

# From Insights to Actions: The Impact of Interpretability and Analysis Research on NLP

Anonymous ACL submission

## Abstract

Interpretability and analysis (IA) research is a growing subfield within NLP with the goal of developing a deeper understanding of the behavior or inner workings of NLP systems and methods. Despite growing interest in the subfield, a commonly voiced criticism is that it lacks actionable insights and therefore has little impact on NLP. In this paper, we seek to quantify the impact of IA research on the broader field of NLP. We approach this with a mixed-methods analysis of: (1) a citation graph of 185K+ papers built from all papers published at ACL and EMNLP conferences from 2018 to 2023, and (2) a survey of 138 members of the NLP community. Our quantitative results show that IA work is well-cited outside of IA, and central in the NLP citation graph. Through qualitative analysis of survey responses and manual annotation of 556 papers, we find that NLP researchers build on findings from IA work and perceive it is important for progress in NLP, multiple subfields, and rely on its findings and terminology for their own work. Many novel methods are proposed based on IA findings and highly influenced by them, but highly influential non-IA work cites IA findings without being driven by them. We end by summarizing what is missing in IA work today and provide a call to action, to pave the way for a more impactful future of IA research.

## 1 Introduction

The rapid progress made in the development of large language models (LLMs, Devlin et al. (2019); Radford et al. (2019); Raffel et al. (2020); Bommasani et al. (2022); Touvron et al. (2023); OpenAI et al. (2024); Team et al. (2024)) has had a profound impact on the field of natural language processing (NLP) (Gururaja et al., 2023). While these models demonstrate unprecedented performance and novel capabilities (Brown et al., 2020; Wei et al., 2022), and are rapidly finding their way

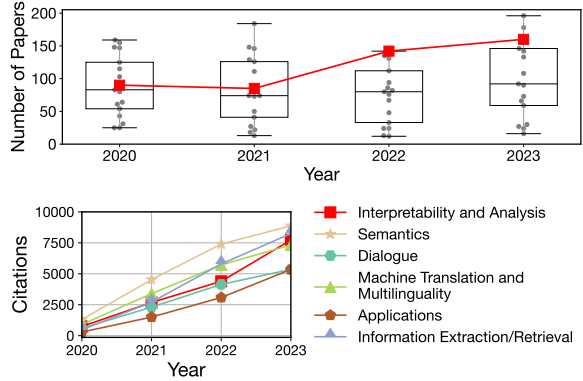


Figure 1: Interpretability and analysis (IA) is an increasingly popular subfield of NLP: (top) Number of IA papers in ACL/EMNLP in comparison to other tracks that have existed since 2020. The number of IA papers has grown considerably, from 90 papers in 2020 to 160 papers in 2023 (a growth rate of 77.8%). This is the highest growth rate among these tracks. (bottom) Citations to IA papers compared to other highly cited tracks.

into real-world applications (OpenAI, 2022; Microsoft, 2023; Google, 2024), they are also opaque and largely treated as black boxes, which does not satisfy other expectations for successful machine learning deployment, such as fairness, trust, accountability, and explainability (Lipton, 2018; Goodman and Flaxman, 2017).

In NLP research, these factors have motivated a large body of work on *interpretability and analysis* (IA), which aims to understand the inner workings of LLMs and explain their predictions (Belinkov and Glass, 2019; Rogers et al., 2020; Rauker et al., 2023, *inter alia*). Researchers in this area are often motivated by the idea that better understanding LLMs is imperative to improve their efficiency, robustness, and trustworthiness, towards successful and safe deployment. IA research has thus witnessed rapid growth in the past few years and is now one of the biggest research areas (in terms of number of publications and citations) at the major NLP conferences (see Figure 1).

063 Despite the rapid growth of IA research (see also  
064 Figure 9), a commonly voiced criticism is that it  
065 often lacks actionable insights, especially for how  
066 to improve models, and therefore has little impact  
067 on how new NLP models are designed and built.  
068 This criticism raises questions about the usefulness  
069 of IA research, and whether its current form is the  
070 right path towards progress in NLP.

071 In this work, we tackle these questions with a  
072 systematic, mixed-methods study of the impact  
073 of IA research on NLP in the past and the present,  
074 and use our findings to inform a vision for the  
075 future of IA. More specifically, we ask: **how does  
076 interpretability and analysis research influence  
077 NLP researchers in what they choose to work  
078 on, what they cite, and how they think about  
079 NLP altogether?**

080 We perform a bibliometric analysis of 185,384  
081 publications based on the two major NLP confer-  
082 ences, ACL and EMNLP, between 2018 and 2023,  
083 and solicit opinions from 138 members of the NLP  
084 community via a survey. In addition to quantitative  
085 results, we perform qualitative analysis of survey  
086 responses and 556 papers. This approach gives us a  
087 holistic view of the impact of IA research on NLP.

088 Our analysis reveals that (1) NLP researchers  
089 build on findings from IA work in their research,  
090 regardless of whether they work on IA themselves  
091 or not (§4), (2) NLP researchers and practitioners  
092 perceive IA work to be important for progress in  
093 NLP, multiple subfields, and their own work, for  
094 various reasons (§5), and (3) many novel non-IA  
095 methods are proposed based on IA findings and  
096 highly influenced by them, for various areas, even  
097 though highly influential non-IA work is not driven  
098 by IA findings despite citing them (§6).

099 While our findings show that IA work presents  
100 insightful observations, there are still opportuni-  
101 ties for greater impact on the rest of NLP. Thus,  
102 based on survey responses, we identify the key in-  
103 gredients that are missing in IA research today —  
104 unification; actionable recommendations; human-  
105 centered, interdisciplinary work; and standardized,  
106 robust methods — and close with a call to action  
107 with recommendations (§7). We hope our work  
108 paves the way towards a more impactful future for  
109 IA research as the field continues to grow.

## 110 2 Methodology

111 We start by discussing what we consider as IA  
112 research and our approach for measuring impact.

## 2.1 Interpretability and analysis (IA) research 113

114 *Interpretability* research has a long tradition in Ma-  
115 chine Learning as well adjacent fields like NLP  
116 (Tishby and Zaslavsky, 2015; Karpathy et al., 2015;  
117 Kim et al., 2018, *inter alia*). There is no single  
118 agreed upon definition of the term *interpretability*  
119 (see Lipton (2018) for a critical discussion), but two  
120 prominent types of interpretability research focus  
121 on post-hoc explainability or increasing the trans-  
122 parency of machine learning methods and models  
123 (Lipton, 2018; Madsen et al., 2024). *Analysis* re-  
124 search is an even broader term and one might argue  
125 that nearly every scientific paper contains some  
126 form of analysis. In NLP, however, many inter-  
127 pretability and analysis papers have in common that  
128 their *primary* contribution is an analysis that aims  
129 to advance our understanding of NLP in some way,  
130 e.g., by analyzing methods, models, or algorithms  
131 (Belinkov and Glass, 2019; Rogers et al., 2020).

132 Here, we adopt a broad definition of interpretabil-  
133 ity and analysis (IA) research in NLP that includes  
134 all papers that aim to **develop a deeper under-  
135 standing** of the behavior or inner workings of  
136 NLP models, methods, or systems. This includes  
137 work on explaining models’ predictions or inter-  
138 nal computations, investigating broader phenom-  
139 ena observed during pre-training or adaptation, and  
140 providing a better understanding of the limitations  
141 and robustness of existing models.

## 2.2 Measuring impact 142

143 Our goal is to measure the *impact* of IA work on  
144 NLP research, which is not trivial to define, let  
145 alone quantify. To get a **holistic view of impact**,  
146 we consider two different, complementary ways of  
147 measuring impact – a bibliometric analysis, and a  
148 survey of the NLP community.

**Citational impact** In scientometrics research, ci-  
149 tation counts are used as a standard measure of  
150 scientific impact (Nicolaisen, 2007; Bornmann and  
151 Daniel, 2008; Chacon et al., 2020, *inter alia*). Thus,  
152 we perform a bibliometric analysis to quantify the  
153 citational impact of IA work on NLP research.<sup>1</sup> We  
154 note that citation behavior is complex and there is a  
155 growing consensus that citation statistics might not  
156 be sufficient for measuring impact (Bornmann and  
157 Daniel, 2008; Zhu et al., 2015; Iqbal et al., 2021).  
158

<sup>1</sup>This choice excludes other forms of impact such as in-  
creasing user trust, influencing policy and regulation, etc. In  
addition, even though IA work impacts other fields, this is  
beyond the scope of our paper.

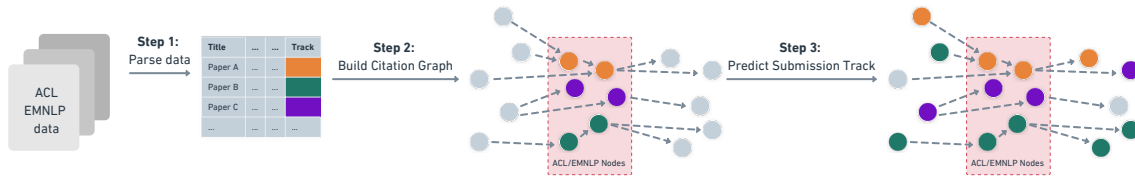


Figure 2: Diagram showing the process of constructing our citation graph. Starting from an initial set of ACL and EMNLP papers we collect citations via the Semantic Scholar API and label papers with a classifier.

**Surveying the NLP community** To incorporate a second dimension of impact beyond citation counts, we survey NLP researchers and practitioners on how they view the impact of IA research on the field. Specifically, we ask respondents about their *perceptions* of IA (its importance in general, for specific subfields, and its impact on progress in NLP), and their *use* of IA (how much they read, are influenced by, and use concepts from IA work). We also solicit opinions on what is missing in IA research and where it should go in the future.

### 3 Citation graph and community survey

Here, we describe the construction of our citation graph for bibliometric analysis, and the design of our survey of the community.

#### 3.1 Citation graph construction

Figure 2 illustrates the process of constructing our citation graph. We start from an initial set of all papers published at ACL and EMNLP from 2018 to 2023. We focus on these two venues as they are leading NLP conferences with a dedicated track for interpretability and analysis research since 2020.<sup>2</sup> Using this initial set of papers, we build a citation graph using the Semantic Scholar API (Kinney et al., 2023). For papers outside our initial set, where we have gold labels, we rely on classifiers to predict submission tracks. More details on all these stages are provided below.

**Collecting ACL and EMNLP papers** We collect paper lists and track information from various sources (see Table 3 in Appendix B), as there is no one source of this data for ACL and EMNLP conferences.<sup>3</sup> Between 2018 and 2023, official names of submission tracks have changed substantially, so we standardize all data to 27 tracks. More details on this process are provided in Appendix B, including summary statistics per track (Table 1).

<sup>2</sup>We discuss this decision in more detail in Section 8.

<sup>3</sup>The ACL Anthology does not contain information on the submission track.

**Building the citation graph** We collect the citations of each paper in our initial set via the Semantic Scholar API (Kinney et al., 2023), resulting in a citation graph of 185,384 papers (see Table 2 in Appendix B for additional statistics). For each node (paper) in the graph, we store its title, abstract, and venue. For each edge (citation), we store information on the citation intent (binary labels for background, use of methods or comparing results), and citation influence (normal vs. highly influential), all of which are provided by Semantic Scholar.

**Labeling the citation graph** To assign all papers in the citation graph to our standardized set of tracks, we train a classifier based on the titles and abstracts from our initial set of papers. We find that some tracks are very hard to predict due to limited training data and the inherent ambiguity of submission tracks. We thus keep 11 well-performing labels (including IA), and introduce an ‘Other’ label to group the remaining papers. More details on classifier construction are provided in Appendix B.

Our final classifier achieves a test micro/macro-F1 score of 0.61/0.61. Although this performance appears rather low, we note that submission tracks have fuzzy boundaries, so papers can often be plausible submissions to multiple tracks. Given that we care primarily about accurately predicting IA compared to other tracks, we evaluate our classifier on two additional gold sets of data (see Appendix B.1) and obtain 78.1% and 87.8% accuracy on each set.

#### 3.2 Surveying the NLP community

To solicit opinions from the NLP community on the impact of IA research, we ran a survey from March 19th to June 7th, 2024, advertising within our networks, on social media, and on NLP mailing lists. The full survey is shown in Appendix C.

To strike a balance between easy scoring and respondent expressivity, we included multiple-choice as well as optional free response questions (Shaughnessy et al., 2015). We refined the survey following

	Phonology, Morphology and Word Segmentation	Syntax	Discourse and Pragmatics	Linguistic Theories and Psycholinguistics	Industry	Social Science	Applications	Sentiment Analysis	Semantics	Machine Translation and Multilinguality	Efficient Methods	Information Extraction/Retrieval	Dialogue	Generation	Multimodality, Speech and Grounding	Resources and Evaluation	Machine Learning	Question Answering	Summarization	Ethics	Large Language Models
2020	83.2%	78.5%	73.5%	81.2%	N/A	68.0%	66.5%	56.2%	60.6%	58.8%	N/A	57.2%	56.1%	55.9%	58.0%	56.6%	59.2%	54.0%	46.9%	N/A	N/A
2021	N/A	69.7%	78.2%	N/A	N/A	59.4%	60.4%	49.8%	52.8%	54.1%	41.9%	47.3%	50.5%	52.1%	47.1%	45.7%	49.9%	46.5%	47.4%	35.3%	N/A
2022	N/A	71.8%	N/A	67.4%	N/A	64.5%	56.5%	57.8%	52.1%	57.2%	66.7%	56.6%	53.0%	56.3%	49.2%	49.4%	43.9%	49.9%	46.5%	40.7%	N/A
2023	N/A	73.9%	64.6%	63.4%	70.0%	57.1%	59.1%	70.6%	62.0%	57.1%	60.5%	58.4%	59.2%	51.0%	56.2%	53.2%	51.4%	50.2%	50.4%	45.8%	29.1%

Figure 3: Interpretability and analysis track CSI scores matrix against other tracks. These represent the probability that a random interpretability and analysis paper published in certain year has more citations than a random paper of other track published the same year.

best practices<sup>4</sup> and with feedback from four senior NLP researchers who filled out a pilot version. We received a total of 138 responses from NLP researchers in academia and practitioners in industry, with 61% of respondents not working on IA themselves (see Appendix C for more statistics).

Two authors performed qualitative coding, an inductive method from the social sciences (Saldana, 2021), to identify themes in answers to the free-response questions. More details on the coding process are provided in Appendix D. We measure inter-coder reliability with percentage agreement (O’Connor and Joffe, 2020), which was above 90% across all subsets of annotation.

#### 4 Researchers build on findings from IA research in their work

We begin by analyzing whether researchers use contributions of IA research in their work. We approach this by analyzing citational use, as well as survey-reported use beyond citations.

##### IA papers are cited more often than other tracks

When comparing papers from different tracks, global counts of citations can be misleading, as a small number of papers can account for most of the citations in a field (Ioannidis et al., 2016). To account for this, we compare citations based on the *Citation Success Index* (CSI; Milojević et al., 2017) metric. Given two groups of papers *A* and *B*, the CSI score computes the probability that a random paper from *A* is more cited than a random paper of *B*. This score is not subject to biases from the skewness of the citation distribution, and it is clearly interpretable; e.g., if we draw random IA and Machine Translation papers from EMNLP or ACL in

<sup>4</sup>We made sure to clarify definitions, avoid leading questions, etc. (Shaughnessy et al., 2015).

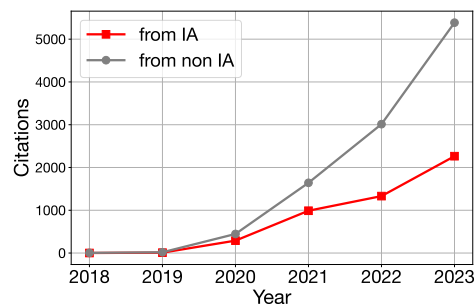


Figure 4: Origin of citations to IA papers.

2023, there is a 57.1% chance that the IA paper is more cited than the Machine Translation paper.

