

# 1 Big Tech influence over AI research revisited: memetic 2 analysis of attribution of ideas to affiliation

3 Stanisław Giziński<sup>a</sup>, Paulina Kaczyńska<sup>a</sup>, Hubert Ruczyński<sup>b</sup>, Emilia  
4 Wiśnios<sup>a,c</sup>, Bartosz Pielniński<sup>d</sup>, Przemysław Biecek<sup>a,b</sup>, Julian Sienkiewicz<sup>e,\*</sup>

<sup>a</sup>*University of Warsaw, Faculty of Mathematics, Informatics and  
Mechanics, Warsaw, Poland*

<sup>b</sup>*Warsaw University of Technology, Faculty of Mathematics and Information  
Science, Warsaw, Poland*

<sup>c</sup>*NASK National Research Institute, Warsaw, Poland*

<sup>d</sup>*University of Warsaw, Faculty of Political Science and International  
Studies, Warsaw, Poland*

<sup>e</sup>*Warsaw University of Technology, Faculty of Physics, Warsaw, Poland*

---

## 5 Abstract

6 There exists a growing discourse around the domination of Big Tech on the  
7 landscape of artificial intelligence (AI) research, yet our comprehension of  
8 this phenomenon remains cursory. This paper aims to broaden and deepen  
9 our understanding of Big Tech’s reach and power within AI research. It  
10 highlights the dominance not merely in terms of sheer publication volume  
11 but rather in the propagation of new ideas or *memes*. Current studies often  
12 oversimplify the concept of influence to the share of affiliations in academic  
13 papers, typically sourced from limited databases such as arXiv or specific  
14 academic conferences.

15 The main goal of this paper is to unravel the specific nuances of such  
16 influence, determining which AI ideas are predominantly driven by Big Tech  
17 entities. By employing network and memetic analysis on AI-oriented paper  
18 abstracts and their citation network, we are able to grasp a deeper insight  
19 into this phenomenon. By utilizing two databases: *OpenAlex* and *S2ORC*,  
20 we are able to perform such analysis on a much bigger scale than previous  
21 attempts.

22 Our findings suggest that while Big Tech-affiliated papers are dispropor-  
23 tionately more cited in some areas, the most cited papers are those affiliated  
24 with both Big Tech and Academia. Focusing on the most contagious memes,  
25 their attribution to specific affiliation groups (Big Tech, Academia, mixed af-

26 filiation) seems equally distributed between those three groups. This suggests  
27 that the notion of Big Tech domination over AI research is oversimplified in  
28 the discourse.

29 *Keywords:* Knowledge Diffusion, Novelty, Affiliation Influence, Big Tech  
30 Impact, Complex Networks, Natural Language Processing

---

## 31 1. Introduction

32 Artificial Intelligence (AI) research is often seen as dominated by Big  
33 Tech companies. These companies, with their immense computing resources  
34 and access to vast amounts of data, have undoubtedly influenced the devel-  
35 opment of the field. However, while this heavy influence has been beneficial  
36 in certain aspects, it has also raised concerns that the strong position of Big  
37 Tech could influence the direction and character of AI, resulting in significant  
38 losses for society and science. One of the issues raised is the reproducibility of  
39 Big Tech’s research (Ebell et al., 2021). Given their proprietary datasets and  
40 custom software, it is often difficult to independently replicate their studies.  
41 Conflict of interest is another issue that has been raised by researchers (Ab-  
42 dalla and Abdalla, 2021; Young et al., 2022; Hagendorff and Meding, 2021).  
43 One notable case that exposed these concerns was that of Dr. Timnit Ge-  
44 bru, formerly part of Google’s AI ethics team. Her dismissal highlighted the  
45 conflicts of interest that can arise when industry-aligned research agendas  
46 collide with ethical considerations. There is also evidence that the longer  
47 researchers remain in the industry, the more their high-impact work becomes  
48 privatized, limiting the dissemination of knowledge that could otherwise ben-  
49 efit the academic community (Jurowetzki et al., 2021).

50 While there are many concerns about potential threats to research in-  
51 tegrity due to this influence, it is equally important to recognize the positive  
52 impact of Big Tech companies. First and foremost, their computational re-  
53 sources and data, previously inaccessible on such a scale, provide unprece-  
54 dented opportunities for the discovery and development of novel Machine  
55 Learning (ML) algorithms. In addition, Big Tech–Academia collaboration  
56 may benefit scientific endeavors for reasons other than access to resources.  
57 Evans (2010) argues that the industry’s relative lack of interest, in theory,  
58 could result in their academic partners being more likely to produce novel  
59 and theoretically unexpected experiments than academics working by them-  
60 selves.

61 The matters mentioned above indicate that the influence of Big Tech  
62 (and more broadly the private sector, also referred to as *company* in our  
63 paper) on AI research needs to be carefully examined. Several such studies  
64 in this regard exist, including analysis of Big Tech’s funding of academic AI  
65 researchers (Abdalla and Abdalla, 2021), changes in the proportion of Big  
66 Tech-affiliated papers at top conferences (Ahmed and Wahed, 2020), the  
67 flow of researchers from Academia to the private sector (Jurowetzki et al.,  
68 2021), the focus of the private sector on specific sub-fields of AI (Klinger  
69 et al., 2020), and the comparison of values encoded in the papers between  
70 the private sector and Academia (Birhane et al., 2021).

71 Unfortunately, approaches to this topic mentioned above are limited. The  
72 methods used so far to quantify Big Tech’s dominance tend to oversimplify  
73 the concept of influence by perceiving it as the share of affiliations in papers  
74 regarding a specific topic. Methods are limited to tools such as topic modeling  
75 and keyword analysis. Current research also suffers from a lack of diversity  
76 in data sources, relying heavily on preprint servers such as arXiv or specific  
77 conferences to collect papers. In the previous studies mentioned above, the  
78 association of papers with affiliations was binary. That is, each paper was  
79 associated with either Big Tech (or the private sector in general) or Academia,  
80 depending on the first author or the proportion of affiliations in the paper.  
81 This limits conclusions that could be drawn from studies, as the Company-  
82 Academia collaboration in the area of AI is widespread.

83 All of the above limitations, combined with strong statements (e.g. re-  
84 searchers on this topic do not shy away from comparing Big Tech to Big  
85 Tobacco, Abdalla and Abdalla, 2021), could result in a limited, and possibly  
86 biased view of the matter, which could increase the tension between Big Tech  
87 and Academia, resulting in losses for science and society. Through our study,  
88 we aim to facilitate a more nuanced understanding of the dynamics between  
89 Big Tech and Academia and to highlight potential strengths and weaknesses  
90 within this relationship. We critically evaluate existing methodologies, advo-  
91 cate for more holistic measures of influence, and explore the implications of  
92 this dominance for the field of AI. We hope that this work will foster further  
93 discourse about the symbiosis between Academia and industry, and help to  
94 shape a more balanced AI research landscape.

95 In this study, along with a standard network analysis aimed at quantifying  
96 the influence of Big Tech and Academia on the AI papers citation network, we  
97 performed memetic analysis. A meme is understood here as a piece of an idea  
98 that is transmitted through culture. In the context of research papers, this

transmission occurs through citation. Using *meme score*, which quantifies the replicating power of a specific meme, we measure the spread of particular ideas in AI research. In addition, we quantify the probability of particular memes being replicated, conditioned on the affiliation of the authors of the paper containing the meme. This approach allows us to understand how the “spreading power” of specific ideas depends on affiliation groups (Big Tech, Academia, Companies, and mixed), which sheds light on which ideas Big Tech has more influence on spreading. Moreover, we investigated how papers with joint Big Tech–Academia affiliations differ from papers authored purely by authors affiliated with one category. In summary, we address the following research questions:

1. Do papers affiliated with both Big Tech and Academia differ in citation distributions from papers affiliated only with Big Tech or Academia? Do Big Tech affiliations differ from Company ones (i.e., industrial but not attributed to Big Tech)?
2. What ideas in AI research are the most contagious?
3. Does the contagiousness of a meme differ depending on the affiliation of paper authors?
4. What ideas are more contagious when discussed by Big Tech?

## 2. Related Work

### 2.1. Big tech influence over AI research

Several studies have attempted to quantify Big Tech’s influence in AI research, using a variety of methods and data sources. Abdalla and Abdalla (2021) examine the funding provided by Big Tech to academic AI researchers. The study identifies a recurring pattern where private companies increase their support for academic institutions when their public image declines, often due to media incidents such as the Cambridge Analytica scandal. The authors compare this funding pattern to the tactics employed by the tobacco industry. The study reveals that 59% of papers published in top journals that address the ethical and societal implications of AI include at least one author with financial ties to a Big Tech company.

Ahmed and Wahed (2020) examine the changing dynamics of participation in major AI conferences following the rise of deep learning in 2012. The authors analyze 171,394 papers from 57 computer science conferences. They

133 find an increase in the participation of large technology companies, partic-  
134 ularly since 2012. In addition, the paper uses term frequency analysis of  
135 abstracts to uncover distinct research areas among different organizations.

136 [Jurowetzki et al. \(2021\)](#) use bibliographic data to measure the flow of  
137 researchers from Academia to industry and examine the factors driving it.  
138 They find that 25% of the AI researchers moved to industry from institutions  
139 at the top 5 of the Nature Index. This observation suggests that industry  
140 tends to attract AI researchers from elite institutions, possibly reflecting a  
141 search for current and potential superstar talent or a narrow focus on high-  
142 prestige sources of talent.

143 [Zhang et al. \(2019\)](#) employs altmetrics to assess the influence of AI pub-  
144 lications, revealing an increased public interest in AI research findings since  
145 2011. However, the literature suggests that altmetrics may not be suited to  
146 measuring societal impact ([Bornmann, 2014](#)).

147 The recent article of [Färber and Tampakis \(2023\)](#) strikes to measure  
148 the relationship between authors’ affiliations with the private sector and the  
149 article’s popularity, measured by citations and attention score. They perform  
150 the keyword analysis with respect to the extent of association with the private  
151 sector. Their quantitative analysis shows the domination of the private sector  
152 in the AI research domain.

153 [Krieger et al. \(2021\)](#) note an upward trend in private sector publications  
154 in fundamental research journals as opposed to those centered on applied re-  
155 search. They also highlight a rise in collaborative publications with academic  
156 entities in contrast to solo publishing efforts.

157 Several studies have tried to analyze the content of AI papers to study the  
158 influence of Big Tech. These include topic modeling of the abstracts. [Klinger  
159 et al. \(2018\)](#) combine data from arXiv, GRID (Global Research Identifier)  
160 and MAG (Microsoft Academic Graph) to create a geocoded dataset of re-  
161 search activity in computer science disciplines. They identify deep learning  
162 papers using topic modeling. Finally, the authors measure the relatedness  
163 of computer science subjects based on their co-occurrence in arXiv papers.  
164 [Klinger et al. \(2020\)](#) analyze the field of study assigned to the papers based on  
165 arXiv. They find that companies focus more on applications of deep learning  
166 and on research advancing the computational infrastructure. AI techniques  
167 outside deep learning and broader AI applications, are of less interest to the  
168 private sector. [Birhane et al. \(2021\)](#) perform textual analysis in order to ex-  
169 tract values encoded in papers. They analyze affiliations and funding sources  
170 in papers and find that the presence of Big Tech is increasing.

