

Towards Reliable Paper Contributions Annotation in the ACL Rolling Review

Anonymous ACL submission

Abstract

With the rapid growth of scientific publications, researchers struggle to efficiently assess the relevance of numerous papers. Identifying the types of contributions an article makes can help readers quickly grasp its significance. The ACL Rolling Review (ARR) introduced a typology requiring authors to specify their contributions to improve review quality and fairness. However, the current typology lacks clear definitions and guidance, leading to inconsistent labeling and raising concerns about its reliability. Our extensive re-annotation campaign reveals substantial disagreement between authors and domain experts. Evaluation of LLMs on paper contribution identification shows competitive performance relative to authors against annotator consensus, highlighting a potential path toward more reliable annotation.

1 Introduction

As the volume of research articles continues to grow, researchers face increasing difficulty in keeping up with the literature given limited time and resources. Staying up to date is essential for advancing research, yet evaluating the relevance and significance of numerous evolving works is time-consuming and demands substantial expertise. Rapidly assessing an article’s content and identifying key trends is crucial for keeping pace with developments in the field, and can be aided by identifying the paper contributions. By categorizing the nature of the research, researchers can more efficiently assess its relevance and impact.

Identifying the type of contributions an article makes can offer several useful applications such as *deciding what to read* by curating reading lists for learners based on pedagogical value and concept structure (Gordon et al., 2017), *setting expectations before reading* by anticipating the kind of contributions presented by a research paper (Wobbrock and Kientz, 2016; Rogers et al., 2023), *summarizing ar-*

icles using disentangled contributions from paper content (Hayashi et al., 2023; Liu et al., 2023), *automatically identifying and structuring knowledge* across different paper contributions types within a collection of articles (Auer et al., 2018; D’Souza et al., 2021), *grasping emerging trends* by classifying research contributions of a collection of articles to conduct analysis over the field (Pramanick et al., 2025; Kaltenhauser et al., 2025).

Several typologies have been proposed to categorize research contributions, but no clear consensus exists due to differences in scope and granularity (see Appendix A). The most widely adopted typology in NLP comes from the ACL Rolling Review (ARR), which requires authors to explicitly state their contributions to support clearer and fairer evaluation (Bawden, 2019; Rogers et al., 2023). In this paper, we critically examine the ARR typology, focusing on the lack of clear annotation guidelines and its consequences for annotation reliability. Through a dedicated annotation campaign, we assess the reliability of author-assigned contribution labels and explore the potential of LLMs to perform this annotation task.

In summary, our contributions are:

- We conduct a controlled re-annotation of a subset of ARR-accepted papers by domain experts, enabling direct comparison with author-assigned contribution labels.
- We quantitatively assess the reliability of ARR contribution labels through inter-annotator agreement analysis and link the main sources of disagreement to inconsistencies in the typology and the lack of clear guidelines.
- We investigate the ability of LLMs to reproduce contribution annotations and show that they perform competitively with authors, highlighting their potential to assist in annotating contributions within the ARR process.

Our code, data and model are available on [GitHub](#).

2 Data Collection

We focus on articles submitted to the ARR via OpenReview, specifically on accepted papers that include a *contribution type* following the introduction of the new guidelines in 2023¹. The ARR typology defines 11 contribution types (Rogers et al., 2023), from which authors are asked to select one or more labels that best characterize their submission: 1) *NLP engineering experiment* (most papers proposing methods to improve state-of-the-art), 2) *approaches for low-compute settings, efficiency*, 3) *approaches for low-resource settings*, 4) *data resources*, 5) *data analysis*, 6) *model analysis and interpretability*, 7) *reproduction studies*, 8) *position papers*, 9) *surveys*, 10) *theory*, and 11) *publicly available software and pre-trained models*.

For each paper, we retain only the most recent version and consider only conferences or tracks with more than ten accepted ARR articles in order to limit sampling bias. The resulting dataset contains 2,050 articles annotated with author-assigned contribution labels. We split the data into training, validation, and test sets (80–10–10) using multi-label stratification (Sechidis et al., 2011). On average, each paper is associated with 2.11 labels ($\sigma = 1.06$), highlighting the multi-faceted nature of many submissions. Figure 1 shows the label distribution across splits, revealing substantial class imbalance, largely driven by the predominance of method- and data-oriented papers. Additional distribution (§C.1) and correlation analysis (§C.2) on the dataset are provided within Appendix.

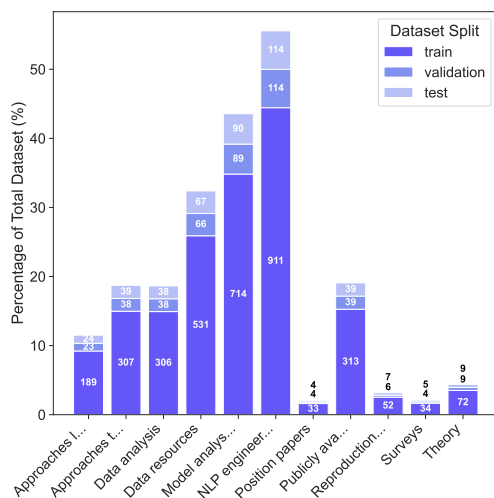


Figure 1: Contribution distribution across data splits.

¹<https://aclrollingreview.org/cfp>

3 How Reliable Are Author-Assigned Paper Contributions?

3.1 Motivations & Re-annotation Task

A major limitation of the ARR typology is the lack of precise definitions for contribution categories, which may lead to inconsistent author-assigned annotations. To assess this issue, we introduce refined definitions to support more consistent annotation, organized into a systematic framework with four complementary components:

1. A general definition of the paper contribution, intended to guide readers unfamiliar with this category of work.
2. Methodologies or techniques typically implemented in such papers.
3. Clarifying criteria that distinguish the contribution from other types in cases of ambiguity.
4. Excerpts or language patterns frequently associated with this contribution.

An example of a paper contribution with our DICE-enhanced definition is provided below. See §D.3 in appendix for extensive definitions.

Label

Approaches to low-resource settings

Description

Papers investigating scenarios where labeled data, computational resources, or linguistic tools are limited.

Implementation

These works typically employ techniques such as transfer learning, unsupervised or semi-supervised learning, or data augmentation to address these limitations.

Clarification

Submissions focusing on well-resourced settings or small improvements that don't address major resource limitations are excluded.

Examples

"Given constraints on computing resources in our deployment environment, we fine-tuned a distilled model to perform efficient and accurate intent classification.", ...

3.2 Annotation Framework

Using the refined contribution definitions, four domain experts (1 senior and 3 junior researchers) annotated the 207-document test split. Given the time-consuming nature of the task (experts required an average of 8 minutes per document), expert annotations are directly compared against author labels. Since papers are associated with few labels on average (2.11, $\sigma = 1.06$), we consider a Jaccard index below $2/3$ as indicative of meaningful disagreement, triggering the assignment of an additional expert annotator. As a result, 72% of the samples were double-annotated. Additional details

and the full annotation guidelines are provided in Appendix D.

3.3 Authors Reliability Analysis

We assess the reliability of author-provided annotations by comparing them to expert annotations, treated as the reference, using Krippendorff’s α (Krippendorff, 2004; Castro, 2017). On the test set, experts achieve an agreement of $\alpha = 0.55$, with substantial variability in pairwise agreement, reflecting the inherent complexity of the task. Agreement between authors and experts is lower but remains moderate ($\alpha = 0.47$; see Figure 2).

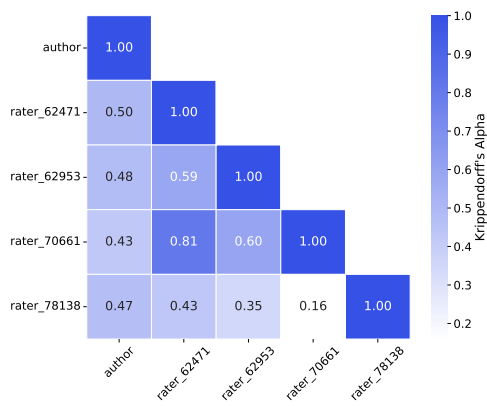


Figure 2: Pairwise Krippendorff’s α between authors and domain experts on the test set.