Figure 3 shows that CSI scores for the IA track are often favorable (CSI > 50%) when compared to other tracks. In 2023, only the Ethics and the Large Language Models tracks had favorable CSI scores against IA. This shows that IA papers have higher citational impact than other tracks, particularly in recent conferences.

##### IA papers are well cited outside of IA

While high CSI scores tell us that IA papers are cited well, they do not tell us where these citations are coming from, i.e., are IA papers mostly cited by other IA papers or by papers outside of IA? To evaluate the impact of IA work outside of IA, we compare citations within the same track, which we call *intra-track* citations, to *extra-track* citations, i.e., citations from outside the track.

Figure 4 shows that most citations to IA papers are predicted to be extra-track citations. The proportion of references to IA papers differs considerably by citing track, with papers about Efficient Methods, Machine Learning, and Large Language Models citing IA research more frequently than others (see Figure 11 for a visualization of this trend).

295 While the IA track does not stand out in terms of  
 296 its *extra-track* citations compared to other tracks  
 297 (see Figure 12), these results still demonstrate that  
 298 the citational impact of IA research extends well  
 299 beyond the IA track itself.

300 **IA papers are central in NLP** Next, we assess  
 301 whether IA papers are impacting NLP as a whole  
 302 rather than just specific tracks. We quantify this  
 303 with the *Betweenness Centrality* (BC) metric, a  
 304 measure of *interdisciplinarity* (Leydesdorff, 2007;  
 305 Barnett et al., 2011; Leydesdorff et al., 2018). BC  
 306 quantifies the extent to which a node in the graph  
 307 acts as a *bridge* along the shortest path between two  
 308 other nodes (Golbeck, 2015); nodes with higher BC  
 309 are considered more important as more information  
 310 passes through them.<sup>5</sup> Therefore, we interpret pa-  
 311 pers with a high BC as *important* papers that are es-  
 312 sential for the connectivity of the citation network.

313 We compute the BC for every paper in EMNLP  
 314 and ACL since the IA track started (2020), and find  
 315 that the median BC of IA papers is higher than most  
 316 other tracks, at  $1.23 \times 10^{-7}$ . Notably, IA ranks as  
 317 the second most central track overall, following the  
 318 Large Language Models track, which has a median  
 319 BC of  $1.95 \times 10^{-7}$ . These results (shown in Fig-  
 320 ure 10) provide further evidence that IA work plays  
 321 a central role in the ACL/EMNLP citation network.

322 **IA influences the work of NLP researchers** For  
 323 a complementary view of impact beyond citations,  
 324 we survey NLP community members on how often  
 325 they use concepts from IA in their day-to-day work,  
 326 and more broadly, how IA influences their work.

327 As Figure 5 shows, the median rating for use of  
 328 IA concepts by respondents who work on IA is *of-*  
 329 *ten*, while even the median respondent who doesn't  
 330 work on IA uses concepts from IA *sometimes*. In  
 331 both groups of respondents, there are people who  
 332 *always* use IA concepts in their day-to-day work.  
 333 Beyond this, IA work influences respondents in dif-  
 334 ferent ways: it provides respondents with research  
 335 ideas (91% of respondents who work on IA; 60% of  
 336 respondents who don't), changes mental models of  
 337 model capabilities and limitations (77%; 65%), and  
 338 helps ground explanations of respondents' results  
 339 (64%; 59%). Notably, only 9 (6.5%) respondents  
 340 state that IA does not affect their work. These re-  
 341 sults complement our citation-based findings by  
 342 providing further evidence that IA work impacts  
 343 both IA and non-IA researchers and their research.

<sup>5</sup>We provide further discussion of BC in Appendix B.1.

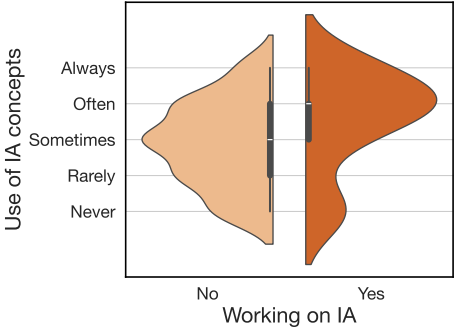


Figure 5: Survey responses on the frequency of using concepts from IA research, split by whether the respondents work in this field or not. Higher values indicate more frequent usage.

344 **5 Researchers find IA work important**

345 We continue by surveying the *perceived importance*  
 346 of IA work by the NLP community. We consider  
 347 various perspectives, such as the perceived impor-  
 348 tance of IA research on overall progress in NLP  
 349 as well as on individual subfields. 133 out of 138  
 350 respondents consider **IA work important**, and per-  
 351 ceive it as important **for progress in NLP, multiple**  
 352 **subfields, and for various reasons.**

353 **Perceived importance for progress in NLP** Fig-  
 354 ure 6 shows that most respondents agree that with-  
 355 out IA findings, progress in NLP in the last 5 years  
 356 (2019 to 2024) would have been slower, but not im-  
 357 possible. Surprisingly, it appears that *people who*  
 358 *are more deeply engaged with interpretability are*  
 359 *more critical of it*. Respondents who read more IA  
 360 work than other topics in NLP, respondents who of-  
 361 ten or always use concepts from IA literature, and  
 362 respondents who work on IA themselves all rate IA  
 363 as having a lower impact on progress in NLP than  
 364 those who read less IA, use related concepts less  
 365 frequently, and who work on other topics.

366 It is plausible that respondents who are more  
 367 engaged with IA work know it better and thus  
 368 give better-calibrated impressions of the field as a  
 369 whole, which happen to be more critical. However,  
 370 it is worth noting that they are perhaps forming  
 371 their opinions from a different sample of papers  
 372 (i.e., the average paper from a large body of work)  
 373 than those who are less engaged with IA work,  
 374 whose reading might be skewed towards IA work  
 375 that is more highly cited and influential. This also  
 376 raises the question of how IA or indeed any sub-  
 377 field *should* be evaluated – by the average paper in  
 378 it, or by the ones that stand out?

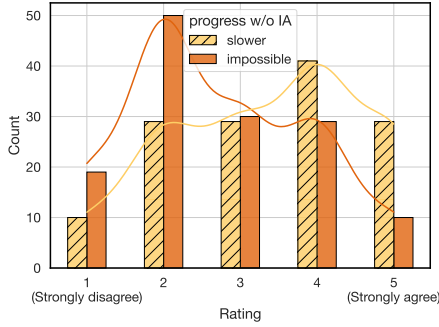


Figure 6: Survey responses (N=138) on whether progress in NLP in the last 5 years would have been *slower* or *impossible* without findings from interpretability and analysis research.

There are many other factors that could also influence the results we see, e.g., that respondents in different categories are reading IA papers that deal with different topics, that they have different levels of research experience, and that they have different definitions of “progress” in NLP. See §8 for a discussion of these factors.

### Perceived importance for different subfields

Figure 7 shows that the IA work is perceived as being important to differing extents for other subfields within NLP. The modal response is that IA work is *somewhat* important for work on multilinguality (52% of responses), multimodal learning (47%) and engineering for large language models (47%), and that it is *very* important for work on reasoning (63%) and bias (72%). Of the five subfields we consider, engineering for LLMs is perceived to be least impacted by IA work, with 31% of respondents indicating that they think IA work is not important for it. These findings are consistent with the themes we find in papers that are highly influenced by IA research, where bias and reasoning are well-represented, and pre-training and architectural advancements appear less frequently.

**Reasons for importance** When asked whether they thought IA work was important and if so, why, respondents overwhelmingly (133/138) consider it important, citing a variety of reasons, the most popular of which were: understanding model limitations and capabilities (90% of respondents), explainability for users (66%), improving model trustworthiness (59%), and improving model capabilities (50%). While a small percentage (4.3%) of respondents indicated that they thought it was not important (possibly also due to selection bias

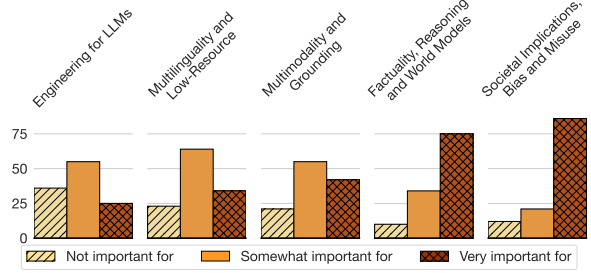


Figure 7: Survey responses (N=138) on how important interpretability and analysis research is to work in different subfields.

in our survey), we found that they voice the same concerns as those who do find it important, e.g., a lack of actionability, results that don’t scale, and a lack of impact on the most capable models of today. In our recommendations for the future of the field (§7), we go into these in more detail.

## 6 A closer look at influential papers

So far we have discussed findings about IA as a whole, either by considering the role of IA papers in the ACL/EMNLP citation graph or the perception of IA work within the community. In this section, we zoom in on specific influential papers sourced from both our survey and citation graph. We seek to answer: What are these papers about? What kind of work are they impacting, and how?

To this end, we inductively obtain the themes of a total of 585 papers, through qualitative coding of their titles and abstracts by two authors (Saldana, 2021). The 585 papers include: (1) All papers mentioned more than once as having influenced survey respondents’ work (N=29); (2) highly-cited IA papers from our citation graph (N=50); (3) highly-cited non-IA papers from our citation graph (N=50); (4) non-IA papers that cite and are highly influenced by the top-10 most-cited IA papers (N=456). The resulting themes are mostly descriptive, including topics (e.g., *in-context learning*, *training dynamics*) and contribution types (e.g., *novel method*, *analysis*). Percentage agreement on our coded themes is above 90% for each subset of papers. See Appendix D for more details.

Our analysis reveals that beyond background citations, IA work influences the development of many novel models and metrics outside of IA work, and affects work in domains such as question answering (QA), reasoning, and bias.

**What are influential IA papers about?** Of the papers that survey respondents submitted as examples of work that has directly influenced their own work, representation analysis appears in over a third of the papers, novel methods for interpretability (e.g., causality, interventions, steering, neuron/activation analysis, etc.) are proposed in nearly a quarter of them, and probing also appears in 24% of these papers.

In contrast, the top-50 most cited IA papers are more often about the *analysis* component of IA (40%). Novel methods (for analysis, evaluation, linguistics, probing) are proposed in 26% of papers, and evaluation is a main contribution of 32%. As expected, the most cited non-IA papers in our citation graph mostly consist of highly influential datasets, models, and methods, e.g., HotpotQA, BART, prefix-tuning (Yang et al., 2018; Lewis et al., 2020; Li and Liang, 2021). More top themes are shown with the percentage of papers in Table 5 in Appendix E.

We also find evidence that many IA papers create novel metaphors to understand models — e.g., seeing feed-forward layers as key-value memories (Geva et al., 2021), or reading from and writing to the “residual stream” (Elhage et al., 2021), and many analysis papers highlight the limits of models. As survey respondents cited these very reasons for why they perceive IA work as important, these themes corroborate why these papers would be particularly influential. In addition, many of the qualities that survey respondents feel are currently lacking in IA research (see §7) appear in these papers, such as moving beyond toy models (Wang et al., 2023), and providing actionable methods (Meng et al., 2022).

**Why are influential IA papers cited?** As citations can have a variety of reasons (Zhu et al., 2015; Tahamtan and Bornmann, 2019), we examine three types of citational intent — background, methods and results citations (see Figure 13 in Appendix E). Overall, we find that influential IA papers are cited most often as background citations, then as methods citations, and least frequently when comparing results. In comparison, highly cited papers that are *not* about IA tend to be cited most frequently for methods. This is expected, as many of these papers are about popular datasets and models, as described above.

**What are the citing papers about?** Despite the large number of background citations, however,

there is plenty of work—including non-IA work—that is highly influenced (according to Semantic Scholar) by IA research. For a closer look at what these citing papers do, we analyze all 456 papers with a highly influential citations to one of the top 10 most-cited IA papers, and annotate their themes based on titles and abstracts.

Unsurprisingly, many of the papers have themes in common with what they cite, e.g., papers that analyze multilingual models are frequently cited by papers on cross-lingual transfer. We thus focus on the *difference* in themes between citing papers and cited papers, and find that **over 33% of non-IA papers that are highly influenced by IA work propose novel methods**, e.g., many novel ICL methods cite analysis work on demonstrations (Min et al., 2022) and similarly, many novel methods for bias mitigation cite datasets for stereotype evaluation such as Nangia et al. (2020) and Nadeem et al. (2021). These provide concrete counterexamples to the claim that IA work does not influence modeling improvements.

### Is IA work impacting highly cited non-IA work?

Looking at the highly-cited non-IA papers, we find that these too tend to cite IA work frequently. 22 out of the top 50 most cited non-IA papers are even highly influenced by some IA work, but 28 are not highly influenced by *any* IA work. These results show that while highly influential non-IA work does acknowledge IA findings, it is likely not driven by them.

## 7 Main takeaways and discussion

We end by discussing our main findings and recommendations on how to move IA research forward.

**Main takeaways** In §4, we saw that *IA research plays a central role in NLP* and researchers build on findings from IA work in their research, regardless of whether they work on IA themselves or not. In section §5, we saw that *NLP researchers and practitioners perceive IA work to be important* for progress in NLP, and multiple subfields. They also find it important for their own work for a variety of reasons, regardless of whether they work on IA themselves. Finally, we took a closer look at the most influential IA papers in §6 and found that *many novel methods are proposed based on IA findings and highly influenced by them*, for various areas, in particular, work on reasoning, factual knowledge, and bias. All these findings present a

550 very positive view of IA research and its role within  
551 NLP in the past and the present. In the remainder  
552 of this section, we turn to the future of IA research.

553 **What is missing?** To understand what the NLP  
554 community believes to be important for the future  
555 of IA work, we asked survey respondents what they  
556 feel is missing in current IA work and what should  
557 be different going forward. 25% of the responses  
558 to this question mentioned a lack of big picture and  
559 unified understanding in IA work. For example,  
560 one respondent said:

561 *“I think the focus should be on climb-*  
562 *ing the right hill towards a higher level*  
563 *understanding instead of focusing on in-*  
564 *teresting individual behaviors.”*

565 The next three most frequent concerns are a lack of  
566 utility (i.e., not being useful in practice), modeling  
567 improvements and actionability—concerns that are  
568 also echoed by the respondents who do not find IA  
569 research useful for their own work. Interestingly, a  
570 commonly voiced opinion among these participants  
571 is that they believe that scale and performance are  
572 all that is needed for good NLP models, and that  
573 IA work only has importance for understanding  
574 models rather than for building them. Addition-  
575 ally, respondents mention that IA work could use  
576 more interdisciplinary connections, through col-  
577 laboration with domain experts, user studies, and  
578 human-centered approaches to computing.

579 Finally, we note another theme appearing in 10%  
580 of responses: as IA has a lack of consensus on  
581 reliable and trustworthy methods, it is unclear how  
582 such work should be evaluated. Although this is  
583 not a new concern (Belinkov and Glass, 2019), it  
584 remains relevant for the impact of IA on NLP.

585 **A call for action** Based on our findings, we make  
586 the following recommendations:

587 **Going forward, IA researchers should:**

- 588 1. Think more about the big picture
- 589 2. Strive for more actionable work
- 590 3. Center humans in your work
- 591 4. Work towards standardized, robust methods

Big-picture thinking involves working towards  
general truths about model architectures or behav-  
iors, rather than model-specific results. Actionable  
work requires thinking about how an IA finding can

propel new ways of building/using NLP systems,  
rather than being merely descriptive. Centering  
humans entails evaluation with realistic and rele-  
vant data and tasks, and performing user studies  
and human evaluation. Human-centered IA work  
can also be enhanced through interdisciplinary  
reading and collaboration. Finally, we urgently  
need to build consensus on using and evaluating IA  
methods. Rigorous, well-motivated methods (e.g.,  
using causality) are critical, rather than correlative  
evidence that may not be correct or faithful.