171 However, there is a significant gap in our understanding of Big Tech’s  
172 influence. All of the above studies have at least one of the following limita-  
173 tions.

174 Firstly, none of the studies described above link the content of papers  
175 to citations in any way. Therefore, the resulting purely fraction-based anal-  
176 ysis of content does not measure the spread of the ideas, but merely their  
177 prevalence. Quantifying both the spread and prevalence of specific ideas by  
178 using *meme score* allows us to find the most contagious ideas, called *memes*.  
179 The impact of affiliation groups on the contagiousness of each meme could  
180 be then modeled, resulting in a more fine-grained view of the influence of  
181 Big Tech, Academia, and other groups. Moreover, this approach allows for  
182 a comparison of the ideas on which contagiousness is most affected by Big  
183 Tech and Academia, which could provide insight into areas where Big Tech  
184 or Academia has the most influence.

185 Secondly, they simplify the notion of influence to the share of affiliations  
186 in papers regarding a specific topic. They use the share of papers affiliated  
187 with a particular group only (with the exception of Färber and Tampakis  
188 (2023)) that get published at prestigious conferences and in top journals as  
189 an indicator of the prevalence of Big Tech in AI research and interpret it as  
190 a proxy for influence measure. However, this approach does not take into  
191 account the number of the paper’s citations, which is a more reliable proxy  
192 for the popularity of the paper and thus closer to measuring actual influence  
193 on the research field.

194 The third limitation is the quantity and diversity of data. The use of  
195 manual annotation limits the scope of the captured categories. Similarly,  
196 using only arXiv or specific conferences as the source of the papers limits the  
197 representativeness of the data.

#### 198 2.1.1. Operationalization of the concept of Big Tech papers

199 Author affiliation is a key characteristic associated with the author rather  
200 than the research paper itself, although even a simple quantitative approach  
201 (i.e., the number of authors) has been proven to be connected to the impact  
202 exerted by an article (Sienkiewicz and Altmann, 2016). When examining the  
203 impact of Big Tech companies on scientific research, it becomes necessary to  
204 determine how to classify papers that are not exclusively authored by indi-  
205 viduals affiliated with Academia or Big Tech. Existing approaches (Klinger  
206 et al., 2020; Ahmed and Wahed, 2020; Birhane et al., 2021) typically assume  
207 that any paper with at least one author affiliated with a Big Tech company

208 should be categorized as a Big Tech paper. Alternatively, they consider the  
209 affiliation of the first author as the defining criterion.

210 However, these approaches oversimplify a more complex situation. Group-  
211 ing papers authored exclusively by researchers from the private sector, along  
212 with papers primarily written by academics, one of whom may have a dual  
213 affiliation, seems counter-intuitive. It has been observed that researchers  
214 associated with Big Tech companies benefit from access to computational  
215 resources, tools, and datasets that would otherwise be unavailable to them  
216 (Whittaker, 2021; Jurowetzki et al., 2021). Considering only the first author’s  
217 affiliation would lead to overlooking these phenomena.

218 Another approach is to look separately at papers that result from col-  
219 laborations between Big Tech companies and other institutions. Articles  
220 describing such collaborations tend to have different characteristics in terms  
221 of format and tone compared to the research papers mentioned above, and  
222 focus on the positive aspects of such collaborations (Popkin, 2019). In recent  
223 work, (Färber and Tampakis, 2023) separately analyzed papers co-authored  
224 by Academia and company-affiliated researchers. To the best of our knowl-  
225 edge, no analyses have specifically examined the influence and topics of in-  
226 terest reflected in papers resulting from collaborations between Big Tech and  
227 academia.

228 Furthermore, we are not aware of any analysis that compares Big Tech-  
229 affiliated papers to those affiliated with the private sector in general. This  
230 limits the ability to assess whether found characteristics are specific to Big  
231 Tech, or are phenomena visible in all research originating from the private  
232 sector. In our study, we examine the differences between several affiliation  
233 groups: purely Big Tech, purely academic, private sector as a whole, and  
234 papers affiliated with both Big Tech/private sector and Academia.

## 235 2.2. Memetic analysis

236 Several measures can characterize the spread of ideas in science (Wag-  
237 ner et al., 2011; Xu et al., 2018) or in business research (Wu et al., 2017).  
238 A straightforward method to identify ideas within a set of documents is to  
239 examine their frequency. Techniques like tf-idf (term frequency–inverse doc-  
240 ument frequency) allow the discovery of potential keywords. Topic modeling  
241 methods, such as Latent Dirichlet Allocation (LDA) or BERTopic (Groo-  
242 tendorst, 2022), help to understand the distribution of ideas consisting of  
243 multiple phrases across a document set. Another option is to focus on top-  
244 ical clusters, which can be built from bibliographic coupling and co-citation



245 networks (Liu et al., 2017, 2019). One can use a citations cascade (Min et al.,  
246 2021) or chains of citations (della Briotta Parolo et al., 2020), which reveal a  
247 higher-order (i.e., reaching further than just the neighborhood of the node)  
248 temporal relationship between the papers.

249 One perspective on the ideas’ propagation is looking at it through the  
250 lens of the evolutionary theory of science (Kantorovich, 2014). Central to  
251 this theory is the concept of *meme*, first introduced in a broader context by  
252 Dawkins (1976). A meme is considered a cultural analog to a biological gene,  
253 a unit of information that replicates itself, mutates, and undergoes selection  
254 in the evolutionary process. Analogous to genes in biology, memes are pieces  
255 of information transmitted between individuals within a given culture.

256 Kuhn et al. (2014) proposed an operationalization of the concept of a  
257 meme in the context of scientific knowledge by introducing the notion of a  
258 *meme score*. The meme score measures how much a given term is a meme by  
259 analyzing the citation network. It takes into account the term’s frequency,  
260 the probability of its transmission through citation, and the probability of it  
261 emerging independently. In this operationalization of memes, a term exhibits  
262 more memetic quality, the higher the meme score of that term is. Since its  
263 proposal, the meme score has been used in several studies. It has been used  
264 to investigate how the gender of researchers influences their positions in the  
265 scientific community (Araújo and Fontainha, 2018), to identify scientific and  
266 technological trajectories of ideas in paper and patent networks (Sun and  
267 Ding, 2018), and to explore diffusion cascades of ideas (Mao et al., 2020).

### 268 3. Materials and methods

#### 269 3.1. Meme score

270 The *meme score* was introduced by Kuhn et al. (2014) as a measure of  
271 the ‘contagiousness’ of meme.

272 The meme score for the given meme is calculated by multiplying the  
273 relative frequency of a term by its *sticking factor* and then dividing it by its  
274 *sparking factor*. The quotient of the sticking factor and the sparking factor  
275 is labeled as a propagation score. Both measures are briefly described below.

276 The sticking factor can be interpreted as the probability of the meme  $m$   
277 being transmitted from one paper to the other paper. It corresponds to the  
278 probability that the meme appears in the abstract of the given paper, given  
279 that it appears in at least one of the abstracts of the papers cited by it. It is  
280 calculated by dividing the number of papers that contain the meme  $m$  and



281 cite at least one paper with this meme by the number of papers that cite at  
 282 least one paper with this meme:  $d_{m \rightarrow m}$  by  $d_{\rightarrow m}$ .

283 The sparking factor can be interpreted as the probability of the meme  
 284 appearing without being present in any of the cited papers - how likely it  
 285 is that the meme will be mentioned in the abstract, given that it does not  
 286 appear in any of the cited papers' abstracts. It is calculated by dividing the  
 287 number of papers with the given meme that do not cite any papers with this  
 288 meme  $d_{m \rightarrow \mathcal{M}}$  by  $d_{\rightarrow \mathcal{M}}$  - all papers that do not cite any papers with the given  
 289 meme.

290 The formula for the meme score can be written as follows:

$$M_m = \frac{N_m}{N} \frac{d_{m \rightarrow m}}{d_{\rightarrow m}} \setminus \frac{d_{m \rightarrow \mathcal{M}}}{d_{\rightarrow \mathcal{M}}} \quad (1)$$

291 where  $N_m$  is the number of papers with the meme  $m$ ,  $N$  is the number  
 292 of all memes in the dataset.  $d_{m \rightarrow m}$  is the number of papers that contain the  
 293 meme and cite at least one paper with the meme and  $d_{\rightarrow m}$  is the number of  
 294 papers that cite at least one paper with the meme. On the other hand  $d_{m \rightarrow \mathcal{M}}$   
 295 is the number of papers containing the meme but not citing a paper with the  
 296 meme, and  $d_{\rightarrow \mathcal{M}}$  is the number of papers that do not cite a paper with the  
 297 given meme.

298 From the three components of the meme score, it is mainly the sticking  
 299 factor that measures the meme's contagiousness. The other factors control if  
 300 the phrase is not merely used as a part of a natural language or if a phrase  
 301 is popular enough. Due to this, the sticking factor on its own could serve as  
 302 a weaker measure of contagiousness.

303 Let us underline that meme score can be interpreted as a ratio of the  
 304 probability of replicating a phrase to the probability of the probability of its  
 305 appearance without earlier presence in the cited papers. To select meaningful  
 306 memes, we must select a threshold of the meme score defining which phrases  
 307 are considered memes and which are not. This is done by selecting the point  
 308 of maximum curvature from the value of the meme score as a function of  
 309 its position in the phrase ranking ordered by the meme score (see Sec. 4.2  
 310 for practical implementation of referred method in this study). Such an  
 311 approach is similar to the so-called "elbow method" (Ketchen and Shook,  
 312 1996) – a commonly used visual technique in assessing the optimal number  
 313 of clusters in the k-means method.

### 3.2. Conditioned sticking factor

To measure how probable a meme is to spread under a given affiliation, we introduce a *conditioned sticking factor* – a modified version of the meme score introduced by Kuhn et al. A similar approach was proposed in the paper of Araújo and Fontainha (2018), where the authors attempted to condition the meme score on the main author’s gender. In contrast, in our approach, we condition only the sticking factor on the authors’ affiliations. The conditioned sticking factor is calculated as

$$\sigma_{m,a} = \frac{d_{m \rightarrow m,a}}{d_{\rightarrow m,a}}. \quad (2)$$

The conditioned sticking factor divides the number of papers  $d_{m \rightarrow m,a}$  that have meme  $m$  and cite at least one paper with this meme and given affiliation  $a$  by the number of papers  $d_{\rightarrow m,a}$  that cite at least one paper with given meme  $m$  and affiliation  $a$ . It can be interpreted as the probability of the meme transitioning during citation under the condition that the meme appeared in the paper with a given affiliation.