As shown in Figure 3, author-assigned labels that are not selected by any expert occur frequently across categories, particularly for meta-level contributions. These discrepancies are also reflected in overall label proportions (Appendix E), with experts and authors assigning certain labels at very different rates. Experts tend to prefer generic labels and use meta-paper labels more conservatively, indicating differing views on the type or role of the annotated articles. Notably, the *Reproduction Study* label is never selected by experts, highlighting a clear mismatch between authors and experts understanding of this category.

These findings suggest that this aspect of the ARR process leads to inconsistent and unreliable labels, reducing their usefulness for reviewers.

4 Can Language Models Serve as Reliable Paper Contributions Annotators?

4.1 Identification of Paper Contributions

We formalize paper contribution identification as a multi-label classification task. Given a document

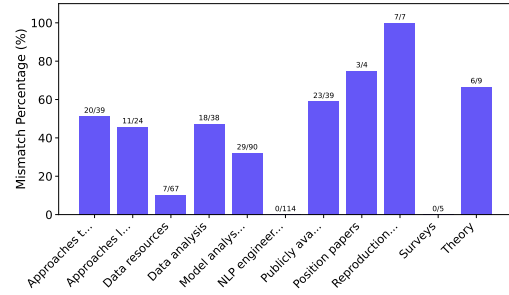


Figure 3: Percentage of author labels not confirmed by experts on test split.

D and a set of paper contributions $C = c_1, \dots, c_n$, the goal is to learn a function $f : D \rightarrow L$, that maps D to one or more applicable labels $L \subseteq C$.

4.1.1 Experimental Setting

We evaluate the ability of generative models to infer paper contributions across a collection of articles. Specifically, we evaluate LLMs such as LLaMA (3.2–3B), Mistral (7B), and GPT-4.1, as well as fine-tuned general-domain Pre-trained Language Models (PLMs) (BERT, RoBERTa), fine-tuned science-focused PLMs (SciBERT, SPECTER2), and TF-IDF as baseline. Detailed models configurations (§F.1), instructions (§F.2) and prompting procedures (§F.3) are provided in Appendix.

In real-world scenarios, articles are often behind paywalls, leaving only metadata accessible. We therefore represent each article using the concatenation of its title and abstract. LLMs are evaluated under both author and expert conditions, using either labels only or the same annotation guidelines and DICE-enhanced typology employed by expert annotators. We further vary the number of few-shot examples to evaluate the models’ ability to generalize from prior annotations, simulating an annotator refining decisions based on previously seen documents. The ground truth for the test set consists of labels on which a majority of domain experts and the authors agree. Model performance is evaluated using micro-averaged F1, reflecting overall labeling accuracy, and macro-averaged F1, to account for class imbalance and assess performance on under-represented labels.

4.2 Experimental Results

Results are presented in Table 1. Overall, the best-performing model is GPT-4.1 using 3-shot prompting. When evaluated against annotators consensus, this model achieves performance comparable to

author annotations in terms of micro-averaged F1.

TF-IDF performance is promising, achieving notably high $F1_{\text{micro}}$ scores, which indicates that lexical cues in the documents strongly support contribution identification. However, relatively low $F1_{\text{macro}}$ scores suggest that certain label types cannot be reliably inferred from lexical information alone. Similar trends are observed for PLMs on $F1_{\text{micro}}$, although $F1_{\text{macro}}$ scores are generally higher with fine-tuned domain-specific models, suggesting that their contextualized embeddings improve detection of underrepresented labels while maintaining stability across models. LLMs show more variable performance across settings. The best configurations of each model generally outperform both baselines and encoder-based approaches, suggesting that effective generalization from contribution definitions is more beneficial than relying solely on labeled training data. Providing models with the refined definitions used by expert annotators (DICE) generally improves alignment with expert labels, particularly when combined with few-shot sampling. Interestingly, for GPT-4.1, few-shot prompting largely mitigates the need for precise definitions, an effect not observed for other models.

Model	Annotators Consensus		
	$F1_{\text{micro}}$	$F1_{\text{macro}}$	
Authors	0.76	0.69	
TF-IDF	0.62	0.27	
BERT Base <small>Uncased</small>	0.64	0.30	
RoBERTa Base	0.64	0.39	
SciBERT <small>SciVocab, Uncased</small>	0.66	0.40	
SPECTER2 <small>Base</small>	0.65	0.36	
LLaMA 3.2 <small>3B Instruct</small>	0-Shot	0.48	0.33
	1-Shot	0.47	0.29
	3-Shot	0.50	0.30
	0-Shot & DICE	0.48	0.32
	1-Shot & DICE	0.65	0.37
	3-Shot & DICE	0.45	0.25
Mistral 7B Instruct <small>v0.3</small>	0-Shot	0.55	0.36
	1-Shot	0.54	0.33
	3-Shot	0.61	0.36
	0-Shot & DICE	0.55	0.36
	1-Shot & DICE	0.64	0.42
	3-Shot & DICE	0.69	0.49
GPT-4.1 <small>2025-04-14</small>	0-Shot	0.67	0.54
	1-Shot	0.70	0.54
	3-Shot	0.74	0.56
	0-Shot & DICE	0.70	0.55
	1-Shot & DICE	0.72	0.55
	3-Shot & DICE	0.72	0.56

Table 1: Paper contribution identification performance (averaged over five seeds ; full results in Appendix G)

GPT-4.1 is the highest overall performing model, yet its per-label F1 remains lower than authors for several contribution types (Figure 4). Conversely, models outperform authors on a few labels (e.g., *NLP engineering experiment*, *approaches for low-*

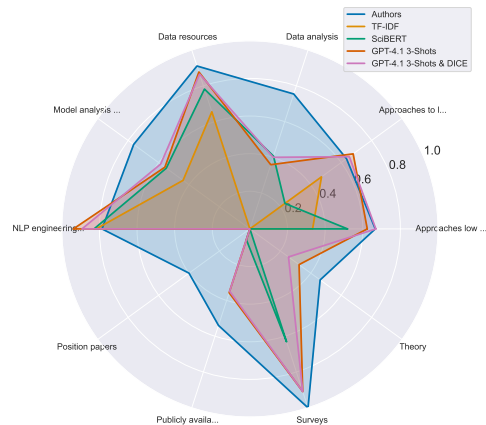


Figure 4: Per-label F1 performance of authors and best-performing AI systems on annotators consensus.

resource settings), indicating that while LLMs can surpass authors on generic categories, some contribution labels remain challenging to automatically identify even with improved guidelines.

These results suggest that, while still below human-level performance, dedicated models could assist authors in annotating contributions within the ARR process. Moreover, the performance gains observed with DICE highlight the value of refined contribution definitions in guiding annotation.

5 Conclusion

Paper contributions, as currently defined and self-assigned, show limited reliability in the ARR process due to inconsistent author labeling. Re-annotation with improved guidelines shows discrepancies between authors and domain experts, highlighting unreliability of current ARR definitions. Our evaluation shows that LLMs are promising annotators: their predictions are relatively close to authors, though they do not yet reach human-level accuracy on all labels. Combining refined typologies with LLM assistance could make contribution annotation more consistent, transparent, and scalable, especially in rapidly evolving fields like NLP. To promote broader adoption of best practices and enhance authors and reviewers perception of paper contributions, we release refined definitions along with a pretrained SciBERT model intended to provide guidance for annotating paper contributions. Moreover, automatically extracting contributions can benefit not only ARR but also a wide range of downstream tasks, paving the way for more structured and accessible scientific knowledge.

281 Limitations

282 **Enhancing ARR Typology Definitions.** We pro- 329
283 posed refinements to the ARR typology definitions 330
284 using the DICE framework (Description, Imple- 331
285 mentation, Clarification, Examples), which is de- 332
286 signed to support maintainability in evolving re- 333
287 search fields. These refined definitions are in- 334
288 formed by typology interpretations derived from 335
289 author-annotated articles and feedback from expert 336
290 annotators. While the present effort represents an 337
291 initial step toward more structured and explicit de- 338
292 finitions, it is possible that ARR editors emphasize 339
293 aspects of the typology that differ from our inter- 340
294 pretation of the papers contributions, which may in 341
295 turn influence agreement with the proposed guide- 342
296 lines. The DICE-based definitions are therefore 343
297 intended to remain adaptable and may be further re- 344
298 fined through more extensive analyses and curated 345
299 to better align with the needs of ARR organizers. 346

300 **Annotation Campaign Scope.** Re-annotation ef- 347
301 forts were focused on the test set, reflecting practi- 348
302 cal constraints related to annotation time (approx- 349
303 imately 8 minutes per document on average). To 350
304 strengthen annotation quality, 72% of the data split 351
305 was double-annotated by experts, and the original 352
306 authors were incorporated as additional annotators 353
307 when computing consensus. While the resulting 354
308 typology analysis is conducted within this setting, 355
309 extending annotation to the full dataset with a larger 356
310 pool of annotators would be a natural direction for 357
311 future work to further consolidate these findings.