592 **IA for its own sake** In closing, we would like to  
593 highlight a viewpoint that came up multiple times  
594 in survey responses, which was to question the  
595 premise of this paper, i.e., to measure the impact of  
596 IA on NLP. Many respondents noted that they see  
597 IA work as being a valuable scientific pursuit in its  
598 own right, stating that *“Without it, we’re not doing*  
599 *science,”* or *“It’s cool! That’s enough for me.”* Re-  
600 spondents further criticized the often performance-  
601 focused definitions of utility, progress, and impact.  
602 One respondent noted that these definitions of util-  
603 ity have been determined *“by extrinsic sociological*  
604 *factors in the broader field of AI”*. We sympathize  
605 with this observation and note that the focus on  
606 performance is a feature of NLP at this point in  
607 time. What we value might change going forward,  
608 especially as NLP systems are increasingly part of  
609 our daily lives, and qualities such as robustness and  
610 fairness become even more important.  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621

## 622 8 Conclusion

623 We contribute a mixed-methods analysis of the im-  
624 pact of interpretability and analysis research on  
625 NLP. By analyzing a citation graph of 185K+ pa-  
626 pers built from all papers published at ACL and  
627 EMNLP from 2018 to 2023, surveying 138 respon-  
628 dents from the NLP community, and manually an-  
629 notating 556 papers, we found that IA work is well-  
630 cited in other subfields of NLP, central to the NLP  
631 citation graph, and highly influential to many novel  
632 methods. NLP researchers and practitioners per-  
633 ceive IA work as important for progress in NLP,  
634 multiple subfields (especially reasoning and fair-  
635 ness), and for their own work. In sum, even though  
636 highly influential models, methods and datasets are  
637 not driven by IA findings, IA work still has a great  
638 impact on NLP in the past and the present. We con-  
639 clude with a call to action based on what is missing  
640 in the subfield, to pave the way for IA work to be  
641 even more impactful in the future.



## 642 Limitations

### 643 Focus on papers published at ACL and EMNLP

644 The starting point of our analysis are all papers published at ACL and EMNLP. Although these are the  
645 most cited \*CL venues (Mohammad, 2020), our  
646 analysis excludes several other big NLP venues, including EACL, NAACL, AACL, TACL, and workshops, including BlackboxNLP, which focuses on  
647 IA work. Additionally, given the growing interest in NLP, and in particular, LLMs, from the broader  
648 machine learning community, there is an increasing  
649 number of IA papers published at machine learning  
650 conferences such as ICLR, NeurIPS, and ICML,  
651 which we also do not consider in our analyses. Similarly, a vast amount of work on mechanistic interpretability has been published as articles (e.g., on  
652 LessWrong<sup>6</sup> and the AI Alignment Forum<sup>7</sup>), and  
653 blog posts (e.g., by Anthropic<sup>8</sup>). Therefore, there is  
654 a risk that our analysis misses potentially influential  
655 IA work published at these venues.

662 This is mitigated to an extent by our survey, where respondents mention some of these papers and blog posts, which we then discuss in our paper. In addition, the set of papers we consider for our analysis is very large (our initial set contains 477 IA papers). This makes us confident that the findings we draw from these papers (and those citing them) are representative of broader trends in the impact of IA research in NLP. We leave it to future work to investigate the impact of IA work published outside of established NLP venues.

673 **Focus on 2018 to 2024** Our analysis focuses on papers published between 2018 and 2024. Our results thus represent a snapshot in time on the scale of research in NLP, where models and methods come and go. The time period that we look at is dominated by transformer-based language models, and a paradigm of using large, general-purpose pre-trained models for many tasks, and thus many IA papers focus on studying these. Understanding this as the context of our analysis and results is important, as they may look completely different in a time period where the most popular models are different or the most popular IA methods are different. This also means that our results cannot speak to the impact of today’s IA work as its true impact might only become clear in the future.

<sup>6</sup><https://www.lesswrong.com/>

<sup>7</sup><https://www.alignmentforum.org/>

<sup>8</sup><https://www.anthropic.com/>

689 **Not all citations are equal** Although our use of  
690 citations is an important component of how we  
691 quantify impact in this paper, we do not consider  
692 citational context or distinguish between types of citations. However, papers can be cited for a number of reasons (Bornmann and Daniel, 2008), not all positive and not all having to do with the conventions of scholarly publishing (Bornmann and Daniel, 2008; Zhu et al., 2015; Bornmann and Marx, 2012). 697

698 **Limitations of our survey** Although we took  
699 steps to get a large number and diversity of survey  
700 responses, and we ensured a minimum of 10 respondents per bucket when reporting disaggregated results, the 138 responses we received may not be representative of the field as a whole. In particular, full professors (N=5, at various career stages), and industry practitioners who are not researchers (N=1) were somewhat underrepresented in our responses, indicating that our results focus more on research impact rather than impact on industry applications, and are overwhelmingly shaped by PhD students (41.3% of respondents), whose interests, incentives, and assessment of impact are sure to be different from respondents at other career stages. 712

713 Some respondents brought up the following concerns: one respondent felt our definition of IA was too broad for their taste, but our inclusion of interpretability *and* analysis was by design (see Section 3). Another respondent noted that we defined IA but not what we meant by “progress,” which was also by design, as we did not want to impose a normative definition of progress on our respondents but rather, get at their own intuitions, regardless of how they might define progress. Finally, one respondent complained that our questions about the usefulness of IA (to various subfields, on one’s own research, etc.) were framed in absolute rather than relative terms, and that just because IA research has some positive impact on our understanding doesn’t mean that it is the best option to pursue given limited time and resources. This paper presents views of absolute *and* relative impact via the survey and citation graph analyses, for a holistic view of IA research that also allows for it to have value for its own sake. Ultimately, we believe that a view of “optimal” impact compared to other options lies in the eye of the beholder, and is one (but not the only) way of interpreting our results. 736

## References

George A Barnett, Catherine Huh, Youngju Kim, and Han Woo Park. 2011. Citations among communication journals and other disciplines: a network analysis. *Scientometrics*, 88(2):449–469.

Yonatan Belinkov and James Glass. 2019. [Analysis methods in neural language processing: A survey](#). *Transactions of the Association for Computational Linguistics*, 7:49–72.

Mariette Bengtsson. 2016. [How to plan and perform a qualitative study using content analysis](#). *NursingPlus Open*, 2:8–14.

Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S. Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, Erik Brynjolfsson, Shyamal Buch, Dallas Card, Rodrigo Castellon, Niladri Chatterji, Annie Chen, Kathleen Creel, Jared Quincy Davis, Dora Demszky, Chris Donahue, Moussa Doumbouya, Esin Durmus, Stefano Ermon, John Etchemendy, Kawin Ethayarajh, Li Fei-Fei, Chelsea Finn, Trevor Gale, Lauren Gillespie, Karan Goel, Noah Goodman, Shelby Grossman, Neel Guha, Tatsunori Hashimoto, Peter Henderson, John Hewitt, Daniel E. Ho, Jenny Hong, Kyle Hsu, Jing Huang, Thomas Icard, Saahil Jain, Dan Jurafsky, Pratyusha Kalluri, Siddharth Karamcheti, Geoff Keeling, Fereshche Khani, Omar Khattab, Pang Wei Koh, Mark Krass, Ranjay Krishna, Rohith Kuditipudi, Ananya Kumar, Faisal Ladhak, Mina Lee, Tony Lee, Jure Leskovec, Isabelle Levent, Xiang Lisa Li, Xuechen Li, Tengyu Ma, Ali Malik, Christopher D. Manning, Suvir Mirchandani, Eric Mitchell, Zanele Muniyikwa, Suraj Nair, Avanika Narayan, Deepak Narayanan, Ben Newman, Allen Nie, Juan Carlos Niebles, Hamed Nilforoshan, Julian Nyarko, Giray Ogut, Laurel Orr, Isabel Papadimitriou, Joon Sung Park, Chris Piech, Eva Portelance, Christopher Potts, Aditi Raghunathan, Rob Reich, Hongyu Ren, Frieda Rong, Yusuf Roohani, Camilo Ruiz, Jack Ryan, Christopher Ré, Dorsa Sadigh, Shiori Sagawa, Keshav Santhanam, Andy Shih, Krishnan Srinivasan, Alex Tamkin, Rohan Taori, Armin W. Thomas, Florian Tramèr, Rose E. Wang, William Wang, Bohan Wu, Jiajun Wu, Yuhuai Wu, Sang Michael Xie, Michihiro Yasunaga, Jiaxuan You, Matei Zaharia, Michael Zhang, Tianyi Zhang, Xikun Zhang, Yuhui Zhang, Lucia Zheng, Kaitlyn Zhou, and Percy Liang. 2022. [On the opportunities and risks of foundation models](#). *Preprint*, arXiv:2108.07258.

Lutz Bornmann and Hans-Dieter Daniel. 2008. [What do citation counts measure? a review of studies on citing behavior](#). *Journal of Documentation*, 64(1):45–80.

Lutz Bornmann and Werner Marx. 2012. [The anna karenina principle: A way of thinking about success in science](#). *Journal of the American Society for Information Science and Technology*, 63(10):2037–2051.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind

Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.

Xiomara S. Q. Chacon, Thiago C. Silva, and Diego R. Amancio. 2020. [Comparing the impact of subfields in scientific journals](#). *Scientometrics*, 125(1):625–639.

Arman Cohan, Waleed Ammar, Madeleine Van Zuylen, and Field Cady. 2019. Structural scaffolds for citation intent classification in scientific publications. *arXiv preprint arXiv:1904.01608*.

Arman Cohan, Sergey Feldman, Iz Beltagy, Doug Downey, and Daniel S Weld. 2020. Specter: Document-level representation learning using citation-informed transformers. *arXiv preprint arXiv:2004.07180*.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. 2021. A mathematical framework for transformer circuits. *Transformer Circuits Thread*. <https://transformer-circuits.pub/2021/framework/index.html>.

Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. 2021. [Transformer feed-forward layers are key-value memories](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 5484–5495, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Jennifer Golbeck. 2015. *Introduction to social media investigation: A hands-on approach*. Syngress.

Bryce Goodman and Seth Flaxman. 2017. [European union regulations on algorithmic decision making and a “right to explanation”](#). *AI Magazine*, 38(3):50–57.

854	Google. 2024. <a href="#">Generative ai in search: Let google do the searching for you.</a>	912
855		913
856	Sireesh Gururaja, Amanda Bertsch, Clara Na, David	914
857	Widder, and Emma Strubell. 2023. <a href="#">To build our</a>	915
858	<a href="#">future, we must know our past: Contextualizing</a>	916
859	<a href="#">paradigm shifts in natural language processing.</a> In	
860	<i>Proceedings of the 2023 Conference on Empirical</i>	
861	<i>Methods in Natural Language Processing</i> , pages	
862	13310–13325, Singapore. Association for Compu-	
863	tational Linguistics.	
864	John PA Ioannidis, Kevin Boyack, and Paul F Wouters.	
865	2016. Citation metrics: a primer on how (not) to	
866	normalize. <i>PLoS biology</i> , 14(9):e1002542.	
867	Sehrish Iqbal, Saeed-Ul Hassan, Naif Radi Aljohani,	
868	Salem Alelyani, Raheel Nawaz, and Lutz Bornmann.	
869	2021. <a href="#">A decade of in-text citation analysis based on</a>	
870	<a href="#">natural language processing and machine learning</a>	
871	<a href="#">techniques: an overview of empirical studies.</a> <i>Scien-</i>	
872	<i>tometrics</i> , 126(8):6551–6599.	
873	Alon Jacovi. 2023. <a href="#">Trends in explainable ai (xai) litera-</a>	
874	<a href="#">ture.</a> <i>ArXiv</i> , abs/2301.05433.	
875	Andrej Karpathy, Justin Johnson, and Li Fei-Fei. 2015.	
876	<a href="#">Visualizing and understanding recurrent networks.</a>	
877	<i>Preprint</i> , arXiv:1506.02078.	
878	Edward Kim, Darryl Hannan, and Garrett Kenyon. 2018.	
879	<a href="#">Deep sparse coding for invariant multimodal halle</a>	
880	<a href="#">berry neurons.</a> <i>Preprint</i> , arXiv:1711.07998.	
881	Rodney Michael Kinney, Chloe Anastasiades, Rus-	
882	sell Authur, Iz Beltagy, Jonathan Bragg, Alexan-	
883	dra Buraczynski, Isabel Cachola, Stefan Candra, Yo-	
884	ganand Chandrasekhar, Arman Cohan, Miles Craw-	
885	ford, Doug Downey, Jason Dunkelberger, Oren Et-	
886	zioni, Rob Evans, Sergey Feldman, Joseph Gorney,	
887	David W. Graham, F.Q. Hu, Regan Huff, Daniel King,	
888	Sebastian Kohlmeier, Bailey Kuehl, Michael Langan,	
889	Daniel Lin, Haokun Liu, Kyle Lo, Jaron Lochner,	
890	Kelsey MacMillan, Tyler C. Murray, Christopher	
891	Newell, Smita R Rao, Shaurya Rohatgi, Paul Sayre,	
892	Zejiang Shen, Amanpreet Singh, Luca Soldaini, Shiv-	
893	ashankar Subramanian, A. Tanaka, Alex D Wade,	
894	Linda M. Wagner, Lucy Lu Wang, Christopher Wil-	
895	helm, Caroline Wu, Jiangjiang Yang, Angele Zamar-	
896	ron, Madeleine van Zuylen, and Daniel S. Weld. 2023.	
897	<a href="#">The semantic scholar open data platform.</a> <i>ArXiv</i> ,	
898	abs/2301.10140.	
899	Mike Lewis, Yinhan Liu, Naman Goyal, Marjan	
900	Ghazvininejad, Abdelrahman Mohamed, Omer Levy,	
901	Veselin Stoyanov, and Luke Zettlemoyer. 2020.	
902	<a href="#">BART: Denoising sequence-to-sequence pre-training</a>	
903	<a href="#">for natural language generation, translation, and com-</a>	
904	<a href="#">prehension.</a> In <i>Proceedings of the 58th Annual Meet-</i>	
905	<i>ing of the Association for Computational Linguistics</i> ,	
906	pages 7871–7880, Online. Association for Computa-	
907	tional Linguistics.	
908	Loet Leydesdorff. 2007. Betweenness centrality as an	
909	indicator of the interdisciplinarity of scientific jour-	
910	nals. <i>Journal of the American Society for Information</i>	
911	<i>Science and Technology</i> , 58(9):1303–1319.	
	Loet Leydesdorff, Caroline S Wagner, and Lutz Born-	912
	mann. 2018. Betweenness and diversity in journal ci-	913
	tation networks as measures of interdisciplinarity—a	914
	tribute to eugene garfield. <i>Scientometrics</i> , 114:567–	915
	592.	916
	Xiang Lisa Li and Percy Liang. 2021. <a href="#">Prefix-tuning:</a>	917
	<a href="#">Optimizing continuous prompts for generation.</a> In	918
	<i>Proceedings of the 59th Annual Meeting of the Asso-</i>	919
	<i>ciation for Computational Linguistics and the 11th</i>	920
	<i>International Joint Conference on Natural Language</i>	921
	<i>Processing (Volume 1: Long Papers)</i> , pages 4582–	922
	4597, Online. Association for Computational Lin-	923
	guistics.	924
	Zachary C. Lipton. 2018. <a href="#">The mythos of model inter-</a>	925
	<a href="#">pretability.</a> <i>Commun. ACM</i> , 61(10):36–43.	926
	Andreas Madsen, Himabindu Lakkaraju, Siva Reddy,	927
	and Sarath Chandar. 2024. <a href="#">Interpretability needs a</a>	928
	<a href="#">new paradigm.</a> <i>ArXiv</i> , abs/2405.05386.	929
	Kevin Meng, David Bau, Alex J Andonian, and Yonatan	930
	Belinkov. 2022. <a href="#">Locating and editing factual associ-</a>	931
	<a href="#">ations in GPT.</a> In <i>Advances in Neural Information</i>	932
	<i>Processing Systems.</i>	933
	Microsoft. 2023. <a href="#">Copilot your everyday ai companion.</a>	934
	Staša Milojević, Filippo Radicchi, and Judit Bar-Ilan.	935
	2017. <a href="#">Citation success index - an intuitive pair-wise</a>	936
	<a href="#">journal comparison metric.</a> <i>Journal of Informetrics</i> ,	937
	11(1):223–231.	938
	Sewon Min, Xinxu Lyu, Ari Holtzman, Mikel Artetxe,	939
	Mike Lewis, Hannaneh Hajishirzi, and Luke Zettle-	940
	moyer. 2022. <a href="#">Rethinking the role of demonstrations:</a>	941
	<a href="#">What makes in-context learning work?</a> In <i>Proceed-</i>	942
	<i>ings of the 2022 Conference on Empirical Methods in</i>	943
	<i>Natural Language Processing</i> , pages 11048–11064,	944
	Abu Dhabi, United Arab Emirates. Association for	945
	Computational Linguistics.	946
	Saif M. Mohammad. 2020. <a href="#">Examining citations of nat-</a>	947
	<a href="#">ural language processing literature.</a> In <i>Proceedings</i>	948
	<i>of the 58th Annual Meeting of the Association for</i>	949
	<i>Computational Linguistics</i> , pages 5199–5209, On-	950
	line. Association for Computational Linguistics.	951
	Moin Nadeem, Anna Bethke, and Siva Reddy. 2021.	952
	<a href="#">StereoSet: Measuring stereotypical bias in pretrained</a>	953
	<a href="#">language models.</a> In <i>Proceedings of the 59th Annual</i>	954
	<i>Meeting of the Association for Computational Lin-</i>	955
	<i>guistics and the 11th International Joint Conference</i>	956
	<i>on Natural Language Processing (Volume 1: Long</i>	957
	<i>Papers)</i> , pages 5356–5371, Online. Association for	958
	Computational Linguistics.	959
	Nikita Nangia, Clara Vania, Rasika Bhalerao, and	960
	Samuel R. Bowman. 2020. <a href="#">CrowS-pairs: A chal-</a>	961
	<a href="#">lenge dataset for measuring social biases in masked</a>	962
	<a href="#">language models.</a> In <i>Proceedings of the 2020 Con-</i>	963
	<i>ference on Empirical Methods in Natural Language</i>	964
	<i>Processing (EMNLP)</i> , pages 1953–1967, Online. As-	965
	sociation for Computational Linguistics.	966

