wahle, T. Ruas, M. Abdalla, B. Gipp, and S. M. Mohammad, "Citation amnesia: Nlp and other academic fields are in a citation age recession," *arXiv preprint arXiv:2402.12046*, 2024.

[44] R. Pranckutė, "Web of science (wos) and scopus: The titans of bibliographic information in today's academic world," *Publications*, vol. 9, no. 1, 2021. [Online]. Available: https://www.mdpi.com/2304-6775/9/1/12

[45] M. Thelwall, S. Simrick, I. Viney, and P. Van den Besselaar, "What is research funding, how does it influence research, and how is it recorded? key dimensions of variation," *Scientometrics*, vol. 128, no. 11, pp. 6085–6106, 2023.

[46] M. Krenn, L. Buffoni, B. Coutinho, S. Eppel, J. G. Foster, A. Gritsevskiy, H. Lee, Y. Lu, J. P. Moutinho, N. Sanjabi, R. Sonthalia, N. M. Tran, F. Valente, Y. Xie, R. Yu, and M. Kopp, "Forecasting the future of artificial intelligence with machine learning-based link prediction in an exponentially growing knowledge network," *Nature Machine Intelligence*, vol. 5, no. 11, p. 1326–1335, Oct. 2023. [Online]. Available: http://dx.doi.org/10.1038/s42256-023-00735-0

[47] N. Maslej, L. Fattorini, E. Brynjolfsson, J. Etchemendy, K. Ligett, T. Lyons, J. Manyika, H. Ngo, J. C. Niebles, V. Parli, Y. Shoham, R. Wald, J. Clark, and R. Perrault, "Artificial intelligence index report 2023," 2023. [Online]. Available: https://arxiv.org/abs/2310.03715

[48] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, A. Rodriguez, A. Joulin, E. Grave, and G. Lample, "Llama: Open and efficient foundation language models," 2023. [Online]. Available: https://arxiv.org/abs/2302.13971

[49] P. Smart, "The evolution, benefits, and challenges of preprints and their interaction with journals," *Science Editing*, vol. 9, no. 1, pp. 79–84, 2022.

[50] J. Morgan, "The future of Big Tech | J.P. Morgan Research," Dec. 2023. [Online]. Available: https://www.jpmorgan.com/insights/global-research/technology/future-of-big-tech

[51] J. Wang and P. Shapira, "Is there a relationship between research sponsorship and publication impact? an analysis of funding acknowledgments in nanotechnology papers," *PloS one*, vol. 10, no. 2, p. e0117727, 2015.

[52] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Köpf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, "Pytorch: An imperative style, high-performance deep learning library," 2019. [Online]. Available: https://arxiv.org/abs/1912.01703

[53] M. T. André, "Ai investment: Eu and global indicators," 2024.

[54] M. Cinelli, G. De Francisci Morales, A. Galeazzi, W. Quattrociocchi, and M. Starnini, "The echo chamber effect on social media," *Proceedings of the National Academy of Sciences*, vol. 118, no. 9, p. e2023301118, 2021.

[55] J. P. Wahle, T. Ruas, M. Abdalla, B. Gipp, and S. M. Mohammad, "We are who we cite: Bridges of influence between natural language processing and other academic fields," *arXiv preprint arXiv:2310.14870*, 2023.

[56] J. E. Hirsch, "An index to quantify an individual's scientific research output," *Proceedings of the National academy of Sciences*, vol. 102, no. 46, pp. 16 569–16 572, 2005.

[57] L. Bornmann and H.-D. Daniel, "What do we know about the h index?" *Journal of the American Society for Information Science and technology*, vol. 58, no. 9, pp. 1381–1385, 2007.

[58] R. Costas and M. Bordons, "The h-index: Advantages, limitations and its relation with other bibliometric indicators at the micro level," *Journal of informetrics*, vol. 1, no. 3, pp. 193–203, 2007.