312 **Processed Documents.** In these experiments, the 358
313 identification of paper contributions is performed 359
314 using article titles and abstracts. This choice re- 360
315 flects both the context length limitations of PLMs 361
316 and the objective of evaluating models under com- 362
317 parable conditions. Nonetheless, restricting the 363
318 input to partial document content may result in 364
319 the omission of relevant contribution information, 365
320 as annotators reported typically relying not only 366
321 on the title and abstract but also prioritizing the 367
322 introduction and conclusion sections, and more 368
323 generally consulting the full article when assessing 369
324 the scope of a paper’s contributions. Although the 370
325 dataset includes full-text content extracted using 371
326 GROBID², this information was not exploited in 372
327 the present study. 373

²<https://github.com/kermitt2/grobid>

Ethical Considerations 328

Annotation Campaign. In this study, we con- 329
ducted internal annotation campaigns with a team 330
of seven NLP experts, including three researchers 331
and four NLP PhD students. Two of the four anno- 332
tators selected for the test set annotations are also 333
co-authors of this paper. Given the nature of the 334
task, the annotation campaign was carried out over 335
a two-week period, with annotators spending an 336
average of approximately 8 minutes assessing the 337
contributions of a single paper. Participants were 338
compensated through non-financial means. 339

Use of generative AI in Scientific NLP. The use 340
of generative AI models for assessing paper contri- 341
butions entails inherent risks associated with this 342
technology. In scholarly contexts, such systems 343
must account for potential false positives and false 344
negatives, overlooked contributions, and inaccurate 345
representations of article content. We therefore 346
advocate for a measured and responsible use of 347
these tools and emphasize that, while they are in- 348
tended to support authors and annotators during 349
the annotation process, they are not designed to re- 350
place or fully automate this component of the ARR 351
workflow. 352

Use of LLMs for the making of this study. In this 353
study, LLMs were used only for writing style, code 354
refactoring, and generating a few generic DICE 355
excerpts. All conceptual contributions, analyses, 356
and experiments were done without AI. 357

References 358

- Liz Allen, Alison O’Connell, and Veronique Kiermer. 359
2019. [How can we ensure visibility and diversity 360
in research contributions? how the contributor role 361
taxonomy \(credit\) is helping the shift from authorship 362
to contributorship.](#) *Learned Publishing*, 32(1):71–74. 363
- Sören Auer, Viktor Kovtun, Manuel Prinz, Anna 364
Kasprzik, Markus Stocker, and Maria Esther Vidal. 365
2018. [Towards a knowledge graph for science.](#) In 366
*Proceedings of the 8th International Conference on 367
Web Intelligence, Mining and Semantics, WIMS ’18,* 368
New York, NY, USA. Association for Computing 369
Machinery. 370
- Rachel Bawden. 2019. [One Paper, Nine Reviews.](#) 371
- Emily M. Bender and Leon Derczynski. 2018. [COLING 372
2018: Paper Types.](#) Accessed: 2025-01-29. 373
- Jordan Boyd-Graber, Naoaki Okazaki, and Anna Rogers. 374
2023. [Paper-reviewer matching at ACL 2023: types 375
of contributions and track sub-areas .](#) 376

377	Amy Brand, Liz Allen, Micah Altman, Marjorie Hlava, and Jo Scott. 2015. Beyond authorship: attribution, contribution, collaboration, and credit . <i>Learned Publishing</i> , 28(2):151–155.	433
378		434
379		435
380		436
381	Santiago Castro. 2017. Fast Krippendorff: Fast computation of Krippendorff’s alpha agreement measure . https://github.com/pln-fing-udelar/fast-krippendorff .	437
382		438
383		439
384		440
385	Haihua Chen, Huyen Nguyen, and Asmaa Alghamdi. 2022. Constructing a high-quality dataset for automated creation of summaries of fundamental contributions of research articles . <i>Scientometrics</i> , 127(12):7061–7075.	441
386		442
387		443
388		444
389		445
390	Liyue Chen, Jielan Ding, Donghuan Song, and Zihao Qu. 2025. Exploring scientific contributions through citation context and division of labor . <i>Scientometrics</i> , 130(5):2901–2921.	446
391		447
392		448
393		449
394	Jennifer D’Souza, Sören Auer, and Ted Pedersen. 2021. SemEval-2021 task 11: NLPContributionGraph - structuring scholarly NLP contributions for a research knowledge graph . In <i>Proceedings of the 15th International Workshop on Semantic Evaluation (SemEval-2021)</i> , pages 364–376, Online. Association for Computational Linguistics.	450
395		451
396		452
397		453
398		454
399		455
400		456
401	Jennifer D’Souza and Sören Auer. 2021. Sentence, Phrase, and Triple Annotations to Build a Knowledge Graph of Natural Language Processing Contributions—A Trial Dataset . <i>Journal of Data and Information Science</i> , 6(3):6–34.	457
402		458
403		459
404		460
405		461
406	Alexander Fabbri, Irene Li, Prawat Trairatvorakul, Yijiao He, Weitai Ting, Robert Tung, Caitlin Westfield, and Dragomir Radev. 2018. TutorialBank: A manually-collected corpus for prerequisite chains, survey extraction and resource recommendation . In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 611–620, Melbourne, Australia. Association for Computational Linguistics.	462
407		463
408		464
409		465
410		466
411		467
412		468
413		469
414		470
415	Jonathan Gordon, Stephen Aguilar, Emily Sheng, and Gully Burns. 2017. Structured generation of technical reading lists . In <i>Proceedings of the 12th Workshop on Innovative Use of NLP for Building Educational Applications</i> , pages 261–270, Copenhagen, Denmark. Association for Computational Linguistics.	471
416		472
417		473
418		474
419		475
420		476
421		477
422	Komal Gupta, Ammaar Ahmad, Tirthankar Ghosal, and Asif Ekbal. 2021. ContriSci: A BERT-Based Multitasking Deep Neural Architecture to Identify Contribution Statements from Research Papers , page 436–452. Springer International Publishing.	478
423		479
424		480
425		481
426		482
427	Komal Gupta, Ammaar Ahmad, Tirthankar Ghosal, and Asif Ekbal. 2024. A bert-based sequential deep neural architecture to identify contribution statements and extract phrases for triplets from scientific publications . <i>International Journal on Digital Libraries</i> , 25(4):1–28.	483
428		484
429		485
430		486
431		487
432		488
	Carolin Haeussler and Henry Sauer mann. 2020. Division of labor in collaborative knowledge production: The role of team size and interdisciplinarity . <i>Research Policy</i> , 49(6):103987.	489
		490
	Hiroaki Hayashi, Wojciech Kryscinski, Bryan McCann, Nazneen Rajani, and Caiming Xiong. 2023. What’s new? summarizing contributions in scientific literature . In <i>Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics</i> , pages 1019–1031, Dubrovnik, Croatia. Association for Computational Linguistics.	491
		492
	Annika Kaltenhauser, Gian-Luca Savino, Nick von Felten, and Johannes Schöning. 2025. CHI’s Greatest Hits: Analyzing the 100 Most-Cited Papers in 43 Years of Research at ACM CHI . <i>Interactions</i> , 32(1):28–33.	493
		494
	Klaus Krippendorff. 2004. Reliability in content analysis . <i>Human Communication Research</i> , 30(3):411–433.	495
		496
	Meng-Huan Liu, An-Zi Yen, Hen-Hsen Huang, and Hsin-Hsi Chen. 2023. Contributionsum: Generating disentangled contributions for scientific papers . In <i>Proceedings of the 32nd ACM International Conference on Information and Knowledge Management, CIKM ’23</i> , page 5351–5355, New York, NY, USA. Association for Computing Machinery.	497
		498
	Aniket Pramanick, Yufang Hou, Saif M. Mohammad, and Iryna Gurevych. 2025. The nature of NLP: Analyzing contributions in NLP papers . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 25169–25191, Vienna, Austria. Association for Computational Linguistics.	499
		500
	Anna Rogers, Marzena Karpinska, Jordan Boyd-Graber, and Naoaki Okazaki. 2023. Program chairs’ report on peer review at acl 2023 . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages xl–lxxv, Toronto, Canada. Association for Computational Linguistics.	501
		502
	Konstantinos Sechidis, Grigorios Tsoumakas, and Ioannis Vlahavas. 2011. On the stratification of multi-label data . In <i>Machine Learning and Knowledge Discovery in Databases</i> , pages 145–158, Berlin, Heidelberg. Springer Berlin Heidelberg.	503
		504
	Emily Sheng, Prem Natarajan, Jonathan Gordon, and Gully Burns. 2017. An investigation into the pedagogical features of documents . In <i>Proceedings of the 12th Workshop on Innovative Use of NLP for Building Educational Applications</i> , pages 109–120, Copenhagen, Denmark. Association for Computational Linguistics.	505
		506
	Jacob O. Wobbrock and Julie A. Kientz. 2016. Research contributions in human-computer interaction . <i>Interactions</i> , 23(3):38–44.	507
		508