967	Jeppe Nicolaisen. 2007. <a href="#">Citation analysis</a> . <i>Annual Review of Information Science and Technology</i> , 41(1):609–641.	
970	OpenAI. 2022. <a href="#">Introducing chatgpt</a> .	
971	OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rameesh Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov,	1029 1030 1031 1032 1033 1034 1035 1036 1037 1038 1039 1040 1041 1042 1043 1044 1045 1046 1047 1048 1049 1050 1051 1052 1053 1054 1055 1056 1057 1058
	Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. 2024. <a href="#">Gpt-4 technical report</a> . <i>Preprint</i> , arXiv:2303.08774.	1029 1030 1031 1032 1033 1034 1035 1036 1037 1038 1039 1040 1041 1042 1043 1044 1045 1046 1047 1048 1049 1050 1051 1052 1053 1054 1055 1056 1057 1058
	Clíodhna O’Connor and Helene Joffe. 2020. <a href="#">Intercoder reliability in qualitative research: Debates and practical guidelines</a> . <i>International Journal of Qualitative Methods</i> , 19:1609406919899220.	1059 1060 1061 1062
	Aniket Pramanick, Yufang Hou, Saif Mohammad, and Iryna Gurevych. 2023. <a href="#">A diachronic analysis of paradigm shifts in NLP research: When, how, and why?</a> In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 2312–2326, Singapore. Association for Computational Linguistics.	1063 1064 1065 1066 1067 1068 1069
	Jason Priem, Heather Piwowar, and Richard Orr. 2022. <a href="#">Openalex: A fully-open index of scholarly works, authors, venues, institutions, and concepts</a> . <i>arXiv preprint arXiv:2205.01833</i> .	1070 1071 1072 1073
	Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. <a href="#">Language models are unsupervised multitask learners</a> .	1074 1075 1076
	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. <a href="#">Exploring the limits of transfer learning with a unified text-to-text transformer</a> . <i>Journal of Machine Learning Research</i> , 21(140):1–67.	1077 1078 1079 1080 1081 1082
	T. Rauker, A. Ho, S. Casper, and D. Hadfield-Menell. 2023. <a href="#">Toward transparent ai: A survey on interpreting the inner structures of deep neural networks</a> . In <i>2023 IEEE Conference on Secure and Trustworthy Machine Learning (SaTML)</i> , pages 464–483, Los Alamitos, CA, USA. IEEE Computer Society.	1083 1084 1085 1086 1087 1088

1089	Anna Rogers, Olga Kovaleva, and Anna Rumshisky.	Barth-Maron, William Wong, Rishabh Joshi, Rahma	1149
1090	2020. <a href="#">A primer in BERTology: What we know about</a>	Chaabouni, Deeni Fatiha, Arun Ahuja, Gaurav Singh	1150
1091	<a href="#">how BERT works</a> . <i>Transactions of the Association</i>	Tomar, Evan Senter, Martin Chadwick, Ilya Kor-	1151
1092	<a href="#">for Computational Linguistics</a> , 8:842–866.	nakov, Nithya Attaluri, Iñaki Iturrate, Ruibo Liu,	1152
1093	Johnny Saldana. 2021. <i>The coding manual for qual-</i>	Yunxuan Li, Sarah Cogan, Jeremy Chen, Chao Jia,	1153
1094	<i>itative researchers</i> , 4 edition. SAGE Publications,	Chenjie Gu, Qiao Zhang, Jordan Grimstad, Ale Jakse	1154
1095	London, England.	Hartman, Xavier Garcia, Thanumalayan Sankara-	1155
1096	John J. Shaughnessy, Eugene B. Zechmeister, and	narayana Pillai, Jacob Devlin, Michael Laskin, Diego	1156
1097	Jeanne S. Zechmeister. 2015. <i>Research methods in</i>	de Las Casas, Dasha Valter, Connie Tao, Lorenzo	1157
1098	<i>psychology</i> , tenth edition edition. McGraw-Hill Edu-	Blanco, Adrià Puigdomènech Badia, David Reitter,	1158
1099	cation, Dubuque.	Mianna Chen, Jenny Brennan, Clara Rivera, Sergey	1159
1100	Janvijay Singh, Mukund Rungta, Diyi Yang, and Saif	Brin, Shariq Iqbal, Gabriela Surita, Jane Labanowski,	1160
1101	Mohammad. 2023. <a href="#">Forgotten knowledge: Examin-</a>	Abhi Rao, Stephanie Winkler, Emilio Parisotto, Yim-	1161
1102	<a href="#">ing the citational amnesia in NLP</a> . In <i>Proceedings</i>	ing Gu, Kate Olszewska, Ravi Addanki, Antoine	1162
1103	<a href="#">of the 61st Annual Meeting of the Association for</a>	Miech, Annie Louis, Denis Tepliyashin, Geoff Brown,	1163
1104	<a href="#">Computational Linguistics (Volume 1: Long Papers)</a> ,	Elliot Catt, Jan Balaguer, Jackie Xiang, Pidong Wang,	1164
1105	pages 6192–6208, Toronto, Canada. Association for	Zoe Ashwood, Anton Briukhov, Albert Webson, San-	1165
1106	Computational Linguistics.	jay Ganapathy, Smit Sanghavi, Ajay Kannan, Ming-	1166
1107	Iman Tahamtan and Lutz Bornmann. 2019. <a href="#">What do ci-</a>	Wei Chang, Axel Stjerngren, Josip Djolonga, Yut-	1167
1108	<a href="#">tation counts measure? an updated review of studies</a>	ing Sun, Ankur Bapna, Matthew Aitchison, Pedram	1168
1109	<a href="#">on citations in scientific documents published be-</a>	Pejman, Henryk Michalewski, Tianhe Yu, Cindy	1169
1110	<a href="#">tween 2006 and 2018</a> . <i>Scientometrics</i> , 121(3):1635–	Wang, Juliette Love, Junwhan Ahn, Dawn Bloxwich,	1170
1111	1684.	Kehang Han, Peter Humphreys, Thibault Sellam,	1171
1112	Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-	James Bradbury, Varun Godbole, Sina Samangooui,	1172
1113	Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan	Bogdan Damoc, Alex Kaskasoli, Sébastien M. R.	1173
1114	Schalkwyk, Andrew M. Dai, Anja Hauth, Katie	Arnold, Vijay Vasudevan, Shubham Agrawal, Jason	1174
1115	Millican, David Silver, Melvin Johnson, Ioannis	Riesa, Dmitry Lepikhin, Richard Tanburn, Srivat-	1175
1116	Antonoglou, Julian Schrittwieser, Amelia Glaese,	san Srinivasan, Hyeontaek Lim, Sarah Hodgkinson,	1176
1117	Jilin Chen, Emily Pitler, Timothy Lillicrap, Ange-	Pranav Shyam, Johan Ferret, Steven Hand, Ankush	1177
1118	liki Lazaridou, Orhan Firat, James Molloy, Michael	Garg, Tom Le Paine, Jian Li, Yujia Li, Minh Gi-	1178
1119	Isard, Paul R. Barham, Tom Hennigan, Benjamin	ang, Alexander Neitz, Zaheer Abbas, Sarah York,	1179
1120	Lee, Fabio Viola, Malcolm Reynolds, Yuanzhong	Machel Reid, Elizabeth Cole, Aakanksha Chowdh-	1180
1121	Xu, Ryan Doherty, Eli Collins, Clemens Meyer, Eliza	ery, Dipanjan Das, Dominika Rogozińska, Vitaliy	1181
1122	Rutherford, Erica Moreira, Kareem Ayoub, Megha	Nikolaev, Pablo Sprechmann, Zachary Nado, Lukas	1182
1123	Goel, Jack Krawczyk, Cosmo Du, Ed Chi, Heng-	Zilka, Flavien Prost, Luheng He, Marianne Mon-	1183
1124	Tze Cheng, Eric Ni, Purvi Shah, Patrick Kane, Betty	teiro, Gaurav Mishra, Chris Welty, Josh Newlan,	1184
1125	Chan, Manaal Faruqui, Aliaksei Severyn, Hanzhao	Dawei Jia, Miltiadis Allamanis, Clara Huiyi Hu,	1185
1126	Lin, YaGuang Li, Yong Cheng, Abe Ittycheriah,	Raoul de Liedekerke, Justin Gilmer, Carl Saroufim,	1186
1127	Mahdis Mahdieh, Mia Chen, Pei Sun, Dustin Tran,	Shruti Rijhwani, Shaobo Hou, Disha Shrivastava,	1187
1128	Sumit Bagri, Balaji Lakshminarayanan, Jeremiah	Anirudh Baddepudi, Alex Goldin, Adnan Ozturel,	1188
1129	Liu, Andras Orban, Fabian Göra, Hao Zhou, Xiny-	Albin Cassirer, Yunhan Xu, Daniel Sohn, Deven-	1189
1130	ing Song, Aurelien Boffy, Harish Ganapathy, Steven	dra Sachan, Reinald Kim Amplayo, Craig Swans-	1190
1131	Zheng, HyunJeong Choe, Ágoston Weisz, Tao Zhu,	son, Dessie Petrova, Shashi Narayan, Arthur Guez,	1191
1132	Yifeng Lu, Siddharth Gopal, Jarrod Kahn, Maciej	Siddhartha Brahma, Jessica Landon, Miteyan Pa-	1192
1133	Kula, Jeff Pitman, Rushin Shah, Emanuel Taropa,	tel, Ruizhe Zhao, Kevin Vilella, Luyu Wang, Wen-	1193
1134	Majd Al Merey, Martin Baeuml, Zhifeng Chen, Lau-	hao Jia, Matthew Rahtz, Mai Giménez, Legg Yeung,	1194
1135	rent El Shafey, Yujing Zhang, Olcan Sercinoglu,	James Keeling, Petko Georgiev, Diana Mincu, Boxi	1195
1136	George Tucker, Enrique Piqueras, Maxim Krikun,	Wu, Salem Haykal, Rachel Saputro, Kiran Vodra-	1196
1137	Iain Barr, Nikolay Savinov, Ivo Danihelka, Becca	halli, James Qin, Zeynep Cankara, Abhanshu Sharma,	1197
1138	Roelofs, Anaïs White, Anders Andreassen, Tamara	Nick Fernando, Will Hawkins, Behnam Neyshabur,	1198
1139	von Glehn, Lakshman Yagati, Mehran Kazemi, Luc-	Solomon Kim, Adrian Hutter, Priyanka Agrawal,	1199
1140	cas Gonzalez, Misha Khalman, Jakub Sygnowski,	Alex Castro-Ros, George van den Driessche, Tao	1200
1141	Alexandre Frechette, Charlotte Smith, Laura Culp,	Wang, Fan Yang, Shuo yiin Chang, Paul Komarek,	1201
1142	Lev Prolev, Yi Luan, Xi Chen, James Lottes, Nathan	Ross McIlroy, Mario Lučić, Guodong Zhang, Wael	1202
1143	Schucher, Federico Lebron, Alban Rrustemi, Na-	Farhan, Michael Sharman, Paul Natsev, Paul Michel,	1203
1144	talie Clay, Phil Crone, Tomas Kocisky, Jeffrey Zhao,	Yamini Bansal, Siyuan Qiao, Kris Cao, Siamak Shak-	1204
1145	Bartek Perz, Dian Yu, Heidi Howard, Adam Blo-	eri, Christina Butterfield, Justin Chung, Paul Kishan	1205
1146	niarz, Jack W. Rae, Han Lu, Laurent Sifre, Mar-	Rubenstein, Shivani Agrawal, Arthur Mensch, Kedar	1206
1147	cello Maggioni, Fred Alcober, Dan Garrette, Megan	Soparkar, Karel Lenc, Timothy Chung, Aedan Pope,	1207
1148	Barnes, Shantanu Thakoor, Jacob Austin, Gabriel	Loren Maggiore, Jackie Kay, Priya Jhakra, Shibo	1208
		Wang, Joshua Maynez, Mary Phuong, Taylor Tobin,	1209
		Andrea Tacchetti, Maja Trebacz, Kevin Robinson,	1210
		Yash Katariya, Sebastian Riedel, Paige Bailey, Kefan	1211
		Xiao, Nimesh Ghelani, Lora Aroyo, Ambrose Slone,	1212