[59] Scopus, "What are Scopus subject area categories and ASJC codes?" Aug. 2024. [Online]. Available: https://service.elsevier.com/app/answers/detail/a_id/12007/supporthub/scopus/

[60] "Publish or perish," *Nature*, vol. 467, no. 7313, pp. 252–252, Sep. 2010. [Online]. Available: https://www.nature.com/articles/467252a

[61] W. Liu, "Accuracy of funding information in scopus: A comparative case study," *Scientometrics*, vol. 124, no. 1, pp. 803–811, 2020.

[62] J. Freyne, L. Coyle, B. Smyth, and P. Cunningham, "Relative status of journal and conference publications in computer science," *Commun. ACM*, vol. 53, no. 11, p. 124–132, nov 2010. [Online]. Available: https://doi.org/10.1145/1839676.1839701

[63] M. Valenzuela, V. Ha, and O. Etzioni, "Identifying meaningful citations," in *Workshops at the twenty-ninth AAAI conference on artificial intelligence*, 2015.

[64] X. Zhu, P. Turney, D. Lemire, and A. Vellino, "Measuring academic influence: Not all citations are equal," *Journal of the Association for Information Science and Technology*, vol. 66, no. 2, pp. 408–427, 2015.

[65] J. P. Ioannidis, J. Baas, R. Klavans, and K. W. Boyack, "A standardized citation metrics author database annotated for scientific field," *PLoS biology*, vol. 17, no. 8, p. e3000384, 2019.

[66] M. W. Nielsen and J. P. Andersen, "Global citation inequality is on the rise," *Proceedings of the National Academy of Sciences*, vol. 118, no. 7, p. e2012208118, 2021.

[67] J. P. Wahle, T. Ruas, S. M. Mohammad, N. Meuschke, and B. Gipp, "Ai usage cards: Responsibly reporting ai-generated content," 2023. [Online]. Available: https://arxiv.org/abs/2303.03886

[68] L. Kaesberg, T. Ruas, J. P. Wahle, and B. Gipp, "CiteAssist: A system for automated preprint citation and BibTeX generation," in *Proceedings of the Fourth Workshop on Scholarly Document Processing (SDP 2024)*, T. Ghosal, A. Singh, A. Waard, P. Mayr, A. Naik, O. Weller, Y. Lee, S. Shen, and Y. Qin, Eds. Bangkok, Thailand: Association for Computational Linguistics, 08 2024, pp. 105–119. [Online]. Available: https://aclanthology.org/2024.sdp-1.10/

APPENDIX

## A. Details of the Limitations

*1) Manual Analysis:* The manual analysis has a few limitations. First, because examining thousands of funding agencies individually is time-intensive, this analysis only includes $5\%$ of the extracted funding agencies. Notably, this $5\%$ covers $74\%$ of all funding occurrences, providing robust overall representation. Second, we identified IFs based on available online information. In cases where insufficient data prevented us from confirming a IF, we marked agencies as non-funded, potentially leading to false negatives in the dataset's metadata. Third, since this analysis was performed solely by one person, interpretation variability may affect the consistency and quality of the data.

*2) Automatic Analysis:* Automated fuzzy text matching with standardized company names can incorrectly match unrelated agencies, leading to false positives when identifying IFs. To mitigate this risk of false positives, a high similarity threshold $(90\%)$ was set, although this conservative threshold could reduce the number of matches and potentially miss some IFs. However, most funding agency names included the standardized identifier (e.g., Google, Google DeepMind, Google Cloud), minimizing the risk of missing relevant IFs.

*3) General Limitations:* Beyond technical limitations, this study faces broader constraints. Industry involvement in AI research today extends beyond direct financial contributions to include access to models, datasets, computational resources, and specialized expertise [2], [14], [15], [17]. This analysis relies exclusively on Scopus funding data, which originates from paper acknowledgments [61] to capture Big Tech influence. The extent of funding information in these acknowledgments varies depending on the details provided in the publications, reflecting disciplinary and regional disparities in funding reporting practices [44]. Notably, studies by [61] and [44] identify errors in Scopus funding acknowledgment text and funding agency fields. It is crucial to interpret the consistency and quality of the identified industry presence with a grain of salt.

This research focuses on publications from prominent AI-related conferences rather than all AI-related academic publications. Although leading conferences shape the academic research agenda [62], this selection excludes vibrant, often non-English AI communities and venues, limiting the generalizability of our findings to the global AI research landscape. Future studies must explore industry influence across diverse sub-communities and venues worldwide.

