

# COUNTERFACTUAL LLM-BASED FRAMEWORK FOR MEASURING RHETORICAL STYLE

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

The rise of AI has fueled growing concerns about “hype” in machine learning papers, yet a reliable way to quantify rhetorical style independently of substantive content has remained elusive. Because strong empirical results can justify stronger claims, it is often unclear whether bold language reflects genuine evidence or merely rhetorical style. We introduce a counterfactual, LLM-based framework to disentangle rhetorical style from substantive content: multiple LLM rhetorical personas generate counterfactual writings from the same substantive content, an LLM judge compares them through pairwise evaluations, and the outcomes are aggregated using a Bradley–Terry model. Applying this method to 8,485 ICLR submissions sampled from 2017 to 2025, we generate more than 250,000 counterfactual writings and provide a large-scale quantification of rhetorical style in ML papers. Visionary framing significantly predicts downstream attention, including citations and media coverage, even after controlling for peer-review evaluations. We also observe a sharp rise in rhetorical strength after 2023, aligning with the growing use of LLM-based writing assistance. The reliability of our framework is validated by its robustness to the choice of personas and the high correlation between LLM judgments and human annotations. Our work demonstrates that LLMs can serve as instruments for improving how ML research is evaluated.

## 1 INTRODUCTION

Machine learning is a field of extraordinary excitement: every year brings breakthroughs in models, applications, and benchmarks. Meanwhile, the ML community has become intensely competitive, with leading conferences now receiving thousands of submissions (e.g., ICLR 2025 received over 10,000). In this crowded environment, visibility is scarce, and authors face strong incentives to emphasize and sometimes overstate the significance of their contributions (Smaldino & McElreath, 2016; McGreivy & Hakim, 2024).

This pattern has fueled concerns about hype, a rhetorical style that exaggerates novelty or impact (Thais, 2024). Concerns about hype are not merely stylistic: rhetorical framing is often suspected of influencing how papers are perceived in peer review and in attracting downstream attention. Despite these concerns, there is no systematic framework for measuring rhetorical style in scientific writing. Addressing this measurement gap is critical. Prior studies document rising use of positive framing and promotional language (Vinkers et al., 2015; Peng et al., 2024; Millar et al., 2024), declining expressions of uncertainty (Yao et al., 2023), and even widespread misrepresentation in published research (Boutron & Ravaud, 2018; McGreivy & Hakim, 2024).

Existing approaches mostly focus on sub-concepts of rhetorical style directly from the text: either by scoring promotional lexicons and indexes (Mishra et al., 2023; Peng et al., 2024; Gentzkow & Shapiro, 2010) or by training models on human-labeled constructs such as “sensationalism” or “uncertainty” (Prabhakaran et al., 2016; Pei & Jurgens, 2021; Wühl et al., 2024). These methods are effective for capturing surface-level tone, but they operate only on the text itself, without accounting for the underlying substantive content. Because strong empirical results can justify stronger claims, it is often unclear whether bold language reflects genuine evidence or merely rhetorical style. The central challenge, then, is to **distinguish a paper’s substantive contribution from its presentation**.

We address this challenge by formalizing rhetorical style, conditioned on fixed substantive content, as a measurement problem. The core idea is simple: hold the substantive content  $X$  (methods, experiments,

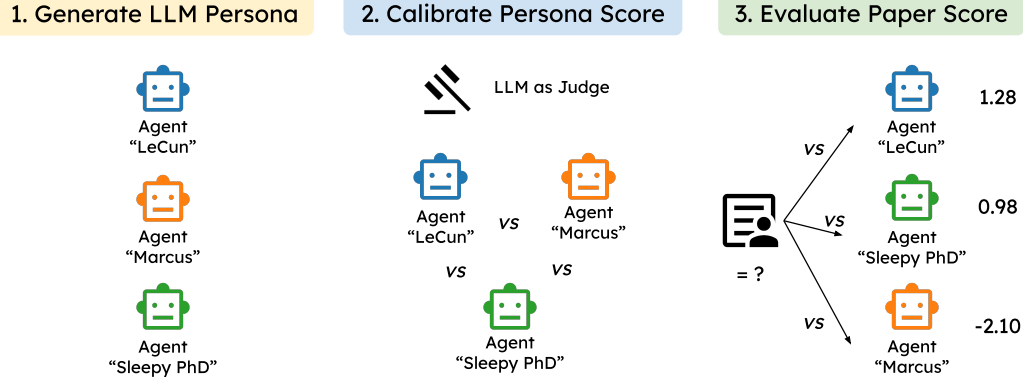


Figure 1: (1) LLM personas generate counterfactual writings in different rhetorical styles based on the same substantive content. (2) We calibrate the LLM personas’ rhetorical scores via pairwise comparisons using an LLM judge. (3) We infer the rhetorical score of any query abstract by comparing it against the calibrated LLM persona panel.

results) constant, vary the language  $Y$ , and infer a latent variable  $Z$  that captures rhetorical strength independently of substantive contribution. Just as benchmarks quantify empirical progress, this formulation provides a quantitative basis for analyzing rhetorical style and understanding how it shapes attention in ML research.

**Our Contributions.** We develop a counterfactual, LLM-based framework that measures rhetorical style independently of substantive content. As shown in Figure 1, we employ a diverse panel of LLM personas that generate counterfactual writings in different rhetorical styles from the same substantive content. These persona outputs are then compared in pairs by an LLM judge, and the pairwise outcomes are aggregated with a Bradley–Terry model (Bradley & Terry, 1952) to assign each persona a calibrated rhetorical score. Finally, the rhetorical score of any new abstract can be inferred by comparing it against this calibrated persona panel.

Conceptually, this parallels a counterfactual design in causal inference: we ask how the evaluation of a paper would change if the rhetorical style were altered while the underlying substantive content remained the same. By coupling controlled generation with pairwise comparison under a classic statistical framework, we introduce a scalable and principled instrument for improving how ML research is evaluated.

- **Methodology:** We introduce a counterfactual evaluation paradigm using a calibrated panel of LLM personas as controlled writers that generate counterfactual abstracts from identical scientific content, with pairwise judgments aggregated via a Bradley–Terry model.
- **Validation:** Unlike prior approaches that infer rhetorical strength directly from text, our approach controls for content, produces fine-grained scores, and is robust to persona choice. We further validate the LLM judgments via human annotations.
- **Findings:** Applying our framework to 8,485 ICLR submissions from 2017 to 2025 and more than 250,000 generated persona writings and 250,000 pairwise comparisons, we provide large-scale evidence that visionary rhetorical style predicts citations and media attention, and that rhetorical strength has risen significantly since 2023.

At a high level, our setup resembles a GAN, where generators are evaluated by discriminators (Goodfellow et al., 2014). It also parallels Reinforcement Learning from Human Feedback (RLHF) (Christiano et al., 2017; Ouyang et al., 2022) and Direct Preference Optimization (DPO) (Rafailov et al., 2023), which compare multiple outputs from the same input to learn preference models. Related approaches in political science likewise leverage multiple texts describing the same event to infer latent attributes such as political leaning via pairwise comparisons (Carlson & Montgomery, 2017; Licht et al., 2025). However, in our case, the focus is on measuring rhetorical style while holding content constant.

## 2 PROBLEM SETTING

### 2.1 PROBLEM FORMULATION

We formalize rhetorical style as a latent variable in a generative model of research writing. Specifically, we denote the text of a paper as  $Y$ , which is influenced by its substantive content  $X$  and a latent rhetorical style variable  $Z$ .

For simplicity, we focus on measuring rhetorical style in abstracts, since they are typically the first section readers encounter and contain the core articulation of a paper’s contributions and significance. Accordingly, we define  $Y$  as the abstract text, while  $X$  refers to the substantive basis, such as preliminaries, methods, experiments, and results, which are comparatively objective and descriptive. The rhetorical style  $Z$  is modeled as a latent continuous scalar capturing the strength of framing. Formally,  $X \in \{1, \dots, T\}^n$  and  $Y \in \{1, \dots, T\}^m$  are token sequences of finite length, and  $Z \in \mathbb{R}$  is a one-dimensional hidden variable representing rhetorical strength.

We posit that the observed writing  $Y$  is generated from a conditional distribution:

$$p(Y \mid X, Z) \quad (1)$$

where  $X$  refers to the substantive content and  $Z$  modulates the framing of the contribution. The objective of this work is to infer  $Z$  independently of  $X$ , yielding a quantitative measure of rhetorical strength disentangled from substantive content.

Traditional approaches to measuring rhetorical style in scientific texts rely on direct linguistic features of  $Y$ , but such measures are confounded by  $X$ , since stronger results naturally justify stronger claims. Our approach instead targets  $Z$ , isolating rhetorical style conditional on fixed substantive content.

### 2.2 RHETORICAL STRENGTH

We assume a total order exists on the space of latent rhetorical style  $Z$ , as  $Z \in \mathbb{R}$ . A higher value of  $Z$  corresponds to a “stronger” rhetorical style (e.g., more visionary), while a lower value suggests a “weaker” rhetorical style (e.g., more conservative).

**Definition 2.1** (Stronger Rhetorical Style). *A **stronger rhetorical style** refers to a rhetorical style that broadly expands its implications. This involves generalizing its applicability, highlighting the significant challenges it solves, and emphasizing its novelty and impact.*

**Definition 2.2** (Weaker Rhetorical Style). *A **weaker rhetorical style** refers to a rhetorical style that narrowly confines its implications, potentially by specifying its applicability and prerequisites, positioning it closely to previous work, and emphasizing its limitations and uncertainties.*

For any two rhetorical styles  $z_1, z_2 \in Z$ , if  $z_1 > z_2$ , then  $z_1$  represents a stronger rhetorical style than  $z_2$ .

**Definition 2.3** (Order of Rhetorical Strength). *Given a fixed substantive basis  $X = x$ , for any two writings  $y_1, y_2$  derived from  $(x, z_1)$  and  $(x, z_2)$ , we say  $y_1$  has a **stronger rhetorical style** than  $y_2$  (denoted  $y_1 \succ_x y_2$ ) if and only if  $z_1 > z_2$ . It implies  $y_1$  presents the writing with greater assertiveness, broader scope, or higher proclaimed impact than  $y_2$ .*

In practice, judgments about whether  $y_1$  is rhetorically stronger than  $y_2$  can be noisy, affected by both uncertainty and observation error. We therefore model these comparisons using the Bradley–Terry framework (Section 3.3).

## 3 COUNTERFACTUAL LLM-BASED FRAMEWORK

In this section, we introduce a counterfactual LLM-based framework for measuring latent rhetorical style. As shown in Figure 1, the method proceeds in two stages. First, we construct a panel of LLM persona writers with specified rhetorical tendencies and calibrate their rhetorical scores via pairwise judgments from an LLM judge. Second, we infer the rhetorical score of any query abstract by comparing it against the calibrated persona panel.

Our framework relies on two assumptions: (i) the persona panel adequately spans the distribution of query abstracts, and (ii) the LLM judge can distinguish rhetorical strength with non-trivial accuracy. We later provide evidence that both assumptions hold in practice in Sections 4.3 and 4.2, respectively.

### 3.1 COUNTERFACTUAL GENERATION WITH LLM PERSONAS

Given the substantive content  $x$ , we sample  $K$  counterfactual abstracts by prompting the LLM with  $K$  different personas that embody diverse rhetorical styles in academic writing:

$$y_{A_k} \sim \text{LLM}(x, \text{prompt}_{A_k}), \quad k = 1, \dots, K \quad (2)$$

This yields a set of diverse counterfactual abstracts  $\{y_{A_1}, y_{A_2}, \dots, y_{A_K}\}$  for the same underlying content  $x^i$ , each differing only in rhetorical framing.

Each persona  $A_k$  is defined by a system prompt that enforces a distinct rhetorical style (e.g., cautious, visionary, technical). The full list of personas and their descriptions is provided in Appendix C.1. Note that the key requirement here is not the specific identity of the personas, but that the panel spans a sufficiently broad range of rhetorical styles. To provide intuition, we generate three versions of our paper’s abstract, each written from a different persona using our framework (Appendix A).

### 3.2 PAIRWISE EVALUATION WITH LLM JUDGES

To establish an ordering over persona-generated abstracts, we use an LLM-based judge that performs pairwise comparisons of rhetorical strength.

Given two abstracts  $y_{A_1}$  and  $y_{A_2}$ , both derived from the same substantive content  $x$  but generated by different personas  $A_1$  and  $A_2$ , the LLM judge is prompted to decide which abstract makes stronger or more sensationalized