In contrast to the earlier approach of Araújo and Fontainha (2018), we have decided not to condition the sparking factor. The original purpose of the sparking factor is to control how often the word appears without any influence and limit the appearance of the most popular words that do not convey specific ideas. The sparking factor, conditioned analogously to the sticking factor, would correspond to the probability that the meme appears in a paper without being present in any cited papers with a given affiliation. However, the meme can appear in the other cited papers with different affiliations. In this way, the conditioned sparking factor would lose its original purpose of excluding phrases not conveying ideas without providing us with any specific additional information. For this reason, we decided to condition only the sticking factor. Finally, once  $N_m \gg d_{m \rightarrow m}$ , i.e., the meme is relatively common but rarely transmitted through the citation network, the meme score can be approximated by the sticking factor. Due to this, we focus primarily on the relationships between conditioned sticking factors for different groups. Figure 1 presents a toy example to explain how the sticking factor, sparking factor, meme score and conditioned sticking factor are calculated.

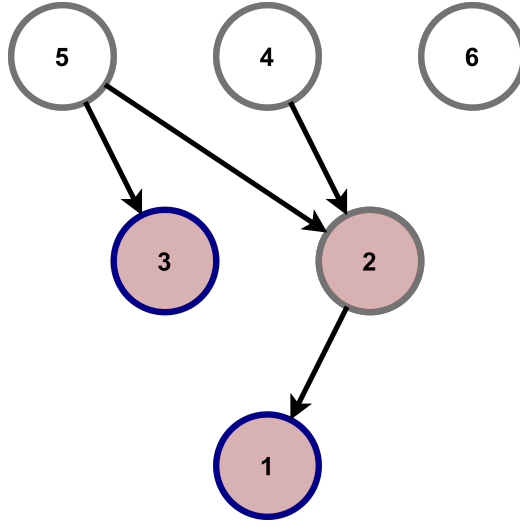


Figure 1: **Toy example of how to calculate the meme score.** Papers 1, 2 and 3 contain the meme (maroon background). Papers 1 and 3 have a given affiliation (dark blue frame) for which we calculate the conditioned sticking factor. The frequency of the meme is  $\frac{3}{6}$ . The unconditioned sticking factor is  $\frac{1}{3}$ , because of three papers that cite papers with the meme (2, 4 and 5), only one has the meme. The sparking factor is  $\frac{2}{3}$ , since the meme appears in two of three papers that do not cite papers with this meme – it appears in 1 and 3, but not in 6. The conditioned sticking factor is  $\frac{1}{2}$ , since meme appears in 1 out of 2 citations of the paper with the meme and given affiliation (2 replicates the meme from 1, 5 does not replicate the meme from 3).

### 3.3. Papers dataset and processing

As a main data source for this study, we used the S2ORC database [Lo et al. \(2020\)](#), which is a large corpus of 81.1 million English-language academic papers spanning many disciplines. The collection comes from various sources, such as in-proceedings of scientific conferences, well-established journals, as well as digital archives (e.g. arXiv). We gathered 557 681 articles from this source by searching for phrases presented in [Appendix A](#)) which occurred in the titles and abstracts of the papers. The aforementioned set of keywords is proposed in [Liu et al. \(2021\)](#). The authors of this paper conducted the study, attempting to find a representative set of keywords for the artificial intelligence domain while maintaining a reasonable balance between recall and precision.

Afterwards, we performed exploratory data analysis. We found out that 138 878 records had no DOI identifiers, which we used to match the papers against the OpenAlex database. As affiliation data provided by those links is crucial for our study, we decided to remove such papers from our dataset. Next, we discovered that 146 525 of the remaining records had no information about citations. Papers like this were useless for citation network analysis, so we decided to remove them from our dataset. Finally, we ended up with 166 455 records that had affiliations matched with OpenAlex.

The whole pipeline of memes extraction is shown in [Fig. A.10](#). All used abstracts were preprocessed first – we lowercased all words and tokenized them. For the extraction of noun chunks, we used the `spaCy` library with the `en_core_web_md` model.

#### Papers' affiliation

First, we look at the category of the company as it is specified in the OpenAlex Database. OpenAlex provides information about the type of institution – it categorizes them into the following categories: company, education, government, facility, healthcare, nonprofit, and other ([Table B.2](#)). By analyzing these categories, we decided to merge education, facility, and government into Academia. We also analyzed the company category separately.

As Big Tech, we take a subset of technological companies that were among the top 60 companies by market capitalization in the years 2016-2019 ([PwC, 2019, 2018, 2017, 2016](#)). The following companies fall into this category: Apple, Google, Microsoft, Facebook, Tencent, Oracle, Intel, IBM, Cisco Systems, TSMC, SAP, Qualcomm, Amazon, Siemens, Alibaba Group, Nvidia, and Samsung. We aimed to construct a list that encompasses companies

previously labelled as Big Techs in media, such as the Big Five, FAANG or MAGA (Economist, 2018), while also being supported by a quantitative measure like market capitalization. To ensure the list’s robustness and account for short-term market changes, we use a union of companies that appeared in the top 60 companies during the four years with the highest number of papers in our dataset. We opt to limit ourselves to the top 60 companies to include companies considered most significant to the field while excluding the lower end of the top 100 companies, where most technological companies appeared only once during the four years we consider.

Papers affiliated with Big Tech are 38% of the papers affiliated with the company category, but they account for 57% of their citations.

## 4. Results

### 4.1. Non-binary distinction between Big Tech and academy-affiliated papers

We provide substantial arguments, advocating for the less coarse distinction between Academia and Big Tech affiliations than the bisection often used in the literature. The most natural way to follow this proposition would be to use a continuous space by simply calculating for each paper the proportion of Big Tech affiliations to the total number of affiliations. However, as shown in Fig. 2, these values are unevenly distributed – the vast majority are simply academic papers (almost 98% of all papers in the examined dataset), around 0.5% belong only to the Big Tech class, and the remaining part is characterized by a mixed affiliation. The last group also suffers from obvious finite-size effects, i.e., due to a limited number of authors, fractions such as  $\frac{1}{4}$ ,  $\frac{1}{3}$ ,  $\frac{1}{2}$ ,  $\frac{2}{3}$ , and  $\frac{3}{4}$  tend to occur more frequently than other values (see Fig. 2). To avoid the above obstacles and to be consistent with the observations presented in Fig. 2, we decided to use a *ternary* affiliation scheme, i.e., Academia–mixed Big Tech–Big Tech.

Let us first explore the implications of the ternary division using simple network metrics: incoming degree and PageRank (Brin and Page, 1998). Both these measures can be used to quantify the importance of a node (paper). In the case of the in-degree, we simply tally the number of citations a paper has received, while PageRank provides more nuanced information. Figure 3 shows distributions of in-degree (panel a) and PageRank (panel c), respectively, grouped according to the introduced ternary affiliation classification (denoted as “Academia”, “Mixed” and “Big Tech”). In both cases, the Kruskal-Wallis test (Kruskal and Wallis, 1952) finds significant differences

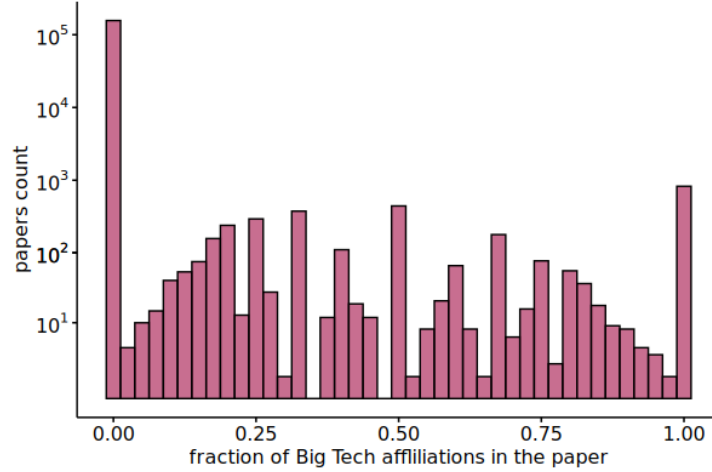


Figure 2: **Histogram (log scale) of the number of papers for different fractions of Big Tech affiliations.** The length of each bin is 0.025. Most papers are affiliated only with Academia, yet we can also witness a visible peak for purely Big Tech papers, which underlines the validity of this group. The plot also suggests that we should consider a third group of mixed affiliations, as the number of articles with Big Tech fraction is even bigger than that of Big Tech papers.

among all groups (p-value < .001), while post-hoc pairwise comparisons using Dunn’s test (Dunn, 1964) with Holm-Bonferroni corrections (Holm, 1979) indicate that the described distinctions are maintained for all pairs of groups (if two or more categories are statistically indistinguishable, they are connected with a solid line on the far right part of the panel). To prove that the division into three categories does not lead to a loss of additional information available in the continuous approach, we divide the dataset into six categories: the first two and the last one are the same as in the previous division, while the remaining three describe cumulative quartiles without exclusive Big Tech affiliations (“>25%”, “>50%” and “>75%”, i.e., “>25%” means that each node in this category has over 1/4 share of Big Tech affiliations but cannot be exclusively Big Tech). Indeed, also at this level, we obtain similar results for the Kruskal-Wallis test (significant differences with p-value < .001) as well as for pairwise comparisons, the only exception being the penultimate category (“>75%”), which is indistinguishable from the pure Big Tech affiliations in the case of PageRank. All the intermediate categories (“Mixed”, “>25%”, “>50%” and “>75%”) are indistinguishable from each other, which proves that performing further divisions in the “Mixed”

category does not bring any additional insights.

To examine the reasons for the clear differences between the elements of our ternary classification, we focus on the in-degree and PageRank probability distributions shown in Fig. 3b and Fig. 3d, respectively. Clearly, mixed and Big Tech groups differ from Academia by being characterized by higher values of probability density for larger values of in-degree or PageRank. On the other hand, the difference between mixed and Big Tech is not as pronounced. The likely main reason for the observed differences is revealed when we focus on the number of zero-degree nodes – the proportion of such nodes in Academia, mixed and Big Tech is 0.648, 0.538, and 0.613, respectively. Although the ratio of zero-degree nodes in the mixed category is visibly lower than in the other two categories, we confirm this observation by performing pairwise two-sample proportion tests with Holm-Bonferroni adjustment, which give p-values of  $p < .0001$  for Academia–mixed Big Tech,  $p = .044$  for Academia – Big Tech, and  $p = .0004$  for mixed Big Tech– Big Tech. In effect, the network with removed zero-degree (isolated) nodes brings different statistical results marked in Fig. 3a and Fig. 3c by the dotted line on the far right – in this case we recover the binary division as only “Academia” category is distinguishable from any other.