A Existing Contribution Typologies

ARR Typology	
Rogers et al. (2023)	Approaches to low-resource settings , Approaches low compute settings-efficiency , Data resources , Data analysis , Model analysis & interpretability , NLP engineering experiment , Publicly available software and/or pre-trained models , Position papers , Reproduction study , Surveys , Theory
ARR Inspirations	
Bender and Derczynski (2018)	Computationally-aided linguistic analysis , NLP engineering experiment paper , Reproduction paper , Resource paper , Position paper , Survey Paper
Boyd-Graber et al. (2023)	Computationally-aided linguistic analysis , NLP engineering experiment , Approaches for data- and compute efficiency , Reproduction study , New data resources , Position papers , Surveys , Theory , Publicly available software and pre-trained models
Other Contributions Typologies	
D'Souza et al. (2021)	Research problem , Approach , Model , Code , Dataset , Experimental setup , Hyper parameters , Baselines , Results , Tasks , Experiments , Ablation analysis .
Chen et al. (2022)	Dataset/Resources creation , Theory proposal , Model construction or optimization , Algorithms/Methods construction or optimization , Performance evaluation , Applications
Liu et al. (2023)	Approach , Analysis , Result , Topic or Resource
Pramanick et al. (2025)	Knowledge (k-Dataset , k-Language , k-Method , k-People , k-Task) , Artifact (a-Dataset , a-Method , a-Task)
Chen et al. (2025)	Theoretical , Methodological , Experimental , Data-based , Other
Other Non-Contributions Typologies	
Brand et al. (2015)	Conceptualization , Methodology , Software , Validation , Formal Analysis , Investigation , Resources , Data curation , Writing – Original Draft , Writing – Review & Editing , Visualization , Supervision , Project Administration , Funding acquisition
Sheng et al. (2017)	Survey , Tutorial , Resource , Reference Work , Empirical Results , Software Manual , Other

Table 2: Colors indicate generic concepts from the ARR typology, as well as similar labels found in other typologies based on the definition provided by their authors. The highlighted categories include: Optimization , Resources , Analysis , Experimental , Models/Softwares , Position paper , Reproduction study , Survey , and Theory . Please note that due to differences between typologies, the highlighted concepts may not exactly match the original labels in every case.

489 **B Related Works**

490 Contributions are commonly defined in the literature as novel scientific advancements attributed to the
491 authors of a research paper (Pramanick et al., 2025). The identification of such contributions is typically
492 performed at the document or statement level. However, the notion of what constitutes a "contribution"
493 frequently extends beyond disciplinary boundaries, presenting similarity with other fields (Sheng et al.,
494 2017; Brand et al., 2015).

495 **Paper Contributions.** Paper contributions identification aims at capturing an holistic view of scope of a
496 paper. This task has been particularly explored in approaches designed to support reviewers during the
497 submission process, helping them more accurately assess the nature of the work they evaluate (Bender and
498 Derczynski, 2018; Bawden, 2019; Boyd-Graber et al., 2023; Rogers et al., 2023). Paper-level identification
499 schemes have also been used to analyze submissions to top-tier conferences, revealing research trends and
500 suggesting ways to adapt to evolving fields (Rogers et al., 2023; Kaltenhauser et al., 2025).

501 **Contributions Statements.** This line of research focuses on identifying individual contribution statements
502 within scientific papers, providing fine-grained insights into the specific achievements of a scientific work.
503 Several works have explored the extraction of contribution-specific data in order to build knowledge
504 graphs, sometime directly from the statement (Gupta et al., 2021, 2024) or combined with contribution
505 types (D'Souza and Auer, 2021; D'Souza et al., 2021). Other studies have leveraged contribution
506 statements to summarize the key contributions of research articles, enabling readers to quickly grasp their
507 main findings (Chen et al., 2022; Hayashi et al., 2023; Liu et al., 2023). More recently, Pramanick et al.
508 (2025) identified and analyzed contribution statements across ACL Anthology papers to find research
509 trends and characterize the nature of NLP research.

510 **Pedagogical Roles of Papers.** The identification of pedagogical roles is closely related to contribution
511 analysis, as it seeks to characterize how useful a document is for individuals aiming to learn specific
512 concepts (Sheng et al., 2017). Although pedagogical roles do not directly represent a paper's contributions,
513 they are often closely aligned with them in existing typologies (see Table 2). These roles have been notably
514 exploited by approaches that model relationships between papers to support resource recommendation,
515 whether focusing on conceptual links (Gordon et al., 2017) or prerequisite chains (Fabbri et al., 2018).

516 **Contributor Roles.** In Science Studies and related fields, researchers have explored the Division of Labor
517 (DOL) in scientific collaboration, analyzing how responsibilities are distributed among contributors (Allen
518 et al., 2019; Haeussler and Sauermann, 2020). Taxonomies such as the Contributor Role Taxonomy
519 (CRediT) (Brand et al., 2015) allow to model the roles of scientists during production and publication of
520 research articles. Interestingly, this taxonomy presented relationships between paper contributions and
521 DOL, as shown by Chen et al. (2025). Their work leveraged contribution statements and citation contexts
522 to examine whether the contributions reported in papers align with authors' self-perceived roles.

C Contributions Analysis

523

C.1 Distribution of Paper Contributions

524

Figure 1 shows the distribution of paper contributions labels across the collection, revealing a strong imbalance in the labels selected by authors. Some labels account for more than 30% of the dataset (e.g., NLP Engineering Experiment, Model Analysis and Interpretability, Data Resources), while others represent less than 5% of the instances (e.g., Theory, Reproduction papers, Surveys, Positions papers). These discrepancies between labels could be explained by multiple factors. First, the ARR typology is designed to make contributions more explicit to reviewers, helping them better recognize and fairly evaluate papers with less common types of contributions (Rogers et al., 2023). Second, certain contribution types may be incidental to others; for example, an experimental paper may include a model analysis, or a data resource paper may present an accompanying dataset analysis. Finally, the general lack of clear definitions or annotation guidelines within the typology may lead authors to favor broader, more generic categories, or to select all labels that appear even loosely relevant. This imbalance motivates the need for clearer contribution definitions and points to possible biases in models trained on this dataset.