Furthermore, this analysis quantifies influence primarily through citations, a method with inherent limitations. Citation counts lack nuance, as not all citations reflect the same level of influence [63], [64]. Additionally, citation patterns are affected by biases [26], [65], [66]. This work also examines citation practices on a large scale, focusing on quantitative trends. Qualitative aspects may reveal why Big Tech-funded research receives more engagement within the AI community, shows growing insularity, and cites recent over older literature. Several factors may contribute to this, such as the volume of recent

publications, the applicability of industry-funded research, and the technical relevance of industry-funded work.

Our analysis also did not address the allocation of financial resources by industry across AI subfields and conferences. Analyzing the cash flows from industry to AI research and their impact over time could reveal whether financial resources markedly drive influential AI research or whether other resources are key. We leave the exploration of cash flows, their impact, and their presence over time for future work.

## B. Details of Ethical Considerations

This study analyzes the scientific literature at an aggregate level and not on individual papers or authors, using data from the Scopus database. The database provides metadata such as titles, authors, funding agencies, and publication years, which are used without infringing copyrighted content.

A critical aspect of this study is its reliance on citation counts as a proxy to characterize funding types. Although citations serve as a convenient metric, this approach raises concerns about potential misinterpretation or misuse of our findings. For example, the observed high number of outgoing citations to Big Tech-funded research should not be used as a rationale for diminishing research funded by non-industry sources or conducted without external funding. A more comprehensive evaluation framework may be beneficial in addressing the risks of oversimplified interpretations. Such a framework would integrate multiple dimensions, including relevance, popularity, resource availability, impact, geographic context, and temporal trends, thus mitigating the problems of shallow analysis.

## C. Details on the Extraction of Companies

We searched for company names and common aliases (e.g., Microsoft, Microsoft Azure, Microsoft Cloud Computing Research Centre) using the fuzzywuzzy python package[7] with a $90\%$ threshold.

TABLE V: Company name standardization.

| Names of funding agencies | Std. Name |
|---|---|
| Microsoft, Microsoft Azure, Microsoft Research | Microsoft |
| Amazon, AWS, Amazon Research | Amazon |
| Google, Google DeepMind, Google Cloud | Google |
| Nvidia, NVIDIA AI Center, NVIDIA Corp | Nvidia |

TABLE VI: The table shows examples of standardized funding agency names.

## D. Details on the Outgoing Relative Citational Prominence (ORCP)

We rely on the *Outgoing Relative Citational Prominence* ($ORCP$) metric by [55] with one key modification: we adjust the notion of a paper being in specific research fields to a paper being funded by specific funding types. If industry-funded research ($IF$) has an ORCP greater than 0 for $f$, then $IF$ cites

---

[7]https://pypi.org/project/fuzzywuzzy/

$f$ more often than other funding types cite $f$ on average. The following equation shows that metric.

$$ORCP_{IF}(f) = X(f) - Y(f) \qquad (4)$$

$$\text{where } X(f) = \frac{C^{IF \rightarrow f}}{\sum_{\forall f_j \in F} C^{IF \rightarrow f_j}}, \qquad (5)$$

$$\text{and } Y(f) = \frac{1}{N} \sum_{i=1}^{N} \frac{C^{f_i \rightarrow f}}{\sum_{\forall f_j \in F} C^{f_i \rightarrow f_j}} \qquad (6)$$

where $F$ is the set of all funding types, $N$ is the number of all funding types, i.e. 3, and $C^{f_i \rightarrow f_j}$ represents the number of citations from papers in funding type $f_i$ to papers in funding type $f_j$.

### E. Details on the mean Age of Citations

The age of the citation ($AoC$) is the difference between the year of publication ($YoP$) of $x$ and $y_i$:

$$AoC(x, y_i) = YoP(x) - YoP(y_i) \qquad (7)$$

We then calculate the AoC for each of the citations of a paper and average them:

$$mAoC(x, y_i) = \frac{1}{N} \sum_{i}^{N} AoC(x, y_i) \qquad (8)$$

where $N$ denotes the total number of references in paper $x$. For example, if a paper $x$ from 2022 cites two papers, one from 2010 and one from 2020, the $mAoC$ of the paper $x$ is 7 years.