Lastly, the decision to introduce the Big Tech category instead of looking at the OpenAlex company category was supported by the intuition that Big Techs differ significantly from other company actors and that extrapolating conclusions from all private companies to Big Tech would be unjustified. To test whether Big Tech are indeed different from the rest of the private companies, we compare the PageRanks of Big Tech–affiliated and company-affiliated papers. We perform the Kruskal-Wallis test between the PageRanks of Big Tech–affiliated papers and papers affiliated with companies that are not Big Tech. The test results confirm that the distribution of PageRanks of Big Tech-affiliated papers is indeed different from that of other companies (p-value  $< 0.05$ ). Nevertheless, analyzing the ternary division from the company’s point of view (instead of Big Tech) does not bring any additional insights, as can be seen in Fig. 4.

#### 4.2. Most contagious ideas in AI research

Using the meme score measure, we were able to identify and analyze which ideas are the most prevalent in AI research in general, which is the first step in differentiating Big Tech from Academia.



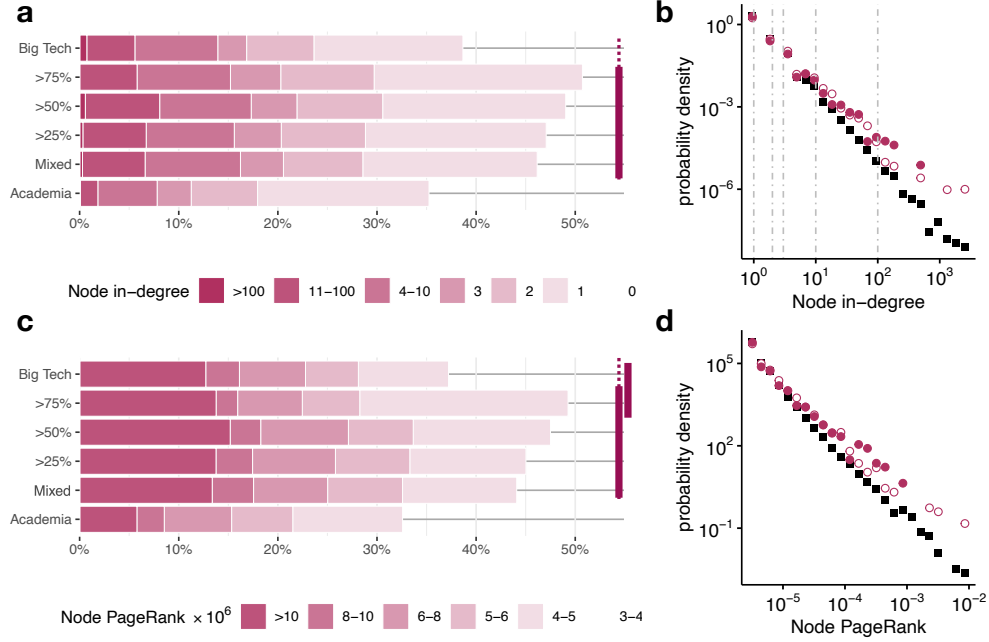


Figure 3: **Differences between Academia-mixed Big Tech-Big Tech as seen from the citation network perspective.** The rows are, respectively in-degrees of the nodes (panels a, b) and the PageRank values of the nodes (panels c, d). The left column shows the distribution plot, presenting the percentage of nodes with a given in-degree or PageRank value for ternary classification (“Academia”, “Mixed”, and “Big Tech”) as well as cumulative quartiles (see text for details). Vertical solid lines on the far right of panels a and c connect statistically indistinguishable categories, while vertical dotted lines extend solid lines for the case of the network with removed isolated nodes. The right column presents probability density distributions associated with the ternary classification on panels a and b: filled squares represent Academia, empty circles – mixed Big Tech, and filled circles – Big Tech (node in-degree is increased by one to overcome log scale issues).

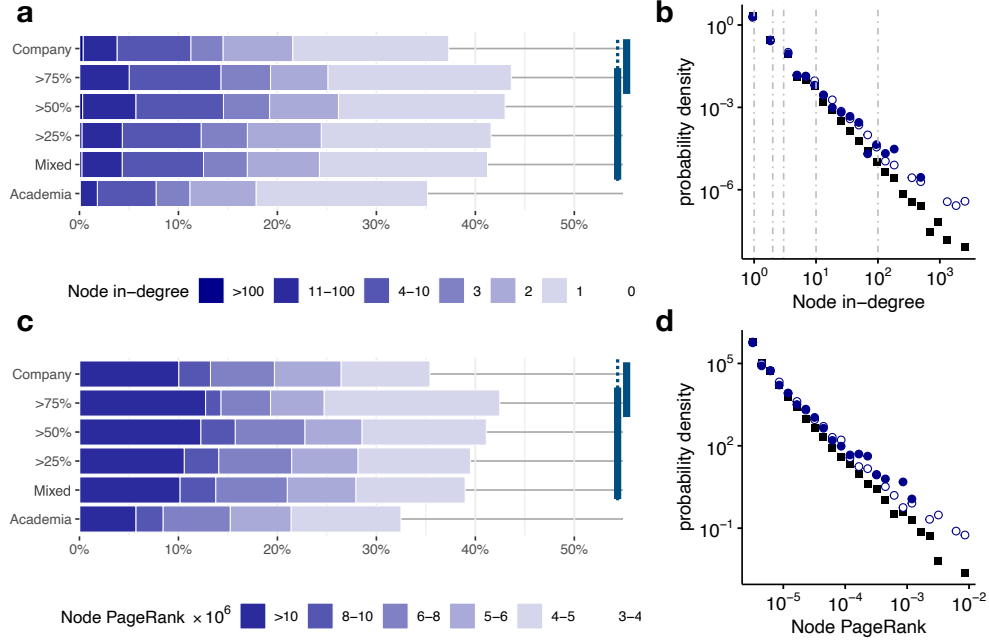


Figure 4: **Differences between Academia-mixed Company-Company as seen from the citation network perspective.** The rows are, respectively, in-degrees of the nodes (panels a, b) and the PageRank values of the nodes (panels c, d). The left column shows the distribution plot, presenting the percentage of nodes with a given in-degree or PageRank value for ternary classification (“Academia”, “Mixed”, and “Company”) as well as cumulative quartiles (see text for details). Vertical solid lines on the far right of panels a and c connect statistically indistinguishable categories, while vertical dotted lines extend solid lines for the case of the network with removed isolated nodes. The right column presents probability density distributions associated with the ternary classification on panels a and b: filled squares represent Academia, empty circles – mixed Company, and filled circles – Company (node in-degree is increased by one to overcome log scale issues).

473 To perform this analysis, we calculated the meme scores and propagation  
 474 scores and counted the occurrences for all memes present in our dataset,  
 475 which resulted in over 1.7 million observations. Such a vast amount of data  
 476 had to be filtered to get the most interesting results, which is why we decided  
 477 to limit our study to phrases where the meme score was greater than 0, which  
 478 resulted in over 60,000 meaningful observations.

479 Since our main goal is to deeply analyze the most spreading ideas, we  
 480 limited the number of memes according to two criteria: the observed occur-  
 481 rences and their meme score values. The first criterion is based on the idea  
 482 that the memes that have high meme scores but only occur once or twice are  
 483 not really meaningful topics, as they are too niche to be considered conta-  
 484 gious. The second condition is obvious since the higher the meme score, the  
 485 more likely it is to spread. To choose a cutoff value, we analyzed the curves  
 486 shown in Fig. 5. Finally, after limiting ourselves to phrases with a meme  
 487 score above 0.25 and at least 20 occurrences, we ended up with 251 observa-  
 488 tions representing the top memes in our dataset. **A manual annotation and**  
 489 **inspection of the selected memes yields 80% precision, which is a comparable**  
 490 **result to the one reported by Kuhn et al. 2014 (i.e., around 81.2% on 150**  
 491 **memes).**

492 The set of chosen top memes was later annotated by the group of domain  
 493 experts, which is briefly described in [Appendix C](#), resulting in the assignment  
 494 of one category per observation. The resulting categories are *ML algorithms*,  
 495 *Medical terms*, *Security*, *Niche ML applications*, *ML Concept*, *Maths*, *Graphs*,  
 496 *Computer Vision*, *Natural Language Processing*, *Data Related*, and *Other*.  
 497 The annotations of the top memes allow us to analyze which high-level topics  
 498 are popular among ML researchers, indicating the most promising research  
 499 areas. According to Table 1, the most popular topic among scientists is  
 500 the medical applications of ML. This category has the most memes and a  
 501 third number of the memes occurrences in our dataset. It is not the only  
 502 application, however, as we can see on the top memes list other important  
 503 topics like security, and a whole category describing niche applications of  
 504 ML solutions. However, this doesn’t mean that today’s research focuses  
 505 only on ML applications, as the second-highest average meme score value,  
 506 and the second sum of the occurrences belong to the topic describing the  
 507 ML algorithms. Additionally, we were able to distinguish some particularly  
 508 popular directions of AI development, such as Computer Vision or NLP.  
 509 Finally, such an in-depth analysis proves the quality of our framework as  
 510 the category Other, which includes unrelated topics, as well as meaningless

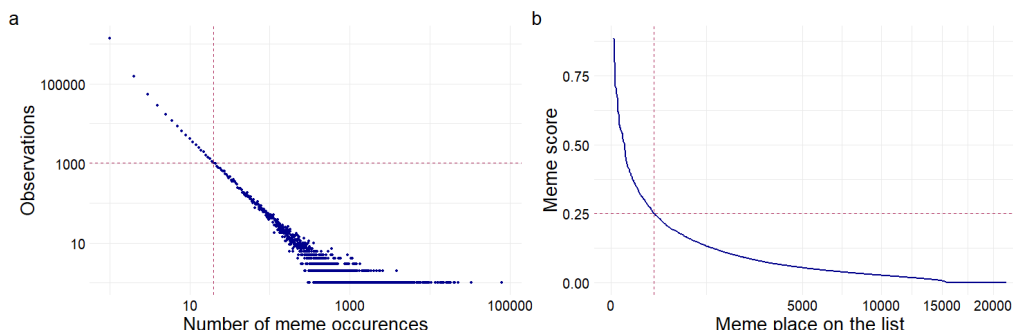
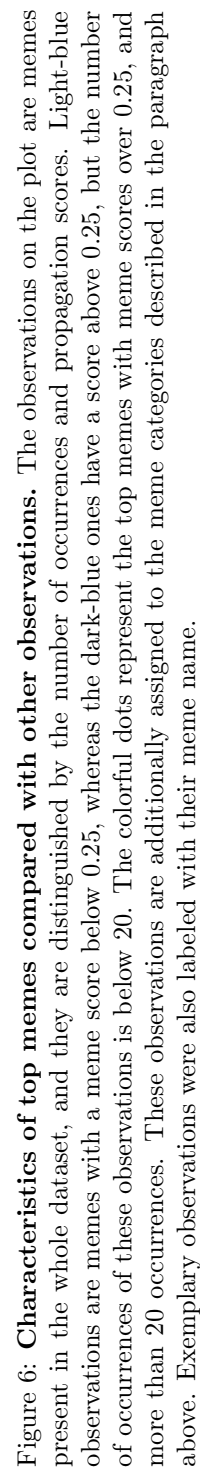


Figure 5: **Top memes selection by thresholds.** The left plot shows the number of remaining phrases, depending on the minimum number of meme occurrences. Since the dispersion of observations occurs around the number 20 on the x-axis, and 1000 on the y-axis (marked as dashed lines), we decided to use this point as a cutoff, which resulted in limiting ourselves to over 20 000 observations that have more than 20 appearances. The right plot shows the meme score value of the observation as a function of its position in the list sorted by the meme score. Since the major flattening occurs around the meme score value of 0.25, we decided to use this point as a cutoff, which resulted in limiting ourselves to 251 observations (marked as dashed lines).