525

526

527

528

529

530

531

532

533

534

535

536

C.2 Correlation Analysis of Contributions

537

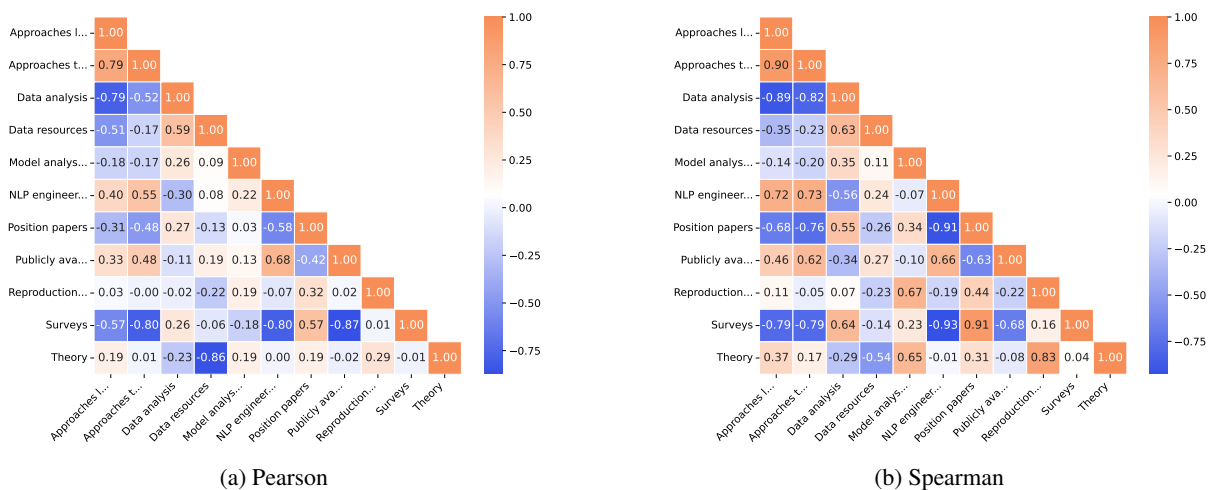


Figure 5: Correlation analysis of paper contribution labels.

We compute correlation matrices over the co-occurrence of contribution labels to identify potential relationships between paper contribution types. Due to the imbalance in label distribution, we normalize the co-occurrence matrix using Pointwise Mutual Information (PMI). Both Pearson’s r and Spearman’s ρ correlation coefficients are computed, showing similar trends (see Figure 5).

538

539

540

541

The results reveal strong associations between certain contribution types: papers focusing on low-resource settings and efficiency are highly correlated ($\rho = 0.91$); surveys and position papers also co-occur frequently, as do theory papers with reproduction studies ($\rho > 0.80$). Conversely, and perhaps unsurprisingly, NLP engineering experiment papers exhibit strong negative correlations ($\rho < -0.90$) with both position and survey papers. Papers providing publicly available models show notable variation, being most positively correlated with experimental papers ($r = 0.68$) and most negatively correlated with survey papers ($r = -0.87$).

542

543

544

545

546

547

548

Overall, the observed correlations indicate that the current typology may benefit from refinement, as ideally each contribution type should capture a distinct facet of a paper with minimal overlap.

549

550

D Annotation Campaign

D.1 Pilot Annotation

Before moving to the main annotation campaign, we conduct a preliminary annotation campaign with seven annotators (three senior and four junior researchers, all specializing in NLP) aimed at refining guidelines and identifying the experts with the highest agreement. This was conducted on 65 full papers, deliberately oversampling underrepresented classes to ensure comprehensive coverage of the typology. The distribution of this pilot is as follows:

- Approaches low compute settings-efficiency: 15.38%
- Approaches to low-resource settings: 23.08%
- Data analysis: 24.62%
- Data resources: 26.15%
- Model analysis and interpretability: 44.62%
- NLP engineering experiment: 55.38%
- Position papers: 6.15%
- Publicly available software and/or pre-trained models: 13.85%
- Reproduction study: 9.23%
- Surveys: 6.15%
- Theory: 15.38%

Each paper is annotated by at least three annotators, with up to two additional annotators added when agreement, measured by the Jaccard index, fell below 30%. Feedback from this campaign is used to refine and consolidate the typology definitions through direct exchanges with annotators, addressing the difficulties they encounter throughout the annotation process. The *Clarification* and *Examples* components of the DICE framework notably emerged from this initial annotation phase. Agreement of this pilot annotations is presented in Table 6, showing Krippendorff’s α among domain-experts. We selected annotators with higher agreement scores in this preliminary campaign for the final annotation campaign.

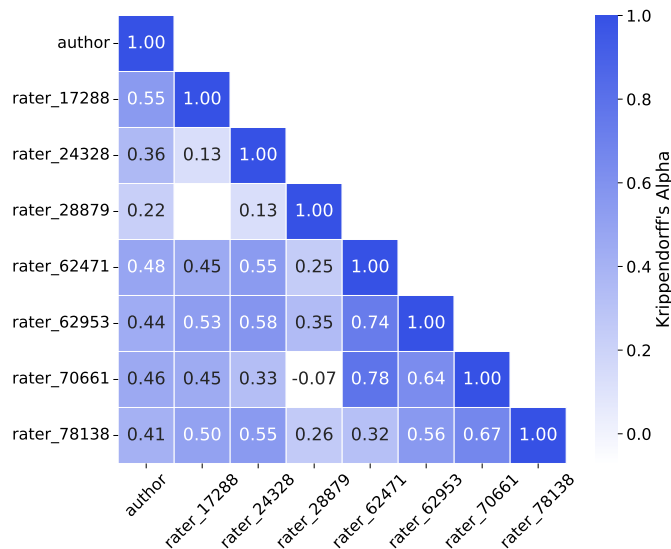


Figure 6: Pairwise Krippendorff’s α between authors and domain experts on the test set within pilot campaign.

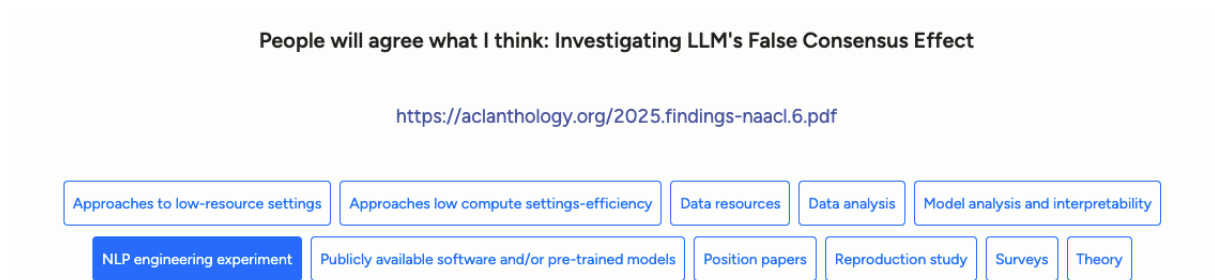


Figure 7: Annotation campaign interface

D.3 Annotation Guidelines

# Guidelines	578
You will be provided with a scientific paper and a typology of potential contribution types.	579
Analyze the paper and classify the type(s) of contributions it makes based on the provided typology.	580
	581
	582
	583
	584
Follow these rules when making your decisions:	585
	586
1. Assign applicable labels. A paper may fall under multiple contribution types, assign those that are clearly present in the paper content.	587
2. Focus on notable contributions. Only assign label if the paper makes a substantial and deliberate contribution of that type, not just a minor mention or incidental inclusion.	588
	589
	590
	591
# Typology	592
	593
Approaches to low-resource settings	594
Papers investigating scenarios where labeled data, computational resources, or linguistic tools are limited.	595
These works typically employ techniques such as transfer learning, unsupervised or semi-supervised learning, or data augmentation to address these limitations. Submissions focusing on well-resourced settings or small improvements that don't address major resource limitations are excluded.	596
	597
	598
Examples:	599
-"Given constraints on computing resources in our deployment environment, we fine-tuned a distilled model to perform efficient and accurate intent classification.	600
-"We propose an adversarial representation alignment model to mitigate performance degradation in low-resource parsing by selectively transferring relevant knowledge from high-resource languages.	601
-"To address limited annotation budgets, we trained a weakly supervised classifier using automatically generated noisy labels based on keyword heuristics."	602
	603
	604
	605
	606
	607
Approaches low compute settings-efficiency	608
Papers focusing on computational efficiency in NLP environments. These works typically employ techniques such as reduced memory usage, optimized training or inference time, or ways to reduce energy consumption, making models more accessible or deployable in low-compute environments. Submissions focusing on minor efficiency gains in high-resource contexts or lacking practical impact on accessibility are excluded.	609
	610
	611
	612
	613
	614
Examples:	615
-"To enable deployment on mobile devices, we reduced model size by pruning and quantization, significantly lowering memory usage without sacrificing accuracy.	616
-"We improve WER on long-form ASR and achieve up to 20x faster inference through batched parallel decoding.	617
-"We propose an optimized training schedule that cuts down GPU hours by 40%, making large-scale NLP model training more feasible for smaller labs."	618
	619
	620
	621
Data resources	622
Papers providing new language-related resources such as datasets, annotated corpora, annotation standards or evaluation benchmarks. These works typically provide detailed documentation of the resource creation process, including methodology and quality assurance. Submissions lacking substantial novelty in resource creation or failing to share them publicly are excluded.	623
	624
	625
	626
	627
Examples:	628
-"We present a corpus of historical legal texts annotated for named entities and syntactic structures, along with an open-access tool for browsing and querying the data.	629
-"This work introduces a benchmark, including a standardized evaluation protocol and carefully curated test sets.	630
-"We propose a taxonomy and release a dataset of ~2,000 annotated NLP paper abstracts, capturing and categorizing their scientific contributions."	631
	632
	633
	634
	635
Data analysis	636
Papers conducting detailed analysis of data resources. These works generally focus on annotation quality, bias, linguistic patterns, or how models interact with data. Submissions should present novel insights that contribute to better data practices or enhanced model design.	637
	638
	639
	640