Section E shows the mean $mAoC$ for papers published between 2018 and 2022, grouped by funding type. Observe how industry-funded papers have the lowest mean $mAoC$ of 4.79, followed closely by non-industry-funded papers with a mean $mAoC$ of 4.92, and non-funded papers at 5.03.

| Funding Type | $mAoC \pm 95\%$ Conf. ($\uparrow$) |
|---|---|
| Industry | $4.79 \pm 0.02$ |
| Non-Industry | $4.92 \pm 0.01$ |
| Non-Funded | $5.08 \pm 0.02$ |

TABLE VII: The $mAoC$ and confidence intervals for different funding types are ordered by increasing $mAoC$.

### F. Supplementary Citation Graph Details

Figure 8 shows an sample of the described citation graph in Section III-A.

### G. Full Conference Names

Section G shows the full names of the matched and unmatched top AI conferences.

### H. Supplemental Experimental Results

In addition to the primary results presented in this study, we describe supplementary results in the form of additional statistics and plots.
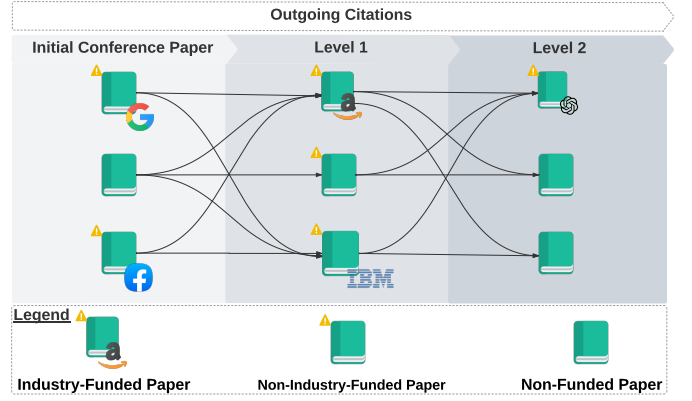


Fig. 8: Example citation graph to identify funding agencies for the search scope.

| Full Name | Acronym | Field |
|---|---|---|
| Advancement of Artificial Intelligence | AAAI | AI |
| International Joint Conference on Artificial Intelligence | IJCAI | |
| Conference on Computer Vision and Pattern Recognition | CVPR | CV |
| International Conference on Computer Vision* | ICCV* | |
| European Conference on Computer Vision† | ECCV† | |
| International Conference on Machine Learning* | ICML* | ML |
| International Conference on Learning Representations | ICLR | |
| Conference and Workshop on Neural Information Processing Systems† | NeurIPS† | |
| Association for Computational Linguistics | ACL | NLP |
| Empirical Methods in Natural Language Processing | EMNLP | |
| International Conference on Web Search and Data Mining* | WSDM* | WIr |
| Conference on Research and Development in Information Retrieval | SIGIR | |
| International World Wide Web Conferences† | WWW† | |

*Replacing conference. †Replaced conference.

TABLE VIII: Full names, acronym, and field of matched and unmatched AI conferences.

*1) Extended Results on Citation Preference Ratio:*
Figure 9 shows the Citation Preference Ratio (CPR) of industry-funded papers, non-industry-funded papers, and non-funded papers to different funding types over time.

*2) Extended Results on Outgoing Relative Citation Prominence:*
Figure 10 shows the ORCP for (a) non-industry-funded papers, and (b) non-funded papers.