1213	Neil Houlsby, Xuehan Xiong, Zhen Yang, Elena Gribovskaya, Jonas Adler, Mateo Wirth, Lisa Lee, Music Li, Thais Kagohara, Jay Pavagadhi, Sophie Bridgers, Anna Bortsova, Sanjay Ghemawat, Zafarali Ahmed, Tianqi Liu, Richard Powell, Vijay Bolina, Mariko Iinuma, Polina Zablotskaia, James Besley, Da-Woon Chung, Timothy Dozat, Ramona Comanescu, Xi-ance Si, Jeremy Greer, Guolong Su, Martin Polacek, Raphaël Lopez Kaufman, Simon Tokumine, Hexiang Hu, Elena Buchatskaya, Yingjie Miao, Mohamed Elhawaty, Aditya Siddhant, Nenad Tomasev, Jinwei Xing, Christina Greer, Helen Miller, Shereen Ashraf, Aurko Roy, Zizhao Zhang, Ada Ma, Angelos Filos, Milos Besta, Rory Blevins, Ted Klimenko, Chih-Kuan Yeh, Soravit Changpinyo, Jiaqi Mu, Oscar Chang, Mantas Pajarskas, Carrie Muir, Vered Cohen, Charline Le Lan, Krishna Haridasan, Amit Marathe, Steven Hansen, Sholto Douglas, Rajkumar Samuel, Mingqiu Wang, Sophia Austin, Chang Lan, Jiepu Jiang, Justin Chiu, Jaime Alonso Lorenzo, Lars Lowe Sjöstrand, Sébastien Cevey, Zach Gleicher, Thi Avrahami, Anudhyan Boral, Hansa Srinivasan, Vittorio Selo, Rhys May, Konstantinos Aisopos, Léonard Hussenot, Livio Baldini Soares, Kate Baumli, Michael B. Chang, Adrià Recasens, Ben Caine, Alexander Pritzel, Filip Pavetic, Fabio Pardo, Anita Gergely, Justin Frye, Vinay Ramasesh, Dan Horgan, Kartikeya Badola, Nora Kassner, Subhrajit Roy, Ethan Dyer, Víctor Campos Campos, Alex Tomala, Yunhao Tang, Dalia El Badawy, Elspeth White, Basil Mustafa, Oran Lang, Abhishek Jindal, Sharad Vikram, Zhitao Gong, Sergi Caelles, Ross Hemsley, Gregory Thornton, Fangxiaoyu Feng, Wojciech Stokowiec, Ce Zheng, Phoebe Thacker, Çağlar Ünlü, Zhishuai Zhang, Mohammad Saleh, James Svensson, Max Bileschi, Piyush Patil, Ankesh Anand, Roman Ring, Katerina Tsihlias, Arpi Vezer, Marco Selvi, Toby Shevlane, Mikel Rodriguez, Tom Kwiatkowski, Samira Daruki, Keran Rong, Allan Dafoe, Nicholas FitzGerald, Keren Gu-Lemberg, Mina Khan, Lisa Anne Hendricks, Marie Pellat, Vladimir Feinberg, James Cobon-Kerr, Tara Sainath, Maribeth Rauh, Sayed Hadi Hashemi, Richard Ives, Yana Hasson, Eric Noland, Yuan Cao, Nathan Byrd, Le Hou, Qingze Wang, Thibault Sottiaux, Michela Paganini, Jean-Baptiste Lespiau, Alexandre Moufaret, Samer Hassan, Kaushik Shivakumar, Joost van Amersfoort, Amol Mandhane, Pratik Joshi, Anirudh Goyal, Matthew Tung, Andrew Brock, Hannah Sheahan, Vedant Misra, Cheng Li, Nemanja Rakićević, Mostafa Dehghani, Fangyu Liu, Sid Mittal, Junhyuk Oh, Seb Noury, Eren Sezener, Fantine Huot, Matthew Lamm, Nicola De Cao, Charlie Chen, Sidharth Mudgal, Romina Stella, Kevin Brooks, Gautam Vasudevan, Chenxi Liu, Mainak Chaitin, Nivedita Melinkeri, Aaron Cohen, Venus Wang, Kristie Seymore, Sergey Zubkov, Rahul Goel, Summer Yue, Sai Krishnakumaran, Brian Albert, Nate Hurley, Motoki Sano, Anhad Mohanane, Jonah Joughin, Egor Filonov, Tomasz Kępa, Yomna Eldawy, Jiawern Lim, Rahul Rishi, Shirin Badiezadegan, Taylor Bos, Jerry Chang, Sanil Jain, Sri Gayatri Sundara Padmanabhan, Subha Puttagunta, Kalpesh Krishna, Leslie Baker, Norbert Kalb, Vamsi Bedapudi, Adam Kurzrok, Shuntong Lei, Anthony Yu, Oren Litvin, Xiang Zhou, Zhichun Wu, Sam Sobell, Andrea Siciliano, Alan Papir, Robby Neale, Jonas Bragagnolo, Tej Toor, Tina Chen, Valentin Anklin, Feiran Wang, Richie Feng, Milad Gholami, Kevin Ling, Lijuan Liu, Jules Walter, Hamid Moghaddam, Arun Kishore, Jakob Adamek, Tyler Mercado, Jonathan Mallinson, Siddhinita Wandekar, Stephen Cagle, Eran Ofek, Guillermo Garrido, Clemens Lombriser, Maksim Mukha, Botu Sun, Hafeezul Rahman Mohammad, Josip Matak, Yadi Qian, Vikas Peswani, Pawel Janus, Quan Yuan, Leif Schelin, Oana David, Ankur Garg, Yifan He, Oleksii Duzhyi, Anton Ålgmyr, Timothée Lottaz, Qi Li, Vikas Yadav, Luyao Xu, Alex Chinien, Rakesh Shivanna, Aleksandr Chuklin, Josie Li, Carrie Spadine, Travis Wolfe, Kareem Mohamed, Subhabrata Das, Zihang Dai, Kyle He, Daniel von Dincklage, Shyam Upadhyay, Akanksha Maurya, Luyan Chi, Sebastian Krause, Khalid Salama, Pam G Rabinovitch, Pavan Kumar Reddy M, Aarush Selvan, Mikhail Dektiarev, Golnaz Ghiasi, Erdem Guven, Himanshu Gupta, Boyi Liu, Deepak Sharma, Idan Heimlich Shtacher, Shachi Paul, Oscar Akerlund, François-Xavier Aubet, Terry Huang, Chen Zhu, Eric Zhu, Elico Teixeira, Matthew Fritze, Francesco Bertolini, Liana-Eleonora Marinescu, Martin Bølle, Dominik Paulus, Khyatti Gupta, Tejasi Latkar, Max Chang, Jason Sanders, Roopa Wilson, Xuewei Wu, Yi-Xuan Tan, Lam Nguyen Thiet, Tulsee Doshi, Sid Lall, Swaroop Mishra, Wanming Chen, Thang Luong, Seth Benjamin, Jasmine Lee, Ewa Andrejczuk, Dominik Rabiej, Vipul Ranjan, Krzysztof Styrz, Pengcheng Yin, Jon Simon, Malcolm Rose Harriott, Mudit Bansal, Alexei Robsky, Geoff Bacon, David Greene, Daniil Mirylenka, Chen Zhou, Obaid Sarvana, Abhimanyu Goyal, Samuel Andermatt, Patrick Siegler, Ben Horn, Assaf Israel, Francesco Pongetti, Chih-Wei "Louis" Chen, Marco Selvatici, Pedro Silva, Kathie Wang, Jackson Tolins, Kelvin Guu, Roey Yogev, Xiaochen Cai, Alessandro Agostini, Maulik Shah, Hung Nguyen, Noah Ó Donnaile, Sébastien Pereira, Linda Friso, Adam Stambler, Adam Kurzrok, Chenkai Kuang, Yan Romanikhin, Mark Geller, ZJ Yan, Kane Jiang, Cheng-Chun Lee, Wojciech Fica, Eric Malmi, Qijun Tan, Dan Banica, Daniel Balle, Ryan Pham, Yanping Huang, Diana Avram, Hongzhi Shi, Jasjot Singh, Chris Hidey, Niharika Ahuja, Pranab Saxena, Dan Dooley, Srividya Pranavi Potharaju, Eileen O'Neill, Anand Gokulchandran, Ryan Foley, Kai Zhao, Mike Dusenberry, Yuan Liu, Pulkit Mehta, Ragha Kotikalapudi, Chalence Safranek-Shrader, Andrew Goodman, Joshua Kessinger, Eran Globen, Praatek Kolhar, Chris Gorgolewski, Ali Ibrahim, Yang Song, Ali Eichenbaum, Thomas Brovelli, Sahitya Potluri, Preethi Lahoti, Cip Baetu, Ali Ghorbani, Charles Chen, Andy Crawford, Shalini Pal, Mukund Sridhar, Petru Gurita, Asier Mujika, Igor Petrovski, Pierre-Louis Cedoz, Chenmei Li, Shiyuan Chen, Niccolò Dal Santo, Siddharth Goyal, Jitesh Punjabi, Karthik Kappaganthu, Chester Kwak, Pallavi LV, Sarmishta Velury, Himadri Choudhury, Jamie Hall, Premal Shah, Ricardo Figueira, Matt Thomas,	1277 1278 1279 1280 1281 1282 1283 1284 1285 1286 1287 1288 1289 1290 1291 1292 1293 1294 1295 1296 1297 1298 1299 1300 1301 1302 1303 1304 1305 1306 1307 1308 1309 1310 1311 1312 1313 1314 1315 1316 1317 1318 1319 1320 1321 1322 1323 1324 1325 1326 1327 1328 1329 1330 1331 1332 1333 1334 1335 1336 1337 1338 1339
------	--	--

1340	Minjie Lu, Ting Zhou, Chintu Kumar, Thomas Jurdi, Sharat Chikkerur, Yenai Ma, Adams Yu, Soo Kwak, Victor Áhdel, Sujeevan Rajayogam, Travis Choma, Fei Liu, Aditya Barua, Colin Ji, Ji Ho Park, Vincent Hellendoorn, Alex Bailey, Taylan Bilal, Huanjie Zhou, Mehrdad Khatir, Charles Sutton, Wojciech Rzadkowski, Fiona Macintosh, Konstantin Shagin, Paul Medina, Chen Liang, Jinjing Zhou, Pararth Shah, Yingying Bi, Attila Dankovics, Shipra Banga, Sabine Lehmann, Marissa Bredesen, Zifan Lin, John Eric Hoffmann, Jonathan Lai, Raynald Chung, Kai Yang, Nihal Balani, Arthur Bražinskis, Andrei Sozanschi, Matthew Hayes, Héctor Fernández Alcalde, Peter Makarov, Will Chen, Antonio Stella, Liselotte Snijders, Michael Mandl, Ante Kärman, Paweł Nowak, Xinyi Wu, Alex Dyck, Krishnan Vaidyanathan, Raghavender R, Jessica Mallet, Mitch Rudominer, Eric Johnston, Sushil Mittal, Akhil Udathu, Janara Christensen, Vishal Verma, Zach Irving, Andreas Santucci, Gamaleldin Elsayed, Elnaz Davoodi, Marin Georgiev, Ian Tenney, Nan Hua, Geoffrey Cideron, Edouard Leurent, Mahmoud Alnahlawi, Ionut Georgescu, Nan Wei, Ivy Zheng, Dylan Scandinaro, Heinrich Jiang, Jasper Snoek, Mukund Sundararajan, Xuezhi Wang, Zack Ontiveros, Itay Karo, Jeremy Cole, Vinu Rajashekhar, Lara Tumeh, Eyal Ben-David, Rishub Jain, Jonathan Uesato, Romina Datta, Oskar Bunyan, Shimu Wu, John Zhang, Piotr Stanczyk, Ye Zhang, David Steiner, Subhajit Naskar, Michael Azzam, Matthew Johnson, Adam Paszke, Chung-Cheng Chiu, Jaume Sanchez Elias, Afroz Mohiuddin, Faizan Muhammad, Jin Miao, Andrew Lee, Nino Vieillard, Jane Park, Jigeng Zhang, Jeff Stanway, Drew Garmon, Abhijit Karmarkar, Zhe Dong, Jong Lee, Aviral Kumar, Luwei Zhou, Jonathan Evens, William Isaac, Geoffrey Irving, Edward Loper, Michael Fink, Isha Arkatkar, Nanxin Chen, Izhak Shafran, Ivan Petrychenko, Zhe Chen, Johnson Jia, Anselm Levskaya, Zhenkai Zhu, Peter Grabowski, Yu Mao, Alberto Magni, Kaisheng Yao, Javier Snaider, Norman Casagrande, Evan Palmer, Paul Suanthan, Alfonso Castaño, Irene Giannoumis, Wooyeol Kim, Mikołaj Rybiński, Ashwin Sreevatsa, Jennifer Prendki, David Soergel, Adrian Goedeckemeyer, Willi Gierke, Mohsen Jafari, Meenu Gaba, Jeremy Wiesner, Diana Gage Wright, Yawen Wei, Harsha Vashisht, Yana Kulizhskaya, Jay Hoover, Maigo Le, Lu Li, Chimezie Iwuanyanwu, Lu Liu, Kevin Ramirez, Andrey Khorlin, Albert Cui, Tian LIN, Marcus Wu, Ricardo Aguilar, Keith Pallo, Abhishek Chakladar, Ginger Perng, Elena Allica Abellan, Mingyang Zhang, Ishita Dasgupta, Nate Kushman, Ivo Penchev, Alena Repina, Xihui Wu, Tom van der Weide, Priya Ponnappalli, Caroline Kaplan, Jiri Simsa, Shuangfeng Li, Olivier Dousse, Fan Yang, Jeff Piper, Nathan Ie, Rama Pasumarthi, Nathan Lintz, Anitha Vijayakumar, Daniel Andor, Pedro Valenzuela, Minnie Lui, Cosmin Paduraru, Daiyi Peng, Katherine Lee, Shuyuan Zhang, Somer Greene, Duc Dung Nguyen, Paula Kurylowicz, Cassidy Hardin, Lucas Dixon, Lili Janzer, Kiam Choo, Ziqiang Feng, Biao Zhang, Achintya Singhal, Dayou Du, Dan McKinnon, Natasha Antropova, Tolga Bolukbasi, Orgad Keller, David Reid, Daniel	Finchelstein, Maria Abi Raad, Remi Crocker, Peter Hawkins, Robert Dadashi, Colin Gaffney, Ken Franko, Anna Bulanova, Rémi Leblond, Shirley Chung, Harry Askham, Luis C. Cobo, Kelvin Xu, Felix Fischer, Jun Xu, Christina Sorokin, Chris Alberti, Chu-Cheng Lin, Colin Evans, Alek Dimitriev, Hannah Forbes, Dylan Banarse, Zora Tung, Mark Omernick, Colton Bishop, Rachel Sterneck, Rohan Jain, Jiawei Xia, Ehsan Amid, Francesco Piccinno, Xingyu Wang, Praseem Banzal, Daniel J. Mankowitz, Alex Polozov, Victoria Krakovna, Sasha Brown, MohammadHossein Bateni, Dennis Duan, Vlad Firoiu, Meghana Thotakuri, Tom Natan, Matthieu Geist, Ser tan Girgin, Hui Li, Jiayu Ye, Ofir Roval, Reiko Tojo, Michael Kwong, James Lee-Thorp, Christopher Yew, Danila Sinopalnikov, Sabela Ramos, John Mellor, Abhishek Sharma, Kathy Wu, David Miller, Nicolas Sonnerat, Denis Vnukov, Rory Greig, Jennifer Beattie, Emily Caveness, Libin Bai, Julian Eischenschlos, Alex Korchemniy, Tomy Tsai, Mimi Jasarevic, Weize Kong, Phuong Dao, Zeyu Zheng, Frederick Liu, Fan Yang, Rui Zhu, Tian Huey Teh, Jason Sanmiya, Evgeny Gladchenko, Nejc Trdin, Daniel Toyama, Evan Rosen, Sasan Tavakkol, Linting Xue, Chen Elkind, Oliver Woodman, John Carpenter, George Papamakarios, Rupert Kemp, Sushant Kafle, Tanya Grunina, Rishika Sinha, Alice Talbert, Diane Wu, Denese Owusu-Afriyie, Cosmo Du, Chloe Thornton, Jordi Pont-Tuset, Pradyumna Narayana, Jing Li, Saaber Fatehi, John Wieting, Omar Ajmeri, Benigno Uria, Yeongil Ko, Laura Knight, Amélie Héliou, Ning Niu, Shane Gu, Chenxi Pang, Yeqing Li, Nir Levine, Ariel Stolovich, Rebecca Santamaria-Fernandez, Sonam Goenka, Wenny Yustalim, Robin Strudel, Ali Elqursh, Charlie Deck, Hyo Lee, Zonglin Li, Kyle Levin, Raphael Hoffmann, Dan Holtmann-Rice, Olivier Bachem, Sho Arora, Christy Koh, Soheil Hassas Yeganeh, Siim Pöder, Mukarram Tariq, Yanhua Sun, Lucian Ionita, Mojtaba Seyedhosseini, Pouya Tafti, Zhiyu Liu, Anmol Gulati, Jasmine Liu, Xinyu Ye, Bart Chrzaszcz, Lily Wang, Nikhil Sethi, Tianrun Li, Ben Brown, Shreya Singh, Wei Fan, Aaron Parisi, Joe Stanton, Vinod Koverkathu, Christopher A. Choquette-Choo, Yunjie Li, TJ Lu, Abe Ittycheriah, Prakash Shroff, Mani Varadarajan, Sanaz Bahargam, Rob Willoughby, David Gaddy, Guillaume Desjardins, Marco Cornero, Brona Robenek, Bhavishya Mittal, Ben Albrecht, Ashish Shenoy, Fedor Moiseev, Henrik Jacobsson, Alireza Ghaffarkhah, Morgane Rivière, Alanna Walton, Clément Crepy, Alicia Parish, Zongwei Zhou, Clement Farabet, Carey Radebaugh, Praveen Srinivasan, Claudia van der Salm, Andreas Fildjeland, Salvatore Scellato, Eri Latorre-Chimoto, Hanna Klimczak-Plucińska, David Bridson, Dario de Cesare, Tom Hudson, Piermaria Mendolicchio, Lexi Walker, Alex Morris, Matthew Mauger, Alexey Guseynov, Alison Reid, Seth Odoom, Lucia Loher, Victor Cotruta, Madhavi Yenugula, Dominik Grewe, Anastasia Petrushkina, Tom Duerig, Antonio Sanchez, Steve Yadlowsky, Amy Shen, Amir Globerson, Lynette Webb, Sahil Dua, Dong Li, Surya Bhupatiraju, Dan Hurt, Haroon Qureshi, Ananth Agarwal, Tomer Shani, Matan Eyal, Anuj	1404 1405 1406 1407 1408 1409 1410 1411 1412 1413 1414 1415 1416 1417 1418 1419 1420 1421 1422 1423 1424 1425 1426 1427 1428 1429 1430 1431 1432 1433 1434 1435 1436 1437 1438 1439 1440 1441 1442 1443 1444 1445 1446 1447 1448 1449 1450 1451 1452 1453 1454 1455 1456 1457 1458 1459 1460 1461 1462 1463 1464 1465 1466 1467
------	---	--	--