641 Examples:

- 642 -"This work propose a comprehensive analysis of gender bias in a widely used sentiment analysis dataset,
643 revealing systematic annotation inconsistencies that affect model predictions.
644 -"We analyze the existing data resources and identify areas for improvement in future iterations.
645 -"Model performances across subsets of a dependency parsing corpus show that annotation errors
646 disproportionately impact low-frequency syntactic constructions."

647 Model analysis & interpretability

648 Papers investigating the internal mechanisms or external behavior of NLP models. These works typically
649 employ techniques such as ablation studies, model probing, or interpretability visualization to make
650 model decisions more transparent and understandable. Submissions should provide novel insights that
651 deepen our understanding of model behavior rather than merely measuring performances.
652

653 Examples:

- 654 -"We probe the model representations across syntactic constructions and show that errors on low-frequency
655 patterns stem from weak internal encoding, rather than data scarcity alone.
656 -"Our analysis identifies the amount of targeted privacy data and the extent of edited privacy neurons as
657 the two key factors contributing to this model behaviour.
658 -"We visualize attention patterns across multiple layers to understand how models resolve coreference,
659 uncovering systematic biases in entity tracking."
660

661 NLP engineering experiment

662 Papers describing the design, implementation, or deployment of NLP systems. These works generally propose
663 new models, enhancements over state-of-the-art methods, or solutions to engineering challenges in NLP
664 systems. Submissions should demonstrate clear technical contributions and real-world applicability.
665

666 Examples:

- 667 -"To address these challenges, we propose a novel fine-tuning method that employs sentence concatenation
668 with augmented random facts to regularize generation.
669 -"Our model achieve state-of-the-art performance in machine translation by introducing efficient attention
670 mechanisms, resulting in faster inference with comparable accuracy.
671 -"We deploy an end-to-end NLP pipeline for document classification, addressing engineering challenges
672 related to scalability and latency."
673

674 Publicly available software and/or pre-trained models

675 Papers providing pretrained models or NLP-related softwares, APIs, or libraries intended for broad community
676 use. These works typically provide public code repositories or access to models. Submissions should
677 demonstrate significant utility or innovation and ensure open availability to the community.
678

679 Examples:

- 680 -"We release a pretrained multilingual language model optimized for low-resource languages, along with a
681 public API for easy integration in downstream applications.
682 -"We provide a comprehensive toolkit for named entity recognition, featuring pretrained models and
683 customizable annotation interfaces, all available via a public repository.
684 -"We will publicly release our code and pre-trained models on the following URL."
685

686 Position papers

687 Papers presenting a strong perspective or argument on existing research. These works challenge existing
688 norms, give a new set of ground rule, or offer visions for the future of the field. Submissions should
689 provide well-founded arguments and contribute meaningfully beyond opinion or commentary.
690

691 Examples:

- 692 -"We argue for a paradigm shift in NLP research, advocating for more human-centered evaluation methods that
693 prioritize interpretability and fairness.
694 -"Recent debates about whether large language models understand text often stem from differing definitions
695 of understanding and views on consciousness, illustrated here by a thought experiment with a high-
696 performing but seemingly non-conscious chatbot.
697 -"We challenge the assumption that larger models inherently lead to better understanding, proposing
698 alternative evaluation metrics that better capture semantic comprehension."
699

700 Reproduction study

701 Papers analyzing, replicating and validating prior work. These works often provide new insights, clarify
702 ambiguities, or expose inconsistencies into existing works, allowing to improve confidence in research
703 findings. Submissions should offer substantial contributions beyond mere repetition or summary of
704 previous results.
705

706 Examples:

- 707 -"We replicate the experiments from the original paper to verify the reported results and assess their
708 robustness across different datasets.
709 -"By reanalyzing a landmark sentiment analysis study, we identify ambiguous evaluation metrics and propose
710 clearer standards for future work.
711 -"Our validation of recent models reveals inconsistencies in reported results."
712

713 Surveys

714 Papers synthesizing and structuring existing literature on a specific topic. These works are generally
715 indicated as such and outline methods, categorize trends, highlight gaps, and suggest future directions
716 , serving as a roadmap for researchers. Submissions should extends well beyond the typical related work
717 section of a research paper.
718

719 Examples:

- 720 -"This survey provides a comprehensive overview of recent advances in transformer-based architectures for
721 natural language processing.
722 -"In this article, we present a state-of-the-art of the main text generation approaches, including
723 evaluation data, methods, and metrics.
724 -"We provide a comprehensive roadmap for explainable NLP, reviewing current techniques, evaluating their
725 applicability, and proposing a unified framework for future research."
726

Theory
 Papers contributing to the formal or mathematical foundations of NLP. These works may include new algorithms, formal grammar models, or computational theories. Submissions should emphasize rigorous theoretical development rather than empirical evaluation.

Examples:
 -"This paper presents a new algorithm with proven computational guarantees for efficient parsing in natural language processing.
 -"We combine theoretical proofs and experimental results to establish a foundation for improving Chain-of-Thought distillation within a multitask learning framework, guided by information-theoretic principles.
 -"We propose a novel formal grammar model that rigorously characterizes syntactic structures without relying on annotated datasets."

728
 729
 730
 731
 732
 733
 734
 735
 736
 737
 738
 739

D.4 Preliminary Annotation Guidelines

741

Context
 In scholarly articles, contributions refer to the scientific achievements or innovations attributed to the authors. These may include additions to existing knowledge, theoretical advancements, methodological innovations, or the development of new artifacts or tools.
 Within the ACL Rolling Review process, authors are encouraged to specify one or more contribution types that their submission addresses. This typology is informed by prevailing research practices and thematic trends within the field of Natural Language Processing (NLP).
 However, this typology is lackluster on ACL RR, no definitions are defined and it is not clear what should constitute a candidate contribution types.

Objective
 The aim of this annotation task is to assess whether independent annotators—that is, individuals other than the original authors—can reliably identify the same contribution types by providing more extensive guidelines.
 These annotations will help evaluate the clarity and communicative effectiveness of scientific writing regarding contribution statements.

Annotation Guidelines
 1. Open the provided PDF file
 2. Using the provided typology and following definitions, please assign labels that are relevant to the document.
 Multiple contribution types can be assigned to the same document if appropriate.
 Please keep in mind that even if some contributions seems to be present in the document, the goal is to asses if the document propose a notable enough contribution so it can be qualified by that type. (e.g. A data resource-oriented paper is not the same as merely making some data available, a paper taking a stance is not the same as a position paper, a paper analysis some results does not necessarily implies a model analysis, etc.)

Typology definitions
 Approaches to low-resource settings
 Papers that propose methods or tasks designed for scenarios where labeled data, computational resources, or linguistic tools are limited. Examples of such works can include approaches leveraging generalizability, transfer learning, data augmentation or unsupervised/semi-supervised learning to address these limitations among other techniques.
 Approaches low compute settings-efficiency
 Papers that propose approaches or techniques to make NLP models more computationally efficient. Examples of such works can include optimizing memory usage, inference/training time or energy consumption to reduce computing costs among other techniques.
 Data resources
 Papers that introduces new datasets, annotated corpora, benchmarks, or other evaluation tools. Such works are often identifiable by detailed descriptions of the data creation process and by making the resources publicly accessible.
 Data analysis
 Papers focused on analyzing trends or patterns in data, such as linguistic phenomena, annotation quality, data biases or how models interact with datasets. Emphasis is on insight rather than system performance.
 Model analysis & interpretability
 Papers investigating how NLP models function internally. Examples of such works can include ablation studies, probing techniques, interpretable evaluations or tools designed to make model behavior more understandable and transparent to humans.
 NLP engineering experiment
 Papers focusing on the design, implementation, or deployment of NLP systems. Examples of such works can include experiments developing new NLP systems, proposing models that improve on the State-of-the-Art, or solving engineering challenges for NLP purpose.