*3) Extended Results on Cited Fields:*
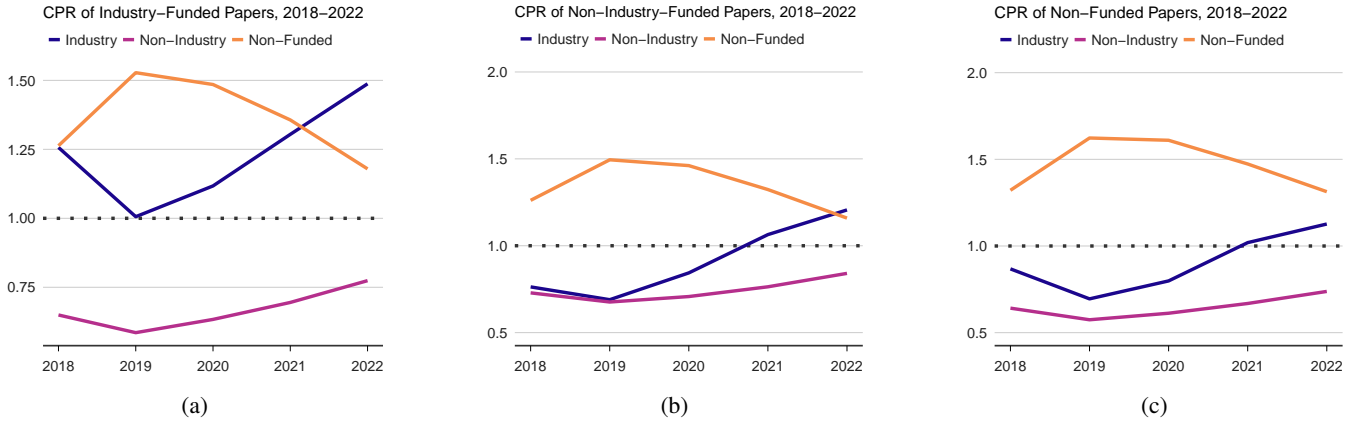Figure 11 shows the top 10 cited fields for (a) non-industry-

Fig. 9: The Citation Preference Ratio (CPR) of (a) industry-funded papers, (b) non-industry-funded papers, and (c) non-funded papers towards all funding types over time.
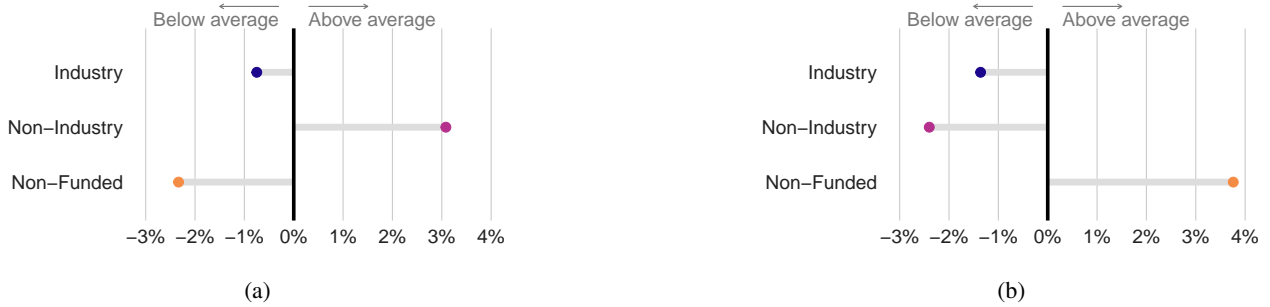


Fig. 10: Outgoing Relative Citational Prominence (ORCP) scores for (a) non-industry-funded papers and (b) non-funded papers.

funded papers and (b) non-funded papers.