1468	Khare, Shreyas Rammohan Belle, Lei Wang, Chetan	Naftali Tishby and Noga Zaslavsky. 2015. <a href="#">Deep learning and the information bottleneck principle</a> . In <i>2015 IEEE Information Theory Workshop (ITW)</i> , pages 1–5.	1532
1469	Tekur, Mihir Sanjay Kale, Jinliang Wei, Ruoxin		1533
1470	Sang, Brennan Saeta, Tyler Liechty, Yi Sun, Yao		1534
1471	Zhao, Stephan Lee, Pandu Nayak, Doug Fritz, Manish		1535
1472	Reddy Vuyyuru, John Aslanides, Nidhi Vyas,	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	1536
1473	Martin Wicke, Xiao Ma, Evgenii Eltyshv, Nina Mar-	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	1537
1474	tin, Hardie Cate, James Manyika, Keyvan Amiri,	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	1538
1475	Yelin Kim, Xi Xiong, Kai Kang, Florian Luisier,	Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton	1539
1476	Nilesh Tripuraneni, David Madras, Mandy Guo,	Ferrer, Moya Chen, Guillem Cucurull, David Esiobu,	1540
1477	Austin Waters, Oliver Wang, Joshua Ainslie, Jason	Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller,	1541
1478	Baldrige, Han Zhang, Garima Pruthi, Jakob Bauer,	Cynthia Gao, Vedanuj Goswami, Naman Goyal, An-	1542
1479	Feng Yang, Riham Mansour, Jason Gelman, Yang Xu,	thony Hartshorn, Saghar Hosseini, Rui Hou, Hakan	1543
1480	George Polovets, Ji Liu, Honglong Cai, Warren Chen,	Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa,	1544
1481	XiangHai Sheng, Emily Xue, Sherjil Ozair, Christof	Isabel Kloumann, Artem Korenev, Punit Singh Koura,	1545
1482	Angermueller, Xiaowei Li, Anoop Sinha, Weiren	Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Di-	1546
1483	Wang, Julia Wiesinger, Emmanouil Koukoumidis,	ana Liskovich, Yinghai Lu, Yuning Mao, Xavier Mar-	1547
1484	Yuan Tian, Anand Iyer, Madhu Gurusurthy, Mark	tinet, Todor Mihaylov, Pushkar Mishra, Igor Moly-	1548
1485	Goldenson, Parashar Shah, MK Blake, Hongkun Yu,	bog, Yixin Nie, Andrew Poulton, Jeremy Reizen-	1549
1486	Anthony Urbanowicz, Jennimaria Palomaki, Chrisan-	stein, Rashi Rungta, Kalyan Saladi, Alan Schelten,	1550
1487	tha Fernando, Ken Durden, Harsh Mehta, Nikola	Ruan Silva, Eric Michael Smith, Ranjan Subrama-	1551
1488	Momchev, Elahe Rahimtoroghi, Maria Georgaki,	nian, Xiaoqing Ellen Tan, Binh Tang, Ross Tay-	1552
1489	Amit Raul, Sebastian Ruder, Morgan Redshaw, Jin-	lor, Adina Williams, Jian Xiang Kuan, Puxin Xu,	1553
1490	hyuk Lee, Denny Zhou, Komal Jalan, Dinghua Li,	Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan,	1554
1491	Blake Hechtman, Parker Schuh, Milad Nasr, Kieran	Melanie Kambadur, Sharan Narang, Aurelien Rod-	1555
1492	Milan, Vladimir Mikulik, Juliana Franco, Tim Green,	riguez, Robert Stojnic, Sergey Edunov, and Thomas	1556
1493	Nam Nguyen, Joe Kelley, Aroma Mahendru, Andrea	Scialom. 2023. <a href="#">Llama 2: Open foundation and fine-tuned chat models</a> . <i>Preprint</i> , arXiv:2307.09288.	1557
1494	Hu, Joshua Howland, Ben Vargas, Jeffrey Hui, Kshij		1558
1495	tij Bansal, Vikram Rao, Rakesh Ghiya, Emma Wang,	Marco Valenzuela, Vu Ha, and Oren Etzioni. 2015. Identi-	1559
1496	Ke Ye, Jean Michel Sarr, Melanie Moranski Preston,	fying meaningful citations. In <i>Workshops at the</i>	1560
1497	Madeleine Elish, Steve Li, Aakash Kaku, Jigar Gupta,	<i>twenty-ninth AAAI conference on artificial intelli-</i>	1561
1498	Ice Pasupat, Da-Cheng Juan, Milan Someswar, Tejvi	<i>gence</i> .	1562
1499	M., Xinyun Chen, Aida Amini, Alex Fabrikant, Eric	Jan Philip Wahle, Terry Ruas, Mohamed Abdalla, Bela	1563
1500	Chu, Xuanyi Dong, Amruta Muthal, Senaka Buth-	Gipp, and Saif M Mohammad. 2023. We are who we	1564
1501	pitiya, Sarthak Jauhari, Nan Hua, Urvashi Khan-	cite: Bridges of influence between natural language	1565
1502	delwal, Ayal Hitron, Jie Ren, Larissa Rinaldi, Sha-	processing and other academic fields. <i>arXiv preprint</i>	1566
1503	har Drath, Avigail Dabush, Nan-Jiang Jiang, Har-	<i>arXiv:2310.14870</i> .	1567
1504	shal Godhia, Uli Sachs, Anthony Chen, Yicheng	Kevin Ro Wang, Alexandre Variengien, Arthur Conmy,	1568
1505	Fan, Hagai Taitelbaum, Hila Noga, Zhuyun Dai,	Buck Shlegeris, and Jacob Steinhardt. 2023. <a href="#">Inter-</a>	1569
1506	James Wang, Chen Liang, Jenny Hamer, Chun-Sung	<a href="#">pretability in the wild: a circuit for indirect object</a>	1570
1507	Ferng, Chenel Elkind, Aviel Atias, Paulina Lee, Vít	<a href="#">identification in GPT-2 small</a> . In <i>The Eleventh Inter-</i>	1571
1508	Listík, Mathias Carlen, Jan van de Kerkhof, Marcin	<i>national Conference on Learning Representations</i> .	1572
1509	Pikus, Krunoslav Zaher, Paul Müller, Sasha Zykova,	Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel,	1573
1510	Richard Stefanec, Vitaly Gatsko, Christoph Hirn-	Barret Zoph, Sebastian Borgeaud, Dani Yogatama,	1574
1511	schall, Ashwin Sethi, Xingyu Federico Xu, Chetan	Maarten Bosma, Denny Zhou, Donald Metzler, Ed H.	1575
1512	Ahuja, Beth Tsai, Anca Stefanoiu, Bo Feng, Ke-	Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy	1576
1513	shav Dhandhanania, Manish Katyal, Akshay Gupta,	Liang, Jeff Dean, and William Fedus. 2022. <a href="#">Emer-</a>	1577
1514	Atharva Parulekar, Divya Pitta, Jing Zhao, Vivaan	<a href="#">gent abilities of large language models</a> . <i>Transactions</i>	1578
1515	Bhatia, Yashodha Bhavnani, Omar Alhadlaq, Xiaolin	<i>on Machine Learning Research</i> . Survey Certifica-	1579
1516	Li, Peter Danenberg, Dennis Tu, Alex Pine, Vera	tion.	1580
1517	Filippova, Abhipso Ghosh, Ben Limonchik, Bhar-	Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio,	1581
1518	gava Urala, Chaitanya Krishna Lanka, Derik Clive,	William Cohen, Ruslan Salakhutdinov, and Christo-	1582
1519	Yi Sun, Edward Li, Hao Wu, Kevin Hongtongsak,	pher D. Manning. 2018. <a href="#">HotpotQA: A dataset for</a>	1583
1520	Ianna Li, Kalind Thakkar, Kuanysh Omarov, Kushal	<a href="#">diverse, explainable multi-hop question answering</a> .	1584
1521	Majmundar, Michael Alverson, Michael Kucharski,	In <i>Proceedings of the 2018 Conference on Empiri-</i>	1585
1522	Mohak Patel, Mudit Jain, Maksim Zabelin, Paolo	<i>cal Methods in Natural Language Processing</i> , pages	1586
1523	Pelagatti, Rohan Kohli, Saurabh Kumar, Joseph Kim,	2369–2380, Brussels, Belgium. Association for Com-	1587
1524	Swetha Sankar, Vineet Shah, Lakshmi Ramachan-	putational Linguistics.	1588
1525	druni, Xiangkai Zeng, Ben Bariach, Laura Weidinger,	Xiaodan Zhu, Peter Turney, Daniel Lemire, and André	1589
1526	Tu Vu, Amar Subramanya, Sissie Hsiao, Demis Hass-	Vellino. 2015. <a href="#">Measuring academic influence: Not</a>	1590
1527	abis, Koray Kavukcuoglu, Adam Sadovsky, Quoc Le,		
1528	Trevor Strohmman, Yonghui Wu, Slav Petrov, Jeffrey		
1529	Dean, and Oriol Vinyals. 2024. <a href="#">Gemini: A fam-</a>		
1530	<a href="#">ily of highly capable multimodal models</a> . <i>Preprint</i> ,		
1531	arXiv:2312.11805.		



all citations are equal. *Journal of the Association for Information Science and Technology*, 66(2):408–427.

## A Related work

The increasing number of IA publications during the last few years has resulted in several survey or position papers that critically discuss existing work, identify common patterns, and provide suggestions for how to go forward. [Lipton \(2018\)](#) critically question common motivations behind interpretability and the lack of definitions in the field. We follow their recommendation and provide a definition of what we consider interpretability and analysis research in §2. [Belinkov and Glass \(2019\)](#) summarize trends in early IA work and discuss recommendations for how to overcome limitations of IA research. Similar to our work, they recommend that future work should think about better ways to evaluate IA research and findings. [Rogers et al. \(2020\)](#) survey and synthesize IA work on BERTology, a subfield of IA work that focuses on encoder-only language models. [Rauker et al. \(2023\)](#) survey a large number of papers that study the internals of language models (transparency), and discuss key challenges in the field. Similar to our work, they also argue for better ways of evaluating IA methods, as well as more actionability and grounding in real-world applications. More recently, [Madsen et al. \(2024\)](#) discuss two prominent trends in interpretability research (post-hoc explanations and intrinsic interpretability) and argue that interpretability (“the study of explaining models in understandable terms to humans”) needs a new paradigm, centered around faithfulness.

Several other works study citational patterns and trends within the broader NLP community. [Mohammad \(2020\)](#) uses citations to measure the impact of NLP publications indexed by the ACL Anthology. Similar to our approach, they compare how well papers from different areas within NLP are cited, and use citation statistics to draw conclusions about the impact of different subfields within NLP. [Singh et al. \(2023\)](#) consider citations as an indicator for how widely the community is reading. They study temporal citations trends and reveal that a majority of cited papers fall within a five year time period before publication of the citing work, demonstrating a recency bias in citation behavior. [Jacovi \(2023\)](#) uses Semantic Scholar to curate a large number of papers focusing on explainability, studying citation trends in the field based on this collection. [Wahle et al. \(2023\)](#) analyze the influ-

ence between NLP and other fields over the years. Also using Semantic Scholar, they rely on citations to conclude that NLP has become more *insular* over time.

Another set of related papers surveys the NLP community for their perceptions and opinions, a method we also use. [Gururaja et al. \(2023\)](#), for example, focus on paradigm shifts and study factors that shape NLP as a field. They conduct interviews with NLP researchers and experts and gather their opinions on critical trends and patterns that emerge in the field. [Pramanick et al. \(2023\)](#) also focus on paradigm shifts and impact, but from a diachronic perspective. They provide a novel framework to study the evolution of research topics within a field to establish what drives research in NLP across time. They find that tasks and methods have a bigger impact on the field than metrics do.

Lastly, there are several related works in the scientometrics literature that study and compare the impact of research using the same metrics as we do: [Chacon et al. \(2020\)](#) apply the citation success index to compare sub-fields in physics, and [Leydesdorff \(2007\)](#) propose the use of Betweenness Centrality as a measure of the interdisciplinarity of journals.

## B Citation graph details

We provide additional details on the creation of our citation graph below.

**Summary statistics** Table 1 shows the number of papers per track in our initial collection. With 477 papers, IA is the 6th largest track in the collection.

**Standardizing submission tracks** The submission tracks of ACL and EMNLP conferences have changed considerably from 2018 to 2023. Some tracks were split into multiple tracks, some tracks appeared (and disappeared), and some were renamed. As we are mostly interested in comparing IA with other tracks, we decided to merge tracks in order to create a consistent set of tracks starting from 2020 (when the IA track was established). This unification makes our analysis more feasible. We manually assigned every track from ACL/EMNLP from 2020 to 2023 into 27 different categories:

- Information Extraction/Retrieval
- Machine Translation and Multilinguality
- Machine Learning
- Applications

Track	Paper Count
Information Extraction/Retrieval	674
Machine Translation and Multilinguality	594
Machine Learning	557
Applications	516
Dialogue	487
Interpretability and Analysis	477
Semantics	456
Resources and Evaluation	423
Multimodality, Speech and Grounding	389
Generation	361
Question Answering	334
Sentiment Analysis	258
Summarization	244
Theme	188
Social Science	178
Ethics	130
Syntax	121
Efficient Methods	113
Linguistic Theories and Psycholinguistics	106
Discourse and Pragmatics	84
Large Language Models	83
Industry	76
Phonology, Morphology and Word Segmentation	72
Commonsense Reasoning	32
Human-Centered NLP	18
Unsupervised and Weakly-Supervised Methods in NLP	17
Theory and Formalism in NLP	6

Table 1: Papers per track in ACL/EMNLP.