742
 743
 744
 745
 746
 747
 748
 749
 750
 751
 752
 753
 754
 755
 756
 757
 758
 759
 760
 761
 762
 763
 764
 765
 766
 767
 768
 769
 770
 771
 772
 773
 774
 775
 776
 777
 778
 779
 780
 781
 782
 783
 784
 785
 786
 787
 788
 789
 790
 791
 792
 793
 794
 795
 796
 797
 798
 799
 800
 801
 802
 803
 804
 805
 806
 807
 808
 809
 810
 811
 812

813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836

Publicly available software and/or pre-trained models
Papers releasing tools, APIs, libraries, or pre-trained models that are intended for broad use by the research and development community. Such works are often identifiable by the inclusion of public repository URLs.

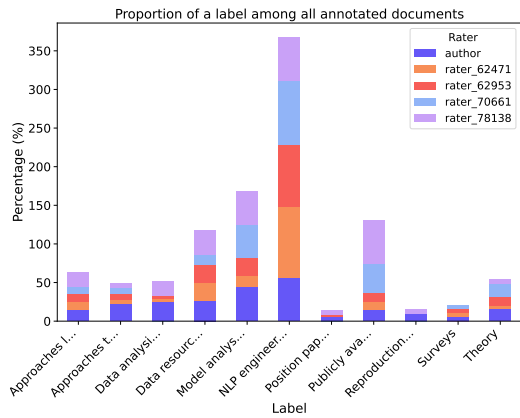
Position papers
Papers that articulate a clear stance or perspective on a significant issue in NLP. Such works may argue for specific changes in methodology, evaluation, ethical standards or future directions, often characterized by a more opinionated or speculative tone.

Reproduction study
Papers that attempts to replicate and validate the results of prior studies. Example of such works can include generalization to new settings or identification of gaps, ambiguities and errors in the original work.

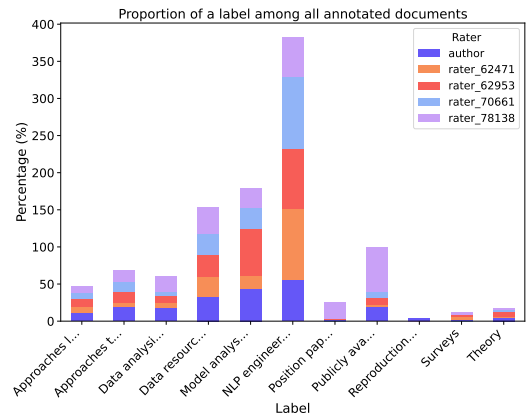
Surveys
Papers that summarize and organize existing literature on a particular topic, method, or subfield. Unlike a standard related work section, surveys aim to synthesize trends, highlight gaps, and provide a roadmap for future research.

Theory
Papers that contribute to the formal or mathematical foundations of NLP. Example of such works can include new algorithms, computational models of grammar, or rules, often without direct empirical results.

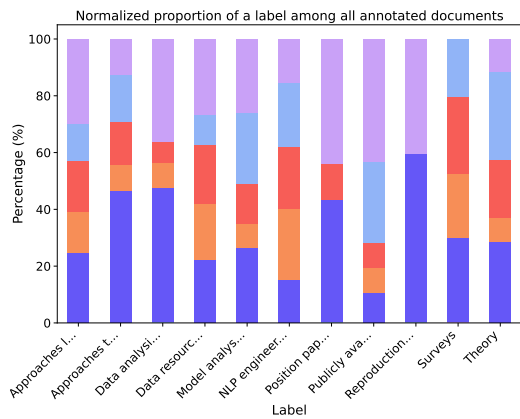
E Additional Annotations Analysis



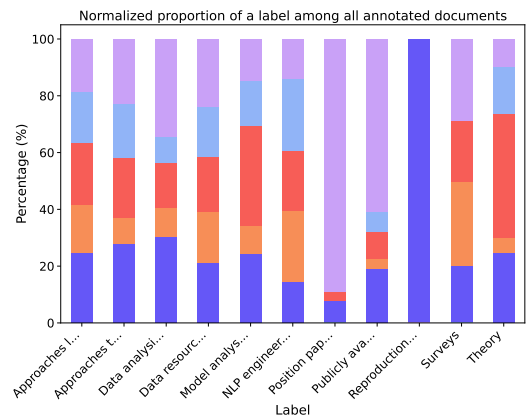
(a) Preliminary



(b) Final



(c) Normalized Preliminary



(d) Normalized Final

Figure 8: Proportion (top) and Normalized proportion (bottom) of labels selected across all annotations by the four best annotators, shown for the preliminary campaign (left) and the testing set campaign (right).

F Architecture Details

F.1 Configurations

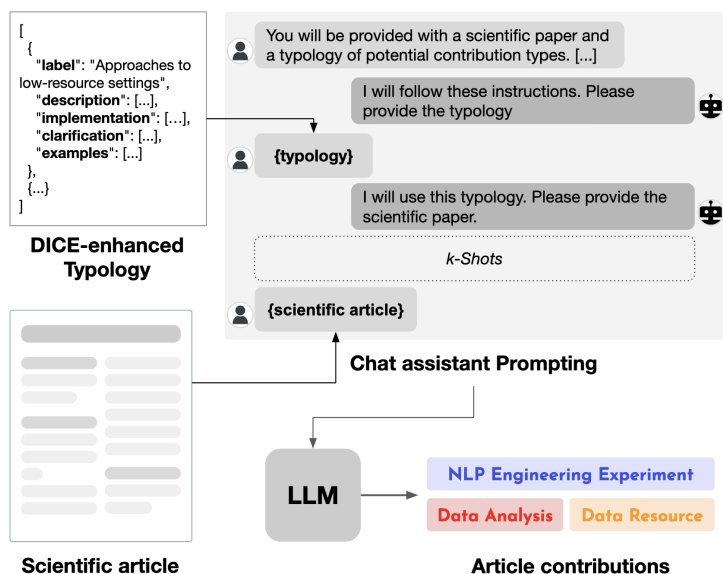


Figure 9: Workflow for identifying paper contributions via LLM prompting and DICE-enhanced definitions.

Large Language Models (LLMs). The LLM architectures included in our study are Llama (3.2–3B)³ and Mistral (7B)⁴. We use the following generation parameters for reproducibility: sampling is disabled (`do_sample=False`), the repetition penalty is set to 1.0, the no repeat n-gram size is 0, and the maximum number of generated tokens is 128.

GPT-4.1. In addition to open-weight models, we also evaluate the closed-source model GPT-4.1⁵ using the OpenAI Completion API. The temperature and presence penalty are both set to 0, and the maximum number of generated tokens is 128.

Pre-trained Language Models (PLMs). The PLMs included in our study includes both general-purpose models (BERT⁶, RoBERTa⁷) and science-focused models (SciBERT⁸, SPECTER2⁹). The models are fine-tuned using a learning rate of 1e-5, a batch size of 16, 20 training epochs with early stopping after 5, and a weight decay of 0.01.

Baselines. The baseline methods include (1) randomly label sampling from all available classes, and (2) a logistic regression classifier trained on TF-IDF features with 1 000 iterations using the liblinear solver.

① Experiments presented in this paper ran for under 50 GPU hours on two NVIDIA RTX 4500 GPUs (24GB each) and cost approximately \$35 in OpenAI credits.

F.2 LLMs Instructions

The annotation guidelines provided to the LLMs are identical to those used in the human annotation campaign.

You will be provided with a scientific paper and a typology of potential contribution types. Analyze the paper and classify the type(s) of contributions it makes based on the provided typology.