## I. Additional Research Questions

**ARQ1.** *Which type of funding is most influenced by industry-funded research? How has this changed over the years?*

**Ans.** To determine the funding types most affected by industry-funded research, we analyze the citation sources to industry-funded research by funding type. Thus, we calculate the average percentage of industry-funded references per paper, i.e., the mean ratio of citations from papers with a given funding type to industry-funded papers relative to the total citations of papers with that funding type. This approach measures how much different funding types rely on or interact with industry-funded research over time.
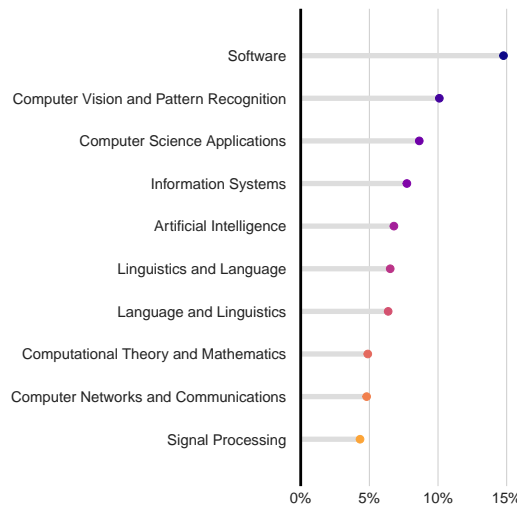
**Results.** Figure 12 shows the proportion of outgoing citations to industry-funded papers per paper and funding type over time. It also shows the macro average of this proportion across all funding types. The share of citations referencing Big Tech-funded research has increased markedly across all funding types since 2018. Non-industry-funded papers showed strong growth, with $44\%$ increase in citations to industry-funded work per paper between 2018 and 2022 after a lower percentage start. This growth surpasses the growth of industry-funded papers ($41\%$) and non-funded papers ($39\%$). Despite this rise in cross-funding-type engagement, industry-funded papers maintained

the highest proportion of outgoing citations per paper to other industry-funded research ($15\%$ in 2022), underlining the self-referential trend of industry-funded research.
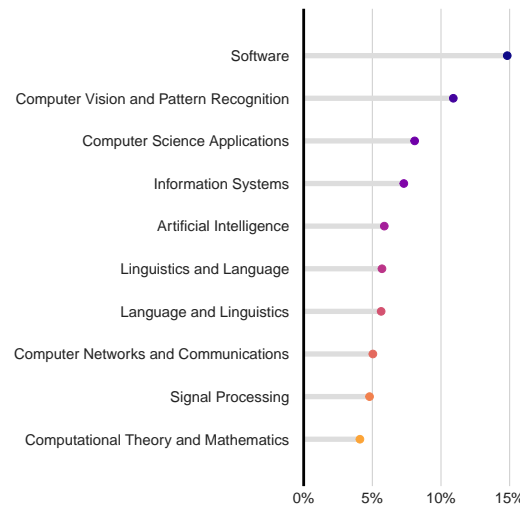
**Discussion.** The rising proportion of outgoing citations to industry-funded papers highlights a marked increase in engagement with industry-funded research across funding types. Non-industry-funded researchers show growing interest in industry-driven topics and methodologies. It is still unclear why non-industry-funded research experienced a substantial increase in engagement with industry-funded work. One possible reason for this engagement is the collaboration between industry and non-industry entities, with industry often providing cutting-edge resources and academia providing the marketplace to identify and recruit talented researchers. [12] supports this perspective by highlighting significant industry-academia collaborations in AI research. They caution that such partnerships may narrow thematic diversity in favour of industry-preferred topics.

## J. AI Usage Card

We report how we used AI assistants such as ChatGPT and Claude for this work in the following standardized card according to [67].

Fig. 11: Percentage of outgoing citations from non-industry-funded papers (a) and non-funded papers (b) to the top 10 cited fields.
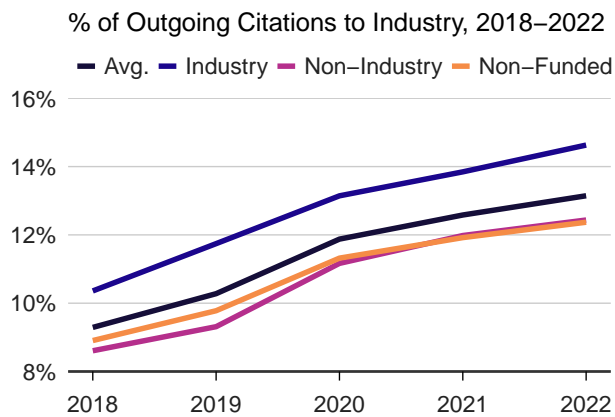


Fig. 12: The average proportion of industry-funded references across various funding types for each paper. The macro-average shows the mean percentage of industry-funded references per paper overall funding types.