1691	•Dialogue
1692	•Semantics
1693	•Interpretability and Analysis
1694	•Resources and Evaluation
1695	•Generation
1696	•Question Answering
1697	•Multimodality, Speech and Grounding
1698	•Summarization
1699	•Sentiment Analysis
1700	•Theme
1701	•Social Science
1702	•Ethics
1703	•Linguistic Theories and Psycholinguistics
1704	•Syntax
1705	•Efficient Methods
1706	•Discourse and Pragmatics
1707	•Large Language Models
1708	•Phonology, Morphology and Word Segmentation
1709	•Industry
1710	•Commonsense Reasoning
1711	•Human-Centered NLP
1712	•Unsupervised and Weakly-Supervised Methods
1713	in NLP
1714	•Theory and Formalism in NLP
1715	

Statistic	Value
Nodes (papers)	185,384
Edges (citations)	786,376
Nodes originally from ACL/EMNLP 2018-2023	9,248
References from ACL/EMNLP 2018-2023 papers	374,857
Citations of ACL/EMNLP 2018-2023 papers	469,580

Table 2: Statistics of the citation graph. As some EMNLP/ACL papers cite other EMNLP/ACL papers, the total number of edges is less than the sum of the references and citations.

We note that we consider the EMNLP 2023 track: Language Modeling and Analysis of Language Models as part of IA. Additionally, we ignore papers from the theme track, as these topics change every year.

**Cleaning the collected data** Since the ACL Anthology does not provide information about the submission track, we obtain our data from a diverse set of sources as listed in Table 3. Since the data comes in very different formats, we performed the following steps to clean it.

We searched for paper titles in the ACL anthology to obtain their DOIs. As some papers were renamed, preventing us from finding the corresponding paper in the ACL Anthology, we queried the Semantic Scholar API for the closest match, with a minimum of 0.85 similarity using the Python `difflib.SequenceMatcher` class. Finally, we manually searched for the remaining papers on Semantic Scholar. After this process, we were left with only 6 papers with no Semantic Scholar ID. We exclude these from our analysis. Finally, for each paper, we queried its citations and its references using the Semantic Scholar API, and constructed the citation graph based on the results.

**Citation intent and influence** For each citation, the Semantic Scholar API provides a label of the intent (e.g. as background information, use of methods, or comparing results) (Cohan et al., 2019), and a label on whether it is a “highly influential” citation for the paper or not (Valenzuela et al., 2015). We rely on the latter label when analyzing the most cited IA papers in Section 6.

**Track classifiers details** We are interested in analyzing how papers from different tracks cite each other. However, as most of the nodes in our citation graph are papers that are not in ACL and EMNLP, we have no ground truth information for the track of these papers. Therefore, we built a classifier

Conference	Data Source
ACL 2018	Conference schedule web page
ACL 2019	Conference schedule web page
ACL 2020	Virtual conference web page
ACL 2021	Conference schedule web page
ACL 2022	Provided by the program chairs
ACL 2023	Github repository to generate webpage
EMNLP 2018	Provided by the program chairs
EMNLP 2019	Conference schedule web page
EMNLP 2020	Github repository to generate webpage
EMNLP 2021	Provided by the program chairs
EMNLP 2022	Provided by the program chairs
EMNLP 2023	Provided by the program chairs

Table 3: Data source for each conference.

to predict the track of a paper, given its title and abstract. The classifier is based on the Specter2 model (Cohan et al., 2020), which takes a title and an abstract of a paper, and outputs an embedding. We add and train a MLP layer on top of this model to obtain our classifier.

We split the data 80/20 using only papers from ACL and EMNLP from 2020 to 2023 (for which we have gold labels), and we trained the classifier for 50 epochs using Adam and a cross entropy loss. We used a learning rate of  $2 * 10^{-3}$  and a learning rate scheduler with exponential decay ( $\gamma = 0.995$ ). We perform upsampling as the number of papers in each track is imbalanced. Additionally, to get an even more diverse set of papers for the interpretability and analysis track, we augment the training data with papers accepted to the BlackboxNLP workshop, which focuses on IA work.

We find that some tracks are more difficult to predict correctly than others (e.g., Efficient Methods). We attribute this to both the limited training data and the ambiguity of submission tracks. We hence restrict ourselves to the 11 tracks (including IA) with the highest classification accuracy, and introduced an ‘Other’ category to group the remaining tracks, which we exclude from our classifier analyses. The final set of tracks in our classifier is:

- Dialogue
- Ethics
- Generation
- Information Extraction/Retrieval
- Interpretability and Analysis
- Machine Learning
- Machine Translation and Multilinguality
- Multimodality, Speech and Grounding
- Question Answering

- Social Science 1791
- Summarization 1792
- Other 1793

On this final set of tracks, our classifier achieves an F1 micro/macro score of 0.61/0.61. Given how noisy submission track labels can be (a paper can often be a plausible candidate for multiple tracks), we find our classifier’s performance to be reasonable. We additionally perform a manual error analysis and expect the classification errors made on the test set; most errors were cases where the paper could have been submitted to the predicted track.

Finally, we label the citation graph using our classifier. We used Semantic Scholar and OpenAlex (Priem et al., 2022) (in accordance with their terms of use) to obtain abstracts. 4.9% of the papers had no abstract in either source; we thus exclude these from our analysis.

## B.1 Sanity checks 1809

**Additional IA track classifier evaluations** As we are mostly interested in the performance of detecting IA papers, we validate our classifier in 2 different ways: using the IA papers suggested by our respondents in the survey, and manual annotation of 556 papers. 1810

For papers suggested by survey respondents (after removing papers included in the training data), we run our classifier and get predicted tracks. The classifier obtained an accuracy of 78.1% (82/105). Considering that these papers are out-of-domain in comparison to the training data (some are even IA papers outside of NLP), we believe this to be a good result. 1811

As for the 556 papers that were manually annotated by two authors, our classifier is 87.8% (488/556) accurate. As this data is biased towards non-IA papers (506/556 papers), we also compute precision, recall and F1 scores. The F1 score is 0.60, precision is 1.0 and recall is 0.42. Since high precision and low recall show that we underselect IA papers, we get a conservative estimate of our positive results rather than an overly generous estimate, which we find acceptable. 1812

**Correlation between betweenness centralities and citation counts** Leydesdorff (2007) find that betweenness centrality can be highly correlated to citation counts. Although this is expected (papers with more citations can also act better as *bridges*), given that BC is being used as a proxy to measure the “interdisciplinarity” of a field, we would 1813

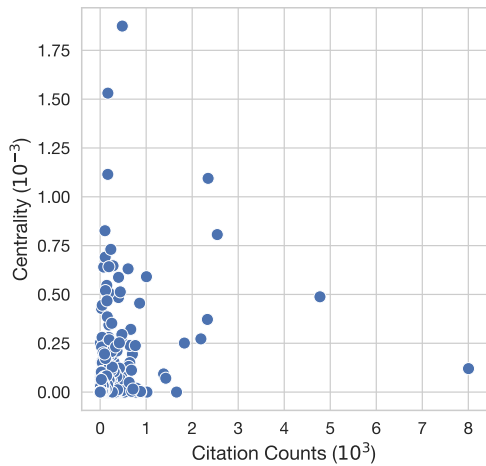


Figure 8: Betweenness centralities versus citation counts for papers in ACL and EMNLP since 2020.

want this metric to be somewhat orthogonal to the citation counts. We compute the correlation between the citation counts and the BC of all nodes in our citation graph. At 0.328 ( $p < 0.001$ ), it is considerably lower than the 0.509 reported by [Leydesdorff \(2007\)](#). Figure 8 provides a visualization of the correlation.

## C Survey details

We outline ethical considerations pertaining to our survey, along with the final version of the survey below.

### C.1 Ethical considerations

Our survey involved research with human participants, thus we report the full text of the survey below, and information about recruitment in Section 3. We determined there to be a negligible risk of harms from participating in our survey, as it contains no offensive or harmful content. As shown in the full survey below, we describe our study objectives and remind respondents that filling out the survey is completely voluntary. We then explicitly ask for their consent to participate, and obtain consent from all 138 survey respondents. For respondents who may not have completed the survey, no data was collected. In lieu of financial compensation, we offered survey respondents the optional opportunity to provide their name or an alias that we would mention in the acknowledgements of any future paper we write with the survey results. To protect respondent privacy and confidentiality, we will not release the original survey responses in full, but only release high-level statis-

tics, annotations from our qualitative coding, and select non-identifying examples in Section 7.

### C.2 Full survey

#### Impact of Model Analysis and Interpretability Research on Progress in NLP

Estimated time to complete the survey: 12 minutes

#### Study description

This project aims to measure the impact that model analysis and interpretability research has on current progress in NLP as well as its possible future impact on the field.

You are encouraged to fill out this survey even if you have no exposure to model analysis and interpretability work.

Filling out this questionnaire is completely voluntary.

By clicking "Yes" below, I am verifying that I have read the description above and I consent to participate in this research study.

- Yes
- No

#### What do we mean by model analysis and interpretability research?

Model analysis and interpretability research in natural language processing (NLP) aims to develop a deeper understanding of and explain the behavior of NLP systems.

This includes (but is not limited to) explaining models' internal computations, investigating broader phenomena observed during pre-training or adaptation, and providing a better understanding of the limitations and robustness of existing models.

Work on topics such as attribution methods, probing, mechanistic interpretability, analysis of embedding spaces, explainability, analysis of training dynamics, analyzing model bias, etc., are additional examples of model analysis and interpretability research.

#### Background questions

##### 1. What is your occupation?

- Bachelor's student
- Master's student

1921	• PhD student/candidate	value pluralism	1973
1922	• Postdoc	• <b>LMs and the world:</b> factuality, retrieval-augmented LMs, knowledge models, common-sense reasoning, theory of mind, social norms, pragmatics, and world models	1974
1923	• Assistant professor	• <b>LMs and embodiment:</b> perception, action, robotics, and multimodality	1975
1924	• Associate professor	• <b>LMs and interaction:</b> conversation, interactive learning, and multi-agents learning	1976
1925	• Full professor	• <b>LMs with tools and code:</b> integration with tools and APIs, LM-driven software engineering	1977
1926	• Junior industry researcher	• <b>LMs on diverse modalities and novel applications:</b> visual LMs, code LMs, math LMs, and so forth, with extra encouragements for less studied modalities or applications such as chemistry, medicine, education, database and beyond	1978
1927	• Senior industry researcher	• <b>NLP applications:</b> sentiment analysis, summarization, question answering, etc.	1979
1928	• NLP practitioner	• <b>Computational linguistics:</b> discourse, pragmatics, phonology, morphology, syntax, semantics	1980
1929	• Other [fill in]	• <b>Information extraction, information retrieval, text mining</b>	1981
1930		• <b>Neurosymbolic approaches</b>	1982
1931	<b>2. What is your area of research?</b>	• <b>Non-neural methods approaches for NLP</b>	1983
1932	Feel free to select multiple options or add missing ones.	• Other [fill in]	1984
1933			1985
1934	<i>(The list below is adapted from the calls for papers of COLM and ARR.)</i>		1986
1935	• <b>LM adaptation:</b> fine-tuning, instruction-tuning, reinforcement learning (with human feedback), prompt tuning, and in-context alignment		1987
1936	• <b>Data for LMs:</b> pre-training data, alignment data, and synthetic data — via manual or algorithmic analysis, curation, and generation		1988
1937	• <b>Evaluation of LMs:</b> benchmarks, simulation environments, scalable oversight, evaluation protocols and metrics, human and/or machine evaluation		1989
1938	• <b>Societal implications:</b> bias, fairness, accountability, transparency, equity, misuse, jobs, climate change, and beyond		1990
1939	• <b>Safety:</b> security, privacy, misinformation, adversarial attacks and defenses		1991
1940	• <b>Science of LMs:</b> scaling laws, fundamental limitations, emergent capabilities, demystification, interpretability, complexity, training dynamics, grokking, learning theory for LMs		1992
1941	• <b>Compute efficient LMs:</b> distillation, compression, quantization, sample efficient methods, memory efficient methods		1993
1942	• <b>Engineering for large LMs:</b> distributed training and inference on different hardware setups, training dynamics, optimization instability		1994
1943	• <b>Learning algorithms:</b> learning, unlearning, meta learning, model mixing methods, continual learning		1995
1944	• <b>Inference algorithms:</b> decoding algorithms, reasoning algorithms, search algorithms, planning algorithms		1996
1945	• <b>Human mind, brain, philosophy, laws and LMs:</b> cognitive science, neuroscience, linguistics, psycholinguistics, philosophical, or legal perspectives on LMs		1997
1946	• <b>LMs for everyone:</b> multilinguality, low-resource languages, vernacular languages, multiculturalism,		1998
1947			1999
1948		<b>[OPTIONAL]</b>	2000
1949		<b>If you would like, provide your name (or an alias) here and we will mention it in the acknowledgements of our future paper.</b> [fill in]	2001
1950			2002
1951		<b>Your take on model analysis and interpretability research</b>	2003
1952			2004
1953		<b>Reminder: What do we mean by model analysis and interpretability research?</b>	2005
1954		Model analysis and interpretability research in natural language processing (NLP) aims to develop a deeper understanding of and explain the behavior of NLP systems.	2006
1955			2007
1956			2008
1957			2009
1958			2010
1959			2011
1960			2012
1961		This includes (but is not limited to) explaining models' internal computations, investigating broader phenomena observed during pre-training or adaptation, and providing a better understanding of the limitations and robustness of existing models.	2013
1962			2014
1963			2015
1964			2016
1965			2017
1966			2018
1967			2019
1968		Work on topics such as attribution methods, probing, mechanistic interpretability, analysis of embedding spaces, explainability, analysis of training dynamics, analyzing model bias, etc., are additional examples of model analysis and	2020
1969			2021
1970			2022
1971			2023
1972			

2024	interpretability research.			2076
2025				2077
2026	<b>3. How much do you agree with the following statement?</b>			2078
2027				2079
2028	The progress in NLP in the last five years would <b>not</b>			2080
2029	<b>have been possible</b> without findings from model			2081
2030	analysis and interpretability research.			2082
2031	• 1: strongly disagree			2083
2032	• 2			2084
2033	• 3			2085
2034	• 4			2086
2035	• 5: strongly agree			2087
2036				2088
2037	<b>4. How much do you agree with the following statement?</b>			2089
2038				2090
2039	The progress in NLP in the last five years would			2091
2040	have been <b>slower</b> without findings from model			2092
2041	analysis and interpretability research.			2093
2042	• 1: strongly disagree			2094
2043	• 2			2095
2044	• 3			2096
2045	• 4			2097
2046	• 5: strongly agree			2098
2047				2099
2048	<b>5. How many model analysis and interpretability works do you read compared to other topics?</b>			2100
2049				2101
2050	• I don't usually read model analysis and inter-			2102
2051	pretability work, but I do read NLP works about			2103
2052	other topics			2104
2053	• I do read some model analysis and interpretability			2105
2054	work, but much less than other topics			2106
2055	• I read model analysis and interpretability work in			2107
2056	about the same volume as other NLP-related topics			2108
2057	• I read model analysis and interpretability work			2109
2058	more than other NLP topics			2110
2059	• Most of the works I read are about model analysis			2111
2060	and interpretability			2112
2061				2113
2062	<b>6. How, if at all, does model analysis and inter-</b>			2114
2063	<b>pretability work influence your own work?</b>			2115
2064	<input type="checkbox"/> It provides me with new research ideas			2116
2065	<input type="checkbox"/> It changes my mental model of what the			2117
2066	capabilities and limitations of models are			2118
2067	<input type="checkbox"/> It helps me ground my explanations of my own			2119
2068	results			2120
2069	<input type="checkbox"/> It adds useful tools for me to visual-			2121
2070	ize/evaluate/understand the behavior of a			2122
2071	model			2123
2072	<input type="checkbox"/> It does not influence my work			2124
2073	<input type="checkbox"/> Other [fill in]			2125
2074				2126
2075	<b>[OPTIONAL]</b>			2127
		<b>7. Provide up to 5 model analysis and inter-</b>		
		<b>pretability papers that have influenced your</b>		
		<b>work (please provide a comma separated list of</b>		
		<b>paper titles or URLs). [fill in]</b>		
		<b>8. In your day-to-day work, do you use con-</b>		
		<b>cepts from model analysis and interpretability</b>		
		<b>research (e.g., probing, residual stream, induc-</b>		
		<b>tion heads, causal interventions, MLP layers as</b>		
		<b>key-value memories, etc.)?</b>		
		• Never		
		• Rarely		
		• Sometimes		
		• Often		
		• Always		
		<b>9. Do you think model analysis and inter-</b>		
		<b>pretability research is important, and if so, why?</b>		
		<input type="checkbox"/> Understanding model limitations and capabili-		
		ties		
		<input type="checkbox"/> Making models more computationally efficient		
		<input type="checkbox"/> Developing safety mechanisms		
		<input type="checkbox"/> Improving model trustworthiness		
		<input type="checkbox"/> Explainability for users		
		<input type="checkbox"/> To fulfill legal requirements (e.g., GDPR)		
		<input type="checkbox"/> Improving model capabilities		
		<input type="checkbox"/> Developing novel architectures		
		<input type="checkbox"/> Developing novel architectures		
		<input type="checkbox"/> I do not think model analysis and interpretability		
		work is important		
		<input type="checkbox"/> Other [fill in]		
		<b>[OPTIONAL]</b>		
		<b>10. If you selected "I do not think model analysis</b>		
		<b>and interpretability research is important"</b>		
		<b>above, please elaborate why. [fill in]</b>		
		<b>[OPTIONAL]</b>		
		<b>11. In your opinion, how important is model</b>		
		<b>analysis and interpretability research to work</b>		
		<b>in the areas below?</b>		
		Work on multilinguality and low-resource lan-		
		guages		
		• Model analysis and interpretability research is		
		not important for		
		• Model analysis and interpretability research is		
		somewhat important for		
		• Model analysis and interpretability research is		
		very important for		
		Work on multimodal learning, grounding, and		