Follow these rules when making your decisions:

³<https://huggingface.co/meta-llama/Llama-3.2-3B>

⁴<https://huggingface.co/mistralai/Mistral-7B-v0.3>

⁵<https://platform.openai.com/docs/models/gpt-4.1>

⁶<https://huggingface.co/google-bert/bert-base-uncased>

⁷<https://huggingface.co/FacebookAI/roberta-base>

⁸https://huggingface.co/allenai/scibert_scivocab_uncased

⁹<https://huggingface.co/allenai/specter2>

1. Assign applicable labels. A paper may fall under multiple contribution types, assign those that are clearly present in the paper content.
2. Focus on notable contributions. Only assign label if the paper makes a substantial and deliberate contribution of that type, not just a minor mention or incidental inclusion.

F3 Prompting Layout

```
messages = [  
    {"role": "system", "content": "Follow the provided instructions. Reply ONLY the  
        list of labels from the typology applicable to the scientific paper, no  
        explanation"},  
    {"role": "user", "content": instructions},  
    {"role": "assistant", "content": "I will follow these instructions. Please  
        provide the typology."},  
    {"role": "user", "content": typology_instructions},  
    {"role": "assistant", "content": "I will use this typology. Please provide the  
        scientific paper."},  
]  
  
#Handle few-shots when requested  
for i in range(nb_shots):  
    messages.extend([  
        {"role": "user", "content": dataset["train"][i]["document"]},  
        {"role": "assistant", "content": f"Here is the list of contribution types  
            present in the scientific paper: {dataset["train"][i]["  
            contribution_types"]}"},  
    ])  
  
#Current document to process  
messages.extend([  
    {"role": "user", "content": doc},  
    {"role": "assistant", "content": "Here is the list of contribution types present  
        in the scientific paper:"},  
])
```

G Additional Experimental Results

	P	Annotators Consensus	
		R	F1 _{micro} F1 _{macro}
Authors	0.75 ±0.00	0.77 ±0.00	0.76 ±0.00 0.69 ±0.00
Baselines			
Random	0.20 ±0.00	0.55 ±0.02	0.29 ±0.01 0.23 ±0.01
TF-IDF	0.81 ±0.00	0.50 ±0.00	0.62 ±0.00 0.27 ±0.00
PLMs			
BERT Base (Uncased)	0.75 ±0.01	0.56 ±0.02	0.64 ±0.01 0.30 ±0.01
RoBERTa Base	0.74 ±0.01	0.57 ±0.01	0.64 ±0.01 0.39 ±0.02
SciBERT (SciVocab, Uncased)	0.75 ±0.00	0.59 ±0.01	0.66 ±0.01 0.40 ±0.01
SPECTER2 Base	0.74 ±0.02	0.58 ±0.01	0.65 ±0.01 0.36 ±0.03
LLMs			
LLaMA 3.2 (3B Instruct) 0-Shot	0.33 ±0.00	0.91 ±0.00	0.48 ±0.00 0.33 ±0.00
+ Definition	0.41 ±0.00	0.83 ±0.00	0.55 ±0.00 0.32 ±0.00
+ Implementation	0.44 ±0.00	0.82 ±0.00	0.57 ±0.00 0.34 ±0.00
+ Clarification	0.39 ±0.00	0.82 ±0.00	0.53 ±0.00 0.33 ±0.00
+ Examples (DICE)	0.34 ±0.00	0.83 ±0.00	0.48 ±0.00 0.32 ±0.00
LLaMA 3.2 (3B Instruct) 1-Shot	0.36 ±0.00	0.65 ±0.00	0.47 ±0.00 0.29 ±0.00
+ Definition	0.52 ±0.00	0.60 ±0.00	0.56 ±0.00 0.37 ±0.00
+ Implementation	0.61 ±0.00	0.68 ±0.00	0.64 ±0.00 0.39 ±0.00
+ Clarification	0.59 ±0.00	0.66 ±0.00	0.62 ±0.00 0.38 ±0.00
+ Examples (DICE)	0.63 ±0.00	0.67 ±0.00	0.65 ±0.00 0.37 ±0.00
LLaMA 3.2 (3B Instruct) 3-Shots	0.47 ±0.00	0.53 ±0.00	0.50 ±0.00 0.30 ±0.00
+ Definition	0.54 ±0.00	0.53 ±0.00	0.54 ±0.00 0.30 ±0.00
+ Implementation	0.59 ±0.00	0.57 ±0.00	0.58 ±0.00 0.33 ±0.00
+ Clarification	0.53 ±0.00	0.49 ±0.00	0.51 ±0.00 0.29 ±0.00
+ Examples (DICE)	0.44 ±0.00	0.45 ±0.00	0.45 ±0.00 0.25 ±0.00
Mistral 7B Instruct (v0.3) 0-Shot	0.44 ±0.00	0.74 ±0.00	0.55 ±0.00 0.36 ±0.00
+ Definition	0.52 ±0.00	0.76 ±0.00	0.62 ±0.00 0.42 ±0.00
+ Implementation	0.52 ±0.00	0.73 ±0.00	0.61 ±0.00 0.37 ±0.00
+ Clarification	0.52 ±0.00	0.71 ±0.00	0.60 ±0.00 0.43 ±0.00
+ Examples (DICE)	0.46 ±0.00	0.68 ±0.00	0.55 ±0.00 0.36 ±0.00
Mistral 7B Instruct (v0.3) 1-Shot	0.46 ±0.00	0.65 ±0.00	0.54 ±0.00 0.33 ±0.00
+ Definition	0.52 ±0.00	0.71 ±0.00	0.60 ±0.00 0.39 ±0.00
+ Implementation	0.57 ±0.00	0.73 ±0.00	0.64 ±0.00 0.43 ±0.00
+ Clarification	0.58 ±0.00	0.69 ±0.00	0.63 ±0.00 0.42 ±0.00
+ Examples (DICE)	0.60 ±0.00	0.69 ±0.00	0.64 ±0.00 0.42 ±0.00
Mistral 7B Instruct (v0.3) 3-Shots	0.53 ±0.00	0.72 ±0.00	0.61 ±0.00 0.36 ±0.00
+ Definition	0.62 ±0.00	0.75 ±0.00	0.68 ±0.00 0.44 ±0.00
+ Implementation	0.66 ±0.00	0.73 ±0.00	0.69 ±0.00 0.47 ±0.00
+ Clarification	0.67 ±0.00	0.74 ±0.00	0.70 ±0.00 0.51 ±0.00
+ Examples (DICE)	0.66 ±0.00	0.73 ±0.00	0.69 ±0.00 0.49 ±0.00
GPT-4.1 (2025-04-14) 0-Shot	0.55 ±0.00	0.85 ±0.01	0.67 ±0.01 0.54 ±0.00
+ Definition	0.60 ±0.01	0.80 ±0.01	0.69 ±0.01 0.55 ±0.00
+ Implementation	0.63 ±0.00	0.78 ±0.01	0.70 ±0.00 0.56 ±0.01
+ Clarification	0.64 ±0.00	0.76 ±0.00	0.69 ±0.00 0.56 ±0.01
+ Examples (DICE)	0.65 ±0.01	0.75 ±0.00	0.70 ±0.00 0.55 ±0.01
GPT-4.1 (2025-04-14) 1-Shot	0.62 ±0.01	0.80 ±0.00	0.70 ±0.00 0.54 ±0.00
+ Definition	0.68 ±0.00	0.79 ±0.01	0.73 ±0.00 0.57 ±0.02
+ Implementation	0.69 ±0.00	0.75 ±0.01	0.72 ±0.01 0.56 ±0.00
+ Clarification	0.69 ±0.00	0.73 ±0.01	0.71 ±0.01 0.55 ±0.01
+ Examples (DICE)	0.71 ±0.01	0.74 ±0.01	0.72 ±0.01 0.55 ±0.01
GPT-4.1 (2025-04-14) 3-Shots	0.68 ±0.00	0.80 ±0.00	0.74 ±0.00 0.56 ±0.01
+ Definition	0.69 ±0.00	0.78 ±0.00	0.73 ±0.00 0.56 ±0.00
+ Implementation	0.69 ±0.00	0.76 ±0.00	0.72 ±0.00 0.56 ±0.01
+ Clarification	0.70 ±0.00	0.74 ±0.00	0.72 ±0.00 0.56 ±0.00
+ Examples (DICE)	0.69 ±0.00	0.75 ±0.00	0.72 ±0.00 0.56 ±0.01

Table 3: Detailed performances of models with ablation comparison of the typology components. Reported scores are averaged across five random seeds.