511 memes, is responsible for only 20% of detected memes. The group has the  
 512 biggest number of occurrences due to the memes such as ‘paper’, ‘results’,  
 513 or ‘art’ (part of ‘state of the art’ phrase), which occur multiple times in the  
 514 majority of the papers.

515 As we already know the big picture of the main topics in the scientist’s  
 516 interests, we can combine it with a more detailed analysis, and provide a  
 517 wider perspective for our top memes group.

518 Figure 6 shows the characteristics of the top memes subset. As the meme  
 519 score value depends on the number of memes occurrences as well as the prop-  
 520 agation score we can easily imagine that the line that stands for the meme  
 521 score equal to 0.25 is used as a cutoff value in this analysis. From the den-  
 522 sity of observations, especially in our group, we can grasp the tendency that  
 523 the highest meme scores were assigned for the observations with a smaller  
 524 number of occurrences. Additionally, the plot helps us detect the outliers  
 525 aforementioned in the previous section, as they yield low propagation score,  
 526 but also have a large number of occurrences. Thanks to the exemplary meme  
 527 annotations we can also track some memes and analyze the reasons why are  
 528 they considered important. For example, ‘fairness’ is considered important



Category	Number of memes	Mean meme score	Mean occurrences	Sum of occurrences
Medical term	70	0.339	141	9865
Other	50	0.347	2620	130993
Niche ML application	24	0.345	56	1336
Security	23	0.346	107	2454
ML Algorithm	20	0.380	581	11616
ML Concept	20	0.343	73	1456
Computer Vision	17	0.331	79	1336
Maths	9	0.348	35	312
Data Related	7	0.405	48	338
NLP	7	0.332	204	1430
Graphs	4	0.384	47	188

Table 1: **A statistical summary of the topics present in the top memes subset.** The table provides information about the number of memes existing within the described category. It also shows additional information describing the average meme score value, the average number of meme occurrences, and the sum of the total occurrences of memes from that topic.

529 mostly due to its high propagation score, whereas ‘cnn’ or ‘deep learning’  
530 benefits from a vast number of occurrences.

531 To finalize the topic analysis, we decided to present all top memes in the  
532 form of a structured word cloud presented in Fig. 7. This way, we are able  
533 to analyze particular memes that build each of the predefined categories. We  
534 should particularly focus on the medical memes, as this topic is definitely the  
535 most important one, due to the abundance of different memes. The memes  
536 creating this topic (presented as orange), cover various medical terms. We  
537 can discover which disease treatments and studies greatly benefit from the  
538 usage of ML (ex. various cancers, epilepsy, diabetic retinopathy, Alzheimer’s  
539 disease, malaria, etc.). We can also say that another major points of interest  
540 are the studies considering genomics (lncRNA, miRNAs, microRNAs, etc.),  
541 and that ML is widely used for various examinations (mammograms, eeg,  
542 eeg, wireless capsule endoscopy, etc.).



22



### 543 4.3. Differences in contagiousness between Companies, Academia, and Big 544 Tech

545 To assess whether contagiousness varies depending on author affiliation,  
546 we conducted a comparative analysis of conditioned sticking factors across  
547 different affiliation groups. A selection of the 15 highest-ranked memes in  
548 each affiliation (Academy, mixed Big Tech, Big Tech, mixed Company, and  
549 Company) ordered according to their conditioned sticking factor is shown  
550 in Fig. 8a. Although this limited picture suggests that *Academia* memes  
551 dominate, we need to underline that this is simply an illustrative excerpt  
552 from the whole distribution.

553 In order to perform a more systematical analysis, we compared all non-  
554 zero conditioned sticking factors across all affiliation groups using Kruskal-  
555 Wallis test. Results indicate a significant difference across groups in terms of  
556 the sticking factor, i.e. affiliation is a factor that creates distinctions among  
557 sticking factor values. Outcomes of post-hoc comparison performed using  
558 Dunn’s test between categories presented in Fig. 8b pinpoint distinctions  
559 among the sticking factor distributions: strikingly, there are no statistical  
560 differences among Academia, Big Tech, and Company. On the other hand,  
561 such pairs as mixed Big Tech – Big Tech or mixed Company – Company can  
562 be distinguished. This is exactly opposite to the relations observed using  
563 network metrics for non-isolated nodes and shown in Figs. 3 and 4, where  
564 differences between mixed Big Tech – Big Tech and mixed Company – Com-  
565 pany were negligible while one could clearly distinguish Academia, Big Tech,  
566 and Company categories.

567 The above analysis is conducted in absolute terms, i.e., the distributions  
568 can be biased by the fact that Academia affiliations are exceedingly more  
569 frequent than Big Tech or Company ones (cf Fig. 2). To overcome this  
570 issue we focused on memes that appear in at least two different affiliation  
571 categories, e.g., “random neural network” (present in mixed Company, Com-  
572 pany, mixed Big Tech, and Big Tech) or “captcha” (present in Company,  
573 mixed Big Tech, and Big Tech) depicted in Fig. 8a. Following, we conducted  
574 Wilcoxon signed-rank tests (Wilcoxon, 1945), comparing sticking factors con-  
575 ditioned on specific groups using matched samples (the same memes in both  
576 groups). In case we were able to reject the null hypothesis (observations in  
577 both groups  $X$  and  $Y$  are exchangeable) in favor of the one-sided alternative  
578 hypothesis (differences  $X - Y$  are stochastically larger than a distribution  
579 symmetric about zero) we noted it with a directed link in Fig. 8c, point-  
580 ing from  $X$  to  $Y$ . The results revealed that memes affiliated with Big Tech

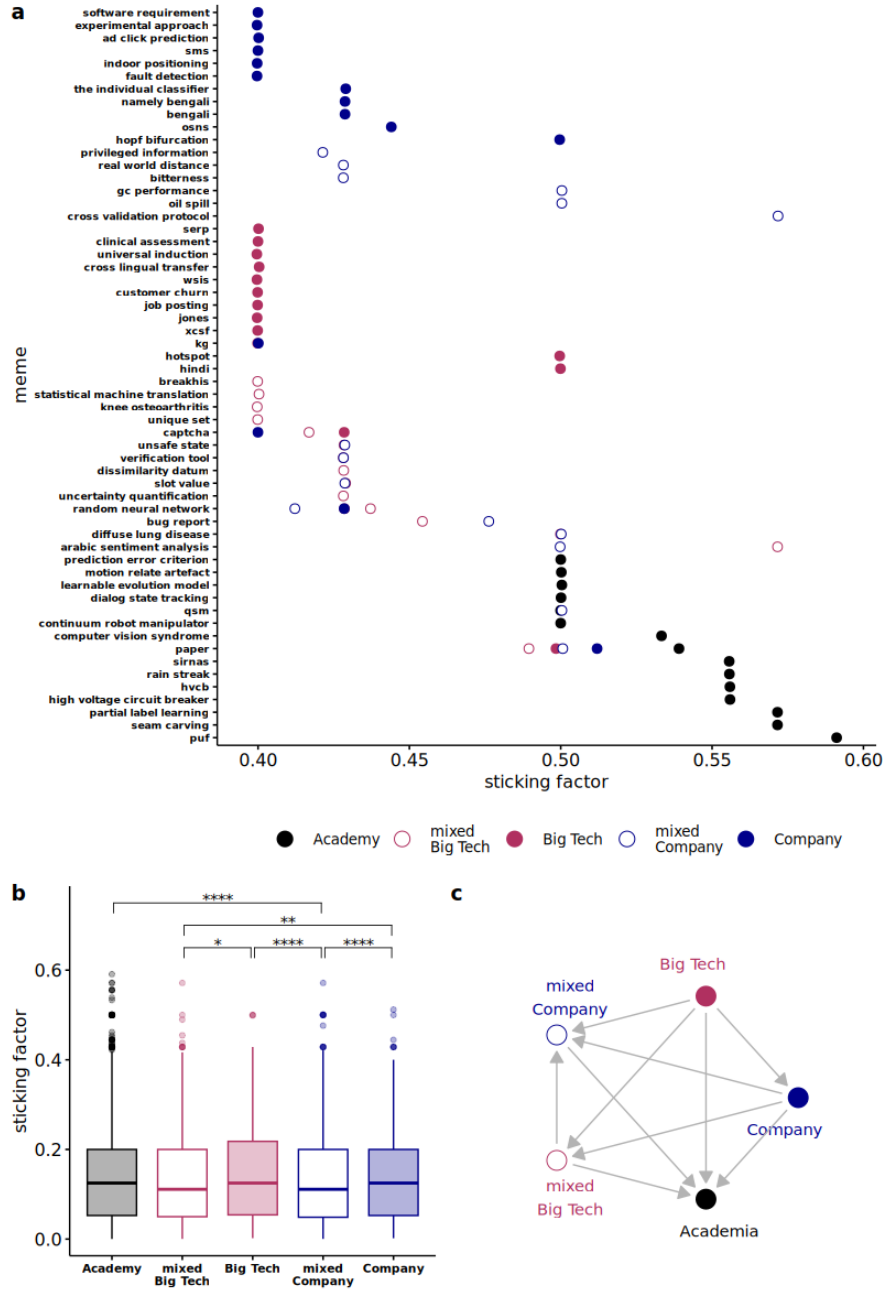


Figure 8: **Differences in contagiousness between Companies, Academia, and Big Tech.** (a) Selection of 15 memes with the highest sticking factor values in each category. (b) Boxplots of sticking factors conditioned on affiliation. The significance of post-hoc comparisons is indicated by \*, \*\*, \*\*\*\*, representing  $.01 < p < .05$ ,  $.001 < p < .01$  and  $p < .0001$ , respectively – insignificant comparisons are not shown. (c) Results of Wilcoxon tests for paired memes – an incoming arrow (pointing from category  $X$  to category  $Y$ ) informs about statistically higher sticking factors of paired memes in  $X$  than in  $Y$ .

581 exhibited higher sticking factors compared to all other affiliation categories,  
582 with company-affiliated memes ranking second. Academic affiliation, on the  
583 other hand, was associated with a decrease in the sticking factor compared  
584 to other affiliations.

585 Following this, we explore how this difference depends on the meme cat-  
586 egory introduced in Table 1. Five hundred memes with the highest meme  
587 scores, excluding those not present at least once in papers from each affil-  
588 iation group, were selected for this comparison. After removing duplicates  
589 and the least common categories of memes, the resulting list has 148 memes.  
590 The resulting comparison of conditioned sticking factors divided by meme  
591 category is presented in Fig. 9.

## 592 5. Discussion and conclusions

593 This work aimed to identify the most contagious ideas (memes) in AI  
594 research and assess to what extent the contagiousness of memes depends on  
595 the affiliation group of authors of the paper, and whether this contagiousness  
596 differs between pure and mixed academic or company affiliation. For this  
597 purpose, we analyzed AI papers’ abstracts, affiliations, and citation networks.  
598 We employed two key metrics, namely the meme score and the sticking factor,  
599 to quantify the contagiousness of these memes. Additionally, we introduced  
600 the concept of conditioned sticking factor to assess the relationship between  
601 affiliations and contagiousness.

602 Our findings reveal several noteworthy insights. Firstly, we observed  
603 that papers affiliated with Big Tech tend to receive disproportionate ci-  
604 tation counts in comparison to all companies. However, the most highly  
605 cited papers are those co-authored by individuals from both Big Tech and  
606 academic institutions, i.e., papers having mixed affiliations. Additionally,  
607 papers co-authored by individuals from both corporate and academic back-  
608 grounds exhibit scientometric distinctions from papers authored exclusively  
609 by Academia or corporate entities. This pattern holds true for both Big  
610 Tech and Company, thus challenging the binary distinctions often assumed  
611 in prior research.

612 The analysis of contagiousness indicates that on average, the contagio-  
613 usness of Academia-authored memes does not differ from those authored by  
614 companies or Big Tech, although when taking into account only memes  
615 present in both affiliations, the same memes have on average higher conta-  
616 giousness when mentioned by Big Tech than Academia. These facts together

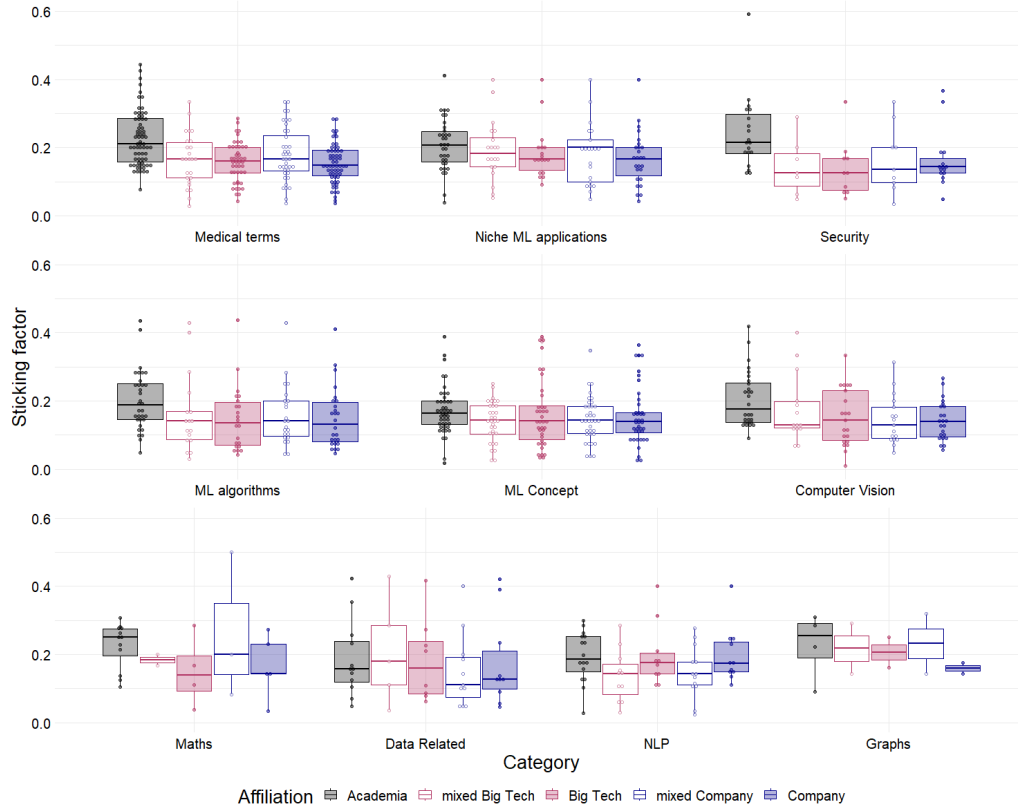


Figure 9: **Conditioned sticking factor across ML categories.** Box-plot depicts the distribution of the conditioned sticking factor of memes across different ML categories, stratified by the originating affiliation. The distribution patterns for most categories mirror the aggregate trends. In the ‘Data Related’ segment, memes generated by Big Tech exhibit a noticeably higher sticking factor, suggesting a dominant role in propagating novel content within this domain—likely a reflection of the substantial resources they possess. For ‘Niche ML Applications’, companies have a higher median sticking factor compared to their performance in other meme categories, implying their effectiveness in spreading memes that may complement their commercial strategies.

617 could be linked with the difference in the diversity of memes produced and  
618 propagated by those two groups. Academia creates far more memes than  
619 Big Tech or Companies (mostly due to the difference in overall papers au-  
620 thored), resulting in mean contagiousness being not significantly different  
621 than Big Tech or Company. When comparing common ideas of interest of  
622 those two groups, the difference in contagiousness becomes apparent. Those  
623 could be due to multiple factors, such as bigger resources available for Big  
624 Tech’s researchers.

625 The impact of affiliation group over the contagiousness of memes could  
626 also be understood by applying the *framing theory*. Framing theory states  
627 that an issue can be perceived from various perspectives, each carrying impli-  
628 cations for multiple values or considerations (Chong and Druckman, 2007).  
629 In the context of memetics, we can assert that each meme possesses the po-  
630 tential for diverse framings. These framings, tied to the affiliations of the  
631 authors, become integral components that propagate alongside the meme.  
632 Consequently, the contagiousness of a meme, when conditioned on a specific  
633 group, signifies the extent of that group’s influence over the framing of the  
634 respective meme. In essence, the more contagious a meme becomes within a  
635 particular group, the more pronounced is that group’s control over the fram-  
636 ing associated with the meme. To assess this further, one could analyze the  
637 co-occurrence of memes, and ties of such memes “mixtures” to affiliations.

638 The analysis we have presented confronts the influence of Big Tech with  
639 that of academic institutions on the propagation of ideas between scientific  
640 papers. In doing so, we offer a unique perspective on tracing the evolution of  
641 ideas influenced by major technological companies and academic institutions.  
642 The conditioned sticking factor, introduced in this work, represents a novel  
643 tool that can extend beyond the confines of this study. It can be employed to  
644 analyze how various modalities, such as country of origin, journal, or insti-  
645 tution, impact the spread of ideas in different contexts. In this way, possible  
646 outcomes could complement or challenge the results obtained in previous  
647 studies that tackled the problem of diffusion of ideas in a given scientific dis-  
648 cipline (Hargreaves Heap and Parikh, 2005), or technology adoption under  
649 different (also institutional) constraints (Galang, 2014). In the same manner,  
650 if provided with additional data, the elaborated approach could be of use to  
651 track relations between organization characteristics and academia-industry  
652 interactions (see, e.g., Scandura and Iammarino, 2022). Last but not least,  
653 the conditioned meme score might be used to help measuring information  
654 overload (IOL, Holyst et al. 2024) in science, which has lately become a

critical issue for researchers.

Further direction of research should also take into account the dynamics of memes contagiousness (i.e., equip the analysis with a dynamical point of view as it was done, e.g., by Kuhn et al. 2014), investigating whether the impact of affiliation on contagiousness – potentially representing a gap between Academia, Companies, and Big Tech – is expanding or contracting would provide valuable insights. The relationship between Big Tech and Academia in AI could also be compared to parallel relations in non-AI computer science and in other disciplines, such as medicine, where the role of “Big Pharma” companies is considered significant (Giunta et al., 2016).

A limitation of this study, inherent in the original meme score calculating technique, is the operationalization of memes as simple noun chunks, overlooking the potential synonymy of phrases. Consequently, there is a possibility that high-spreading memes may be overlooked. Possible solutions could be techniques such as clustering embedded memes.

In conclusion, our results suggest that the prevalent notion of Big Tech’s dominance over AI research is overly simplistic. We advocate for a more nuanced understanding of the roles played by Big Tech companies, corporations, and Academia in shaping the future of AI.

## Acknowledgements

The research was funded by (POB Cybersecurity and Data Science) of Warsaw University of Technology within the Excellence Initiative: Research University (IDUB) programme. This work was also funded by the European Union under the Horizon Europe grant OMINO – Overcoming Multilevel Information Overload (grant number 101086321, <http://ominoproject.eu>). Views and opinions expressed are those of the authors alone and do not necessarily reflect those of the European Union or the European Research Executive Agency. Neither the European Union nor the European Research Executive Agency can be held responsible for them. Computational part of this study was supported in part by the Poznań Supercomputing and Networking Center (grant number 607).

## References

Mohamed Abdalla and Moustafa Abdalla. The Grey Hoodie Project: Big Tobacco, Big Tech, and the Threat on Academic Integrity. In *Proceedings*

- 689 of the 2021 AAAI/ACM Conference on AI, Ethics, and Society, AIES '21,  
690 pages 287–297, New York, NY, USA, July 2021. Association for Computing  
691 Machinery. ISBN 978-1-4503-8473-5. doi:[10.1145/3461702.3462563](https://doi.org/10.1145/3461702.3462563).
- 692 Nur Ahmed and Muntasir Wahed. The De-democratization of AI: Deep  
693 Learning and the Compute Divide in Artificial Intelligence Research.  
694 *arXiv:2010.15581 [cs]*, October 2020. URL [http://arxiv.org/abs/2010.](http://arxiv.org/abs/2010.15581)  
695 [15581](http://arxiv.org/abs/2010.15581).
- 696 Tanya Araújo and Elsa Fontainha. Are scientific memes inherited dif-  
697 ferently from gendered authorship? *Scientometrics*, 117, 08 2018.  
698 doi:[10.1007/s11192-018-2903-7](https://doi.org/10.1007/s11192-018-2903-7).
- 699 Abeba Birhane, Pratyusha Kalluri, Dallas Card, William Agnew, Ravit  
700 Dotan, and Michelle Bao. The Values Encoded in Machine Learning Re-  
701 search. *arXiv:2106.15590 [cs]*, June 2021. URL [http://arxiv.org/abs/](http://arxiv.org/abs/2106.15590)  
702 [2106.15590](http://arxiv.org/abs/2106.15590).
- 703 Lutz Bornmann. Do altmetrics point to the broader impact of research? An  
704 overview of benefits and disadvantages of altmetrics. *Journal of Informet-*  
705 *rics*, 8(4):895–903, October 2014. ISSN 1751-1577. URL [https://www.](https://www.sciencedirect.com/science/article/pii/S1751157714000868)  
706 [sciencedirect.com/science/article/pii/S1751157714000868](https://www.sciencedirect.com/science/article/pii/S1751157714000868).
- 707 Sergey Brin and Lawrence Page. The anatomy of a large-scale hypertextual  
708 web search engine. *Computer Networks and ISDN Systems*, 30(1):107–117,  
709 1998. ISSN 0169-7552. doi:[https://doi.org/10.1016/S0169-7552\(98\)00110-](https://doi.org/10.1016/S0169-7552(98)00110-X)  
710 [X](https://doi.org/10.1016/S0169-7552(98)00110-X). Proceedings of the Seventh International World Wide Web Conference.
- 711 Dennis Chong and James N. Druckman. Framing Theory. *Annual Review of*  
712 *Political Science*, 10(1):103–126, June 2007. ISSN 1094-2939, 1545-1577.  
713 doi:[10.1146/annurev.polisci.10.072805.103054](https://doi.org/10.1146/annurev.polisci.10.072805.103054).
- 714 Richard Dawkins. *The selfish gene*. Oxford University Press, New York,  
715 1976. ISBN 978-0-19-857519-1.
- 716 Pietro della Briotta Parolo, Rainer Kujala, Kimmo Kaski, and Mikko Kivelä.  
717 Tracking the cumulative knowledge spreading in a comprehensive cita-  
718 tion network. *Physical Review Research*, 2(1):013181, February 2020.  
719 doi:[10.1103/PhysRevResearch.2.013181](https://doi.org/10.1103/PhysRevResearch.2.013181).



- 720 Olive Jean Dunn. Multiple comparisons using rank sums. *Technometrics*, 6  
721 (3):241–252, 1964. doi:[10.1080/00401706.1964.10490181](https://doi.org/10.1080/00401706.1964.10490181).
- 722 Christoph Ebell, Ricardo Baeza-Yates, Richard Benjamins, Hengjin Cai,  
723 Mark Coeckelbergh, Tania Duarte, Merve Hickok, Aurelie Jacquet, Angela  
724 Kim, Joris Krijger, John MacIntyre, Piyush B. Madhamshettiwar, Lauren  
725 Maffeo, Jeanna Matthews, Larry Medsker, Peter Smith, Savannah Jennifer  
726 Thais, Savannah Thais, Savannah Thais, and Savannah Thais. Towards  
727 intellectual freedom in an AI ethics global community. *AI and Ethics*, 1  
728 (2):131–138, 2021. doi:[10.1007/s43681-021-00052-5](https://doi.org/10.1007/s43681-021-00052-5).
- 729 The Economist. Move over faang, here comes maga – the tech giants are still  
730 in rude health. *The Economist*, August 2018. Retrieved 28 March, 2024.
- 731 James A. Evans. Industry induces academic science to know less about more.  
732 *American Journal of Sociology*, 116(2):389–452, 2010. doi:[10.1086/653834](https://doi.org/10.1086/653834).
- 733 Michael Färber and Lazaros Tampakis. Analyzing the impact of companies  
734 on AI research based on publications. *Scientometrics*, November 2023.  
735 ISSN 1588-2861. doi:[10.1007/s11192-023-04867-3](https://doi.org/10.1007/s11192-023-04867-3).
- 736 Roberto Martin N. Galang. Divergent diffusion: Understanding the interac-  
737 tion between institutions, firms, networks and knowledge in the interna-  
738 tional adoption of technology. *Journal of World Business*, 49(4):512–521,  
739 2014. ISSN 1090-9516. doi:<https://doi.org/10.1016/j.jwb.2013.12.005>.  
740 URL [https://www.sciencedirect.com/science/article/pii/](https://www.sciencedirect.com/science/article/pii/S109095161300093X)  
741 [S109095161300093X](https://www.sciencedirect.com/science/article/pii/S109095161300093X).
- 742 Anna Giunta, Filippo M. Pericoli, and Eleonora Pierucci. University–  
743 industry collaboration in the biopharmaceuticals: the italian case.  
744 *The Journal of Technology Transfer*, 41(4):818–840, Aug 2016.  
745 doi:[10.1007/s10961-015-9402-2](https://doi.org/10.1007/s10961-015-9402-2). URL [https://doi.org/10.1007/](https://doi.org/10.1007/s10961-015-9402-2)  
746 [s10961-015-9402-2](https://doi.org/10.1007/s10961-015-9402-2).
- 747 Maarten Grootendorst. Bertopic: Neural topic modeling with a class-based  
748 tf-idf procedure, 2022.
- 749 Thilo Hagendorff and Kristof Meding. Ethical considerations and statisti-  
750 cal analysis of industry involvement in machine learning research. *Ai &*  
751 *Society*, pages 1–11, 2021. doi:[10.1007/s00146-021-01284-z](https://doi.org/10.1007/s00146-021-01284-z).

- 752 Shaun P. Hargreaves Heap and Ashok Parikh. The diffusion  
753 of ideas in the academy: A quantitative illustration from eco-  
754 nomics. *Research Policy*, 34(10):1619–1632, 2005. ISSN 0048-7333.  
755 doi:<https://doi.org/10.1016/j.respol.2005.08.005>. URL <https://www.sciencedirect.com/science/article/pii/S0048733305001782>.  
756
- 757 Sture Holm. A simple sequentially rejective multiple test procedure. *Scandi-*  
758 *navian Journal of Statistics*, 6(2):65–70, 1979. URL <http://www.jstor.org/stable/4615733>.  
759
- 760 Janusz A. Hołyst, Philipp Mayr, Michael Thelwall, Ingo Frommholz, Shlomo  
761 Havlin, Alon Sela, Yoed N. Kenett, Denis Helic, Aljoša Rehar, Sebasti-  
762 jan R. Maček, Przemysław Kazienko, Tomasz Kajdanowicz, Przemysław  
763 Biecek, Bolesław K. Szymanski, and Julian Sienkiewicz. Protect our envi-  
764 ronment from information overload. *Nature Human Behaviour*, 8(3):402–  
765 403, Mar 2024. ISSN 2397-3374. doi:[10.1038/s41562-024-01833-8](https://doi.org/10.1038/s41562-024-01833-8). URL  
766 <https://doi.org/10.1038/s41562-024-01833-8>.
- 767 Roman Jurowetzki, Daniel Hain, Juan Mateos-Garcia, and Konstantinos  
768 Stathoulopoulos. The Privatization of AI Research(-ers): Causes and  
769 Potential Consequences – From university-industry interaction to pub-  
770 lic research brain-drain? *arXiv:2102.01648 [cs]*, February 2021. URL  
771 <http://arxiv.org/abs/2102.01648>.
- 772 Aharon Kantorovich. An evolutionary view of science: Imitation and memet-  
773 ics. *Social Science Information*, 53(3):363–373, 2014. ISSN 0539-0184.  
774 doi:[10.1177/0539018414526325](https://doi.org/10.1177/0539018414526325). Publisher: SAGE Publications Ltd.
- 775 David J. Ketchen and Christopher L. Shook. The application of cluster  
776 analysis in strategic management research: an analysis and critique.  
777 *Strategic Management Journal*, 17(6):441–458, 1996. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/%28SICI%291097-0266%28199606%2917%3A6%3C441%3A%3AAID-SMJ819%3E3.0.CO%3B2-G>.  
778  
779
- 780 Joel Klinger, Juan Mateos-Garcia, and Konstantinos Stathoulopoulos. Deep  
781 learning, deep change? Mapping the development of the Artificial Intelli-  
782 gence General Purpose Technology. *arXiv:1808.06355 [cs, econ]*, August  
783 2018. URL <http://arxiv.org/abs/1808.06355>.

- 784 Joel Klinger, Juan Mateos-Garcia, and Konstantinos Stathoulopoulos. A  
785 narrowing of AI research. *Social Science Research Network*, 2020.  
786 doi:[10.2139/ssrn.3698698](https://doi.org/10.2139/ssrn.3698698).
- 787 Bastian Krieger, Maikel Pellens, Knut Blind, Sonia Gruber, and Torben  
788 Schubert. Are firms withdrawing from basic research? An analysis of firm-  
789 level publication behaviour in Germany. *Scientometrics*, 126(12):9677–  
790 9698, December 2021. ISSN 1588-2861. doi:[10.1007/s11192-021-04147-y](https://doi.org/10.1007/s11192-021-04147-y).
- 791 William H. Kruskal and W. Allen Wallis. Use of ranks in one-criterion vari-  
792 ance analysis. *Journal of the American Statistical Association*, 47(260):  
793 583–621, 1952. doi:[10.1080/01621459.1952.10483441](https://doi.org/10.1080/01621459.1952.10483441).
- 794 Tobias Kuhn, Matjaž Perc, and Dirk Helbing. Inheritance Patterns in Cita-  
795 tion Networks Reveal Scientific Memes. *Physical Review X*, 4(4):041036,  
796 November 2014. doi:[10.1103/PhysRevX.4.041036](https://doi.org/10.1103/PhysRevX.4.041036).
- 797 Na Liu, Philip Shapira, and Xiaoxu Yue. Tracking developments in artifi-  
798 cial intelligence research: constructing and applying a new search strat-  
799 egy. *Scientometrics*, 126(4):3153–3192, April 2021. ISSN 1588-2861.  
800 doi:[10.1007/s11192-021-03868-4](https://doi.org/10.1007/s11192-021-03868-4).
- 801 Wenyuan Liu, Andrea Nanetti, and Siew Ann Cheong. Knowledge  
802 evolution in physics research: An analysis of bibliographic coupling  
803 networks. *PLOS ONE*, 12(9):e0184821, 2017. ISSN 1932-6203.  
804 doi:[10.1371/journal.pone.0184821](https://doi.org/10.1371/journal.pone.0184821).
- 805 Wenyuan Liu, Stanisław Saganowski, Przemysław Kazienko, and Siew Ann  
806 Cheong. Predicting the Evolution of Physics Research from a Complex  
807 Network Perspective. *Entropy*, 21(12):1152, December 2019. ISSN 1099-  
808 4300. doi:[10.3390/e21121152](https://doi.org/10.3390/e21121152).
- 809 Kyle Lo, Lucy Lu Wang, Mark Neumann, Rodney Michael Kinney, and  
810 Daniel S. Weld. S2orc: The semantic scholar open research corpus. In  
811 *Annual Meeting of the Association for Computational Linguistics*, 2020.  
812 URL <https://api.semanticscholar.org/CorpusID:215416146>.
- 813 Jin Mao, Zhentao Liang, Yujie Cao, and Gang Li. Quantifying cross-  
814 disciplinary knowledge flow from the perspective of content: Introducing  
815 an approach based on knowledge memes. *Journal of Informetrics*, 14(4):  
816 101092, 2020. ISSN 17511577. doi:[10.1016/j.joi.2020.101092](https://doi.org/10.1016/j.joi.2020.101092).

- 817 Chao Min, Qingyu Chen, Erjia Yan, Yi Bu, and Jianjun Sun. Citation  
818 cascade and the evolution of topic relevance. *Journal of the Association*  
819 *for Information Science and Technology*, 72(1):110–127, 2021. ISSN 2330-  
820 1643. doi:[10.1002/asi.24370](https://doi.org/10.1002/asi.24370).
- 821 Gabriel Popkin. How scientists can team up with big tech. *Nature*, 565  
822 (7741):665–667, January 2019. doi:[10.1038/d41586-019-00290-y](https://doi.org/10.1038/d41586-019-00290-y).
- 823 PwC. Global top 100 companies by market capitalisation 20176. Technical  
824 report, 2016. [Accessed 28-03-2024].
- 825 PwC. Global top 100 companies by market capitalisation 2017. Technical  
826 report, 2017. [Accessed 28-03-2024].
- 827 PwC. Global top 100 companies by market capitalisation 2018. Technical  
828 report, 2018. [Accessed 28-03-2024].
- 829 PwC. Global top 100 companies by market capitalisation 2019. Technical  
830 report, 2019. [Accessed 28-03-2024].
- 831 Alessandra Scandura and Simona Iammarino. Academic engagement with  
832 industry: the role of research quality and experience. *The Journal*  
833 *of Technology Transfer*, 47(4):1000–1036, August 2022. ISSN 1573-  
834 7047. doi:[10.1007/s10961-021-09867-0](https://doi.org/10.1007/s10961-021-09867-0). URL [https://doi.org/10.1007/](https://doi.org/10.1007/s10961-021-09867-0)  
835 [s10961-021-09867-0](https://doi.org/10.1007/s10961-021-09867-0).
- 836 Julian Sienkiewicz and Eduardo G. Altmann. Impact of lexical and sentiment  
837 factors on the popularity of scientific papers. *Royal Society Open Science*,  
838 3(6):160140, 2016. doi:[10.1098/rsos.160140](https://doi.org/10.1098/rsos.160140).
- 839 Xiaoling Sun and Kun Ding. Identifying and tracking scientific and tech-  
840 nological knowledge memes from citation networks of publications and  
841 patents. *Scientometrics*, 116(3):1735–1748, 2018. ISSN 0138-9130, 1588-  
842 2861. doi:[10.1007/s11192-018-2836-1](https://doi.org/10.1007/s11192-018-2836-1).
- 843 Caroline S. Wagner, J. David Roessner, Kamau Bobb, Julie Thomp-  
844 son Klein, Kevin W. Boyack, Joann Keyton, Ismael Rafols, and  
845 Katy Börner. Approaches to understanding and measuring in-  
846 terdisciplinary scientific research (idr): A review of the litera-  
847 ture. *Journal of Informetrics*, 5(1):14–26, 2011. ISSN 1751-  
848 1577. doi:<https://doi.org/10.1016/j.joi.2010.06.004>. URL [https://www.](https://www.sciencedirect.com/science/article/pii/S1751157710000581)  
849 [sciencedirect.com/science/article/pii/S1751157710000581](https://www.sciencedirect.com/science/article/pii/S1751157710000581).

- 850 Meredith Whittaker. The steep cost of capture. *Interactions*, 28(6):50–55,  
851 November 2021. ISSN 1072-5520, 1558-3449. doi:[10.1145/3488666](https://doi.org/10.1145/3488666).
- 852 Frank Wilcoxon. Individual comparisons by ranking methods. *Biometrics*  
853 *Bulletin*, 1(6):80–83, 1945. doi:<https://doi.org/10.2307/3001968>.
- 854 Chaojiang Wu, Chelsey Hill, and Erjia Yan. Disciplinary knowledge diffusion  
855 in business research. *Journal of Informetrics*, 11(2):655–668, 2017. ISSN  
856 1751-1577. doi:<https://doi.org/10.1016/j.joi.2017.04.005>. URL [https://](https://www.sciencedirect.com/science/article/pii/S1751157716303455)  
857 [www.sciencedirect.com/science/article/pii/S1751157716303455](https://www.sciencedirect.com/science/article/pii/S1751157716303455).
- 858 Jian Xu, Yi Bu, Ying Ding, Sinan Yang, Hongli Zhang, Chen Yu, and Lin  
859 Sun. Understanding the formation of interdisciplinary research from the  
860 perspective of keyword evolution: a case study on joint attention. *Sciento-*  
861 *metrics*, 117(2):973–995, Nov 2018. doi:[10.1007/s11192-018-2897-1](https://doi.org/10.1007/s11192-018-2897-1). URL  
862 <https://doi.org/10.1007/s11192-018-2897-1>.
- 863 Meg Young, Michael Katell, and P.M. Krafft. Confronting power and cor-  
864 porate capture at the FAccT conference. In *2022 ACM Conference on*  
865 *Fairness, Accountability, and Transparency*, FAccT ’22, pages 1375–1386.  
866 Association for Computing Machinery, 2022. ISBN 978-1-4503-9352-2.  
867 doi:[10.1145/3531146.3533194](https://doi.org/10.1145/3531146.3533194).
- 868 Xi Zhang, Xianhai Wang, Hongke Zhao, Patricia Ordóñez de Pablos,  
869 Yongqiang Sun, and Hui Xiong. An effectiveness analysis of altmetrics  
870 indices for different levels of artificial intelligence publications. *Scientomet-*  
871 *rics*, 119(3):1311–1344, June 2019. ISSN 1588-2861. doi:[10.1007/s11192-](https://doi.org/10.1007/s11192-019-03088-x)  
872 [019-03088-x](https://doi.org/10.1007/s11192-019-03088-x).

## 873 Appendix A. Data gathering and processing

874 Following Liu et al. (2021) we employed simple search strategy for AI  
875 papers. We searched for following words from titles and abstracts: *artificial*  
876 *intelligence, machine learning, classifier, neural network, deep learning, data*  
877 *science, nlp, machine-learning, computer vision.*

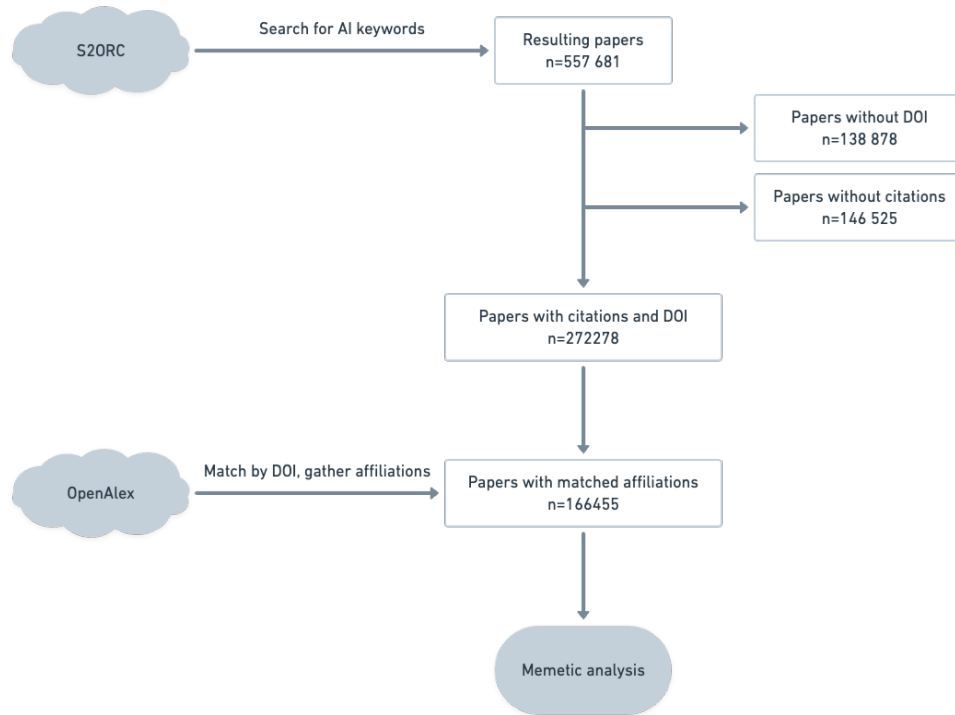


Figure A.10: Diagram of pipeline used for gathering and processing papers.

878 **Appendix B. OpenAlex categories of affiliations treated as Academia**

Category	Most important institutions
education	Tsinghua University
	Shanghai Jiao Tong University
	Carnegie Mellon University
	Stanford University
	Massachusetts Institute of Technology
government	Chinese Academy of Sciences
	French National Centre for Scientific Research
	National Institute of Health
	French Institute for Research in Computer Science and Automation
	Commonwealth Scientific and Industrial Research Organisation
nonprofit	Max Planck Society
	Electric Power Research Institute
	German Research Centre for Artificial Intelligence
	Instituto de Salud Carlos III
	SRI International
facility	Ecole Polytechnique Federale de Lausanne
	National Institute of Informatics
	Electronics and Telecommunications Research Institute
	United States Air Force Research Laboratory
	Italian Institute of Technology
healthcare	Mayo Clinic
	Boston Children’s Hospital
	Brigham and Women’s Hospital
	The University of Texas Southwestern Medical Centre
	Vanderbilt University Medical Centre

Table B.2: Five institutions with the biggest number of publications from each OpenAlex category included by us in the Academia category. Due to overlap in education, government and facility categories we decided to merge them into Academia.

879 **Appendix C. Annotations**

880     The annotation process was conducted for the top memes, defined as the  
881     ones that have a meme score value over 0.25, and more than 20 occurrences,



882 by the team of domain experts. The pipeline for the creation of the annota-  
883 tions consisted of several steps, described below.

- 884 1. Extraction of basic parameters: meme score, meme occurrences, and  
885 meme phrases.
- 886 2. Providing a short description of more complex, and not obvious memes.  
887 The research was focused on providing the descriptions connected to  
888 ML, as we know, that used dataset comes from this domain. Providing  
889 the titles of ML papers found connected to researched memes.
- 890 3. Providing the flags for duplicate or very similar memes. Most of them  
891 were the abbreviations and their elaborations.
- 892 4. Creation of first frequently repeating topics (ML algorithms, Medical  
893 terms, Security, Niche ML application, ML concept, and Other), based  
894 on the annotation process.
- 895 5. Iterative extraction of another topic from the Other group (Maths,  
896 Graphs, Computer Vision, NLP, Data Related).
- 897 6. Revising whole annotations and re-assignment to proper topics.