2128	embodiment	Two authors performed qualitative analysis of all	2178
2129	• Model analysis and interpretability research is	70 open-ended survey responses, and 556 papers	2179
2130	not important for	(based on their titles and abstracts).	2180
2131	• Model analysis and interpretability research is	We began by analyzing the survey responses:	2181
2132	somewhat important for	one round of independent coding was done, based	2182
2133	• Model analysis and interpretability research is	on which we reviewed our codes to normalize terms	2183
2134	very important for	and resolve disagreements. After this, a second	2184
2135		round of annotation was performed.	2185
2136	Work on engineering for large language models	As for the paper annotations, the authors did	2186
2137	• Model analysis and interpretability research is	a combination of independent coding (with dis-	2187
2138	not important for	cussion and re-coding), and co-coding. Through-	2188
2139	• Model analysis and interpretability research is	out the annotation process, the authors followed	2189
2140	somewhat important for	best practices by working closely together to clar-	2190
2141	• Model analysis and interpretability research is	ify the annotation procedure, discuss the emerging	2191
2142	very important for	themes, and re-annotate data that was coded early	2192
2143		on (Bengtsson, 2016).	2193
2144	Work on factuality, reasoning, world models	We iteratively merged codes for related themes	2194
2145	• Model analysis and interpretability research is	(e.g., <i>pre-training trajectories</i> and <i>training dynam-</i>	2195
2146	not important for	<i>ics</i> ), and to resolve inconsistencies from typos (e.g.,	2196
2147	• Model analysis and interpretability research is	<i>in-context learning</i> instead of <i>in-contex learning</i> )	2197
2148	somewhat important for	and to normalize themes (e.g., <i>interventions</i> instead	2198
2149	• Model analysis and interpretability research is	of <i>intervention</i> ), where applicable. All merging op-	2199
2150	very important for	erations are released as part of our code.	2200
2151		We measure inter-coder reliability with percent-	2201
2152	Work on societal implications, bias, misuse, and	age agreement (O'Connor and Joffe, 2020), which	2202
2153	beyond	was above 90% across all subsets of annotation.	2203
2154	• Model analysis and interpretability research is	Summary statistics are shown in Table 4.	2204
2155	not important for		
2156	• Model analysis and interpretability research is	<b>E Additional results</b>	2205
2157	somewhat important for	<b>Relative growth of submission tracks</b> Figure 9	2206
2158	• Model analysis and interpretability research is	shows the the relative growth of the IA track com-	2207
2159	very important for	pared to other tracks that have consistently existed	2208
2160		since 2020. IA is the fastest growing track at ACL	2209
2161	[OPTIONAL]	and EMNLP.	2210
2162	<b>12. In your opinion, what is missing in model</b>	<b>Betweenness centrality</b> Figure 10 shows the be-	2211
2163	<b>analysis and interpretability research right</b>	tweenness centralities for the different tracks we	2212
2164	<b>now? Where should it go in the future and how</b>	consider. We note that for this analysis we only con-	2213
2165	<b>should it be shaped differently?</b> [fill in]	sider the portion of the citation graph for which we	2214
2166		have gold track labels. Our results show that IA has	2215
2167	[OPTIONAL]	the second largest median centrality. This indicates	2216
2168	<b>13. Do you have additional opinions or thoughts</b>	that IA plays a central role in the ACL/EMNLP	2217
2169	<b>on model analysis and interpretability research?</b>	citation graph, in the sense that IA papers often lie	2218
2170	[fill in]	on the shortest path that connects to random papers	2219
2171	<b>D Qualitative coding</b>	of the graph.	2220
2172	Qualitative coding is an inductive methodology	<b>Which tracks cite IA papers</b> Figure 11 shows	2221
2173	from the social sciences (Saldana, 2021), used to	the percentage of references to IA papers across	2222
2174	systematically surface thematic patterns in data	tracks. Efficient Methods, Machine Learning, and	2223
2175	with less structure In the context of this paper,	Large Language Models cite IA papers more often	2224
2176	we use qualitative coding to analyze open-ended	than other tracks.	2225
2177	survey responses, and paper titles and abstracts.		

Data source	Instances	Themes (total)	Themes (per instance)	Agreement
Survey (what’s missing?)	42	44	2.12	91.01
Survey (why not important?)	6	9	1.5	100.00
Survey (additional thoughts)	22	29	1.95	100.00
Papers (survey)	29	59	4.28	100.00
Papers (top-50 IA)	50	115	5.38	97.03
Papers (top-50 non-IA)	50	99	4.46	96.41
Papers (non-IA papers highly influenced by IA)	456	327	4.90	97.49

Table 4: Qualitative coding statistics. For each data source, we list the total number of data instances, the total number of themes assigned, the number of themes per instance, and the percentage agreement between the codes assigned by two annotators.

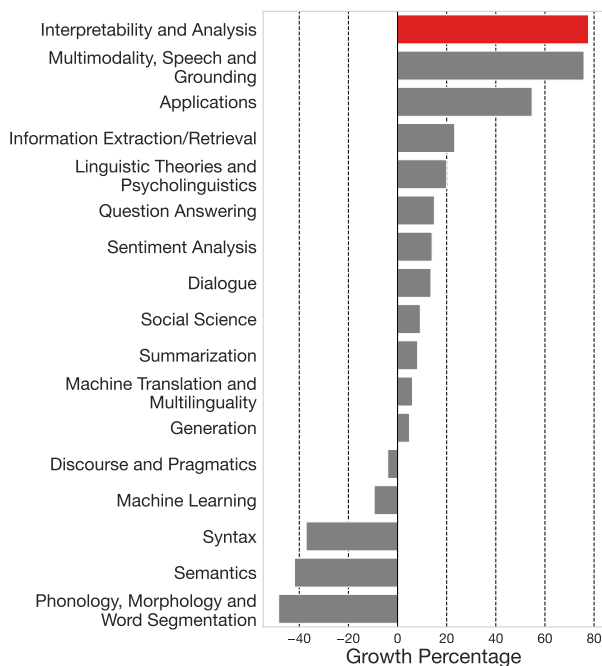


Figure 9: Growth of accepted papers per track in comparing ACL/EMNLP in 2020 vs. in 2023. This considers the tracks that have consistently existed in ACL and EMNLP in both those years.

**Citational intent** Figure 13 shows the distribution of citation intents for three groups: IA papers suggested in our survey responses, the top cited IA papers in ACL/EMNLP, and the overall most cited papers in ACL/EMNLP within our citation graph. Both the IA papers suggested in our survey and the top cited IA papers in ACL/EMNLP are primarily cited as *background information*. In contrast, the overall top cited papers in ACL/EMNLP are mostly cited for their *use of methods*.

2237  
2238  
2239  
2240  
2241  
2242  
2243  
2244  
2245  
2246

**Comparing extra-track ratios** Figure 12 compares the percentage of intra-track citations across tracks. The percentage of intra-track citations of the IA track is positioned roughly in the middle of tracks. This shows that IA is not an outlier in terms of intra-track citations.

**Top themes of highly cited IA papers** Table 5 shows the top themes that appear in (1) the papers mentioned by survey participants; (2) the top-50 most cited IA papers; (3) the top-50 most cited non-IA papers.



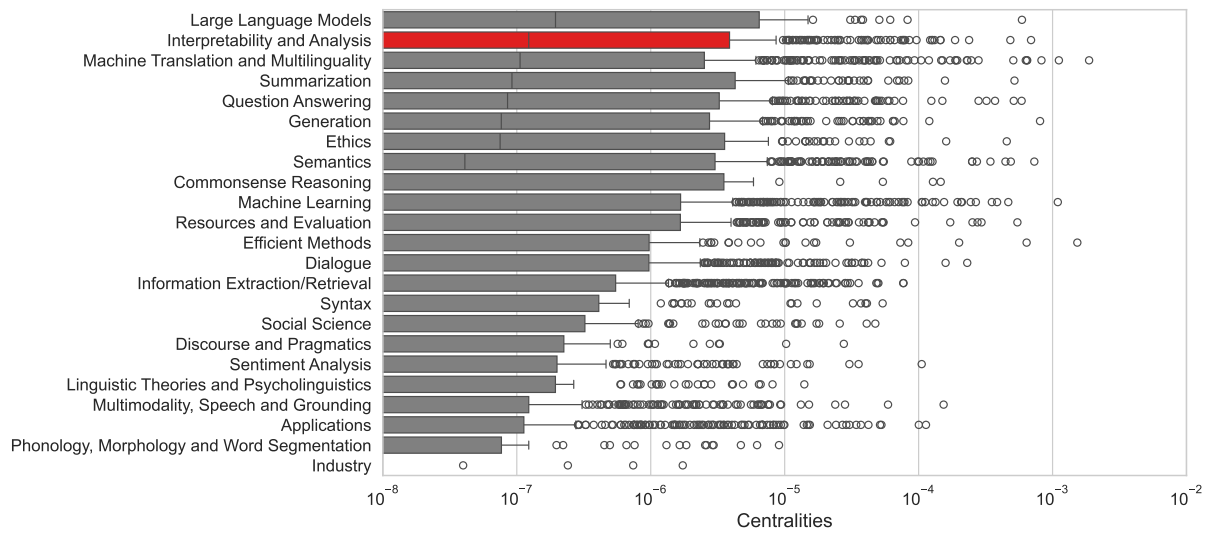


Figure 10: Betweenness centrality of ACL and EMNLP papers since 2020 by track. Lines at the middle of the box represent the medians, but some tracks have their median at 0.

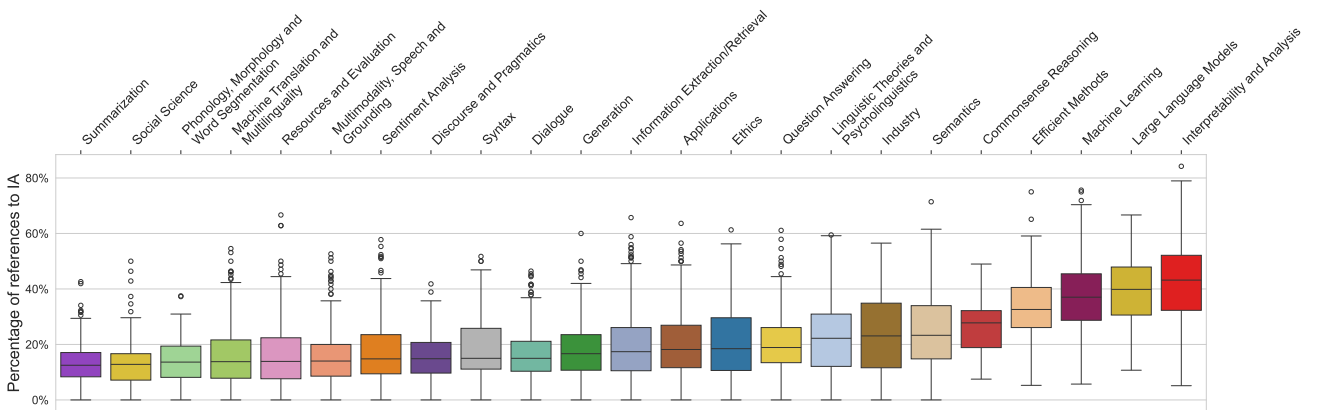


Figure 11: Percentage of references to IA papers according to our classifiers prediction.

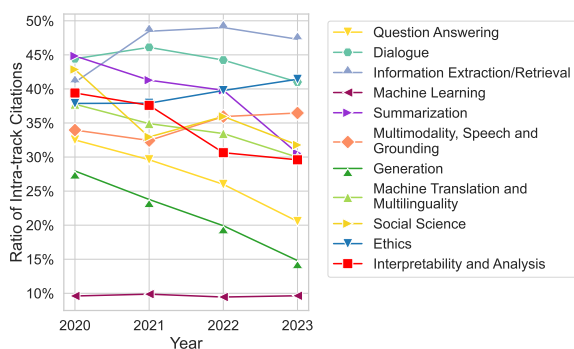


Figure 12: Ratio of intra-track citations according to the predictions of our classifier.

Source	Top themes (% of papers in which the theme appears)
Survey	representation analysis (34%), novel method (24%), probing (24%), attention analysis (21%), interventions (17.2%), mechanistic interp (17.2%), attribution (17.2%)
Top-50 IA	analysis (40%), novel method (36%), evaluation (32%), explainability (20%), linguistics (16%), probing (16%)
Top-50 non-IA	novel model (34%), novel method (32%), novel dataset (24%), analysis (16%)

Table 5: Top themes of highly influential IA papers (mentioned by survey respondents and top-50 most-cited IA papers from the citation graph), compared to the top themes of the top-50 most-cited non-IA papers. Themes are not mutually exclusive.

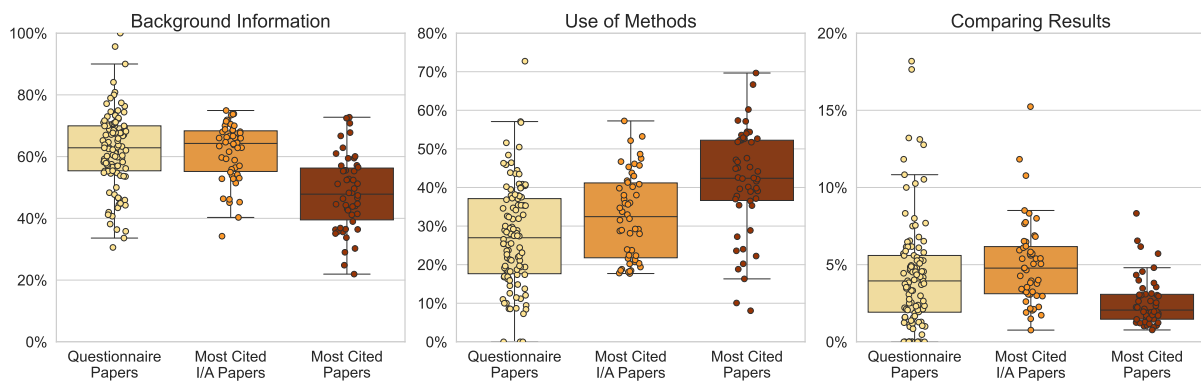


Figure 13: Citation intent percentages for the interpretability and analysis papers suggested in the responses in our survey, the top cited interpretability and analysis papers in ACL/EMNLP, and the top cited papers in ACL/EMNLP for any track.