# AI Usage Card

**CORRESPONDENCE(S)**
Jan Philip Wahle

**CONTACT(S)**
wahle@uni-goettingen.de

**AFFILIATION(S)**
University of Göttingen

**PROJECT NAME**
Big Tech-Funded AI Papers Have Higher Citation Impact, Greater Insularity, and Larger Recency Bias

**KEY APPLICATION(S)**
Scientometrics, Artificial Intelligence, Funding influence, Big Tech impact, Power asymmetries, Monopolization, Echo Chamber, Ethical AI

**MODEL(S)**
ChatGPT
Claude

**DATE(S) USED**
2024-04-01
2024-05-01

**VERSION(S)**
4o, 4o1
3.5 Sonnet

---

**IDEATION**

**GENERATING IDEAS, OUTLINES, AND WORK-FLOWS**
Not used

**IMPROVING EXISTING IDEAS**
Not used

**FINDING GAPS OR COMPARE ASPECTS OF IDEAS**
Not used

**LITERATURE REVIEW**

**FINDING LITERATURE**
Not used

**FINDING EXAMPLES FROM KNOWN LITERATURE**
Not used

**ADDING ADDITIONAL LITERATURE FOR EXISTING STATEMENTS AND FACTS**
Not used

**COMPARING LITERATURE**
Not used

---

**METHODOLOGY**

**PROPOSING NEW SOLUTIONS TO PROBLEMS**
Not used

**FINDING ITERATIVE OPTIMIZATIONS**
Not used

**COMPARING RELATED SOLUTIONS**
Not used

**EXPERIMENTS**

**DESIGNING NEW EXPERIMENTS**
Not used

**EDITING EXISTING EXPERIMENTS**
Not used

**FINDING, COMPARING, AND AGGREGATING RESULTS**
Not used

---

**WRITING**
ChatGPT Claude

**GENERATING NEW TEXT BASED ON INSTRUCTIONS**
Used

**ASSISTING IN IMPROVING OWN CONTENT**
Used

**PARAPHRASING RELATED WORK**
Used

**PUTTING OTHER WORKS IN PERSPECTIVE**
Not used

**PRESENTATION**

**GENERATING NEW ARTIFACTS**
Not used

**IMPROVING THE AESTHETICS OF ARTIFACTS**
Not used

**FINDING RELATIONS BETWEEN OWN OR RELATED ARTIFACTS**
Not used

---

**CODING**
ChatGPT Claude

**GENERATING NEW CODE BASED ON DESCRIPTIONS OR EXISTING CODE**
Used

**REFACTORING AND OPTIMIZING EXISTING CODE**
Used

**COMPARING ASPECTS OF EXISTING CODE**
Not used

**DATA**

**SUGGESTING NEW SOURCES FOR DATA COLLECTION**
Not used

**CLEANING, NORMALIZING, OR STANDARDIZING DATA**
Not used

**FINDING RELATIONS BETWEEN DATA AND COLLECTION METHODS**
Not used

---

**ETHICS**

**WHAT ARE THE IMPLICATIONS OF USING AI FOR THIS PROJECT?**
Generating code and improving the clarity of writing the paper has improved the efficacy of performing this scientific work.

**WHAT STEPS ARE WE TAKING TO MITIGATE ERRORS OF AI FOR THIS PROJECT?**
We manually fact-checked generated texts and inspected source code for potential generated bugs.

**WHAT STEPS ARE WE TAKING TO MINIMIZE THE CHANCE OF HARM OR INAPPROPRIATE USE OF AI FOR THIS PROJECT?**
We did not include text suggestions that had any chance of impacting marginalized groups.

**THE CORRESPONDING AUTHORS VERIFY AND AGREE WITH THE MODIFICATIONS OR GENERATIONS OF THEIR USED AI-GENERATED CONTENT**
Yes

---

# CiteAssist
## CITATION SHEET

Generated with citeassist.uni-goettingen.de
[68]

---

## BibTeX Entry

```
@article{gnewuch2025,
  author={Gnewuch, Max Martin and Wahle, Jan Philip and Ruas, Terry and Gipp, Bela},
  title={Big Tech-Funded AI Papers Have Higher Citation Impact, Greater Insularity, and Larger
      Recency Bias},
  booktitle={2026 International Conference on Artificial Intelligence, Computer, Data Sciences and
      Applications (ACDSA)},
  year={2026},
  month={02}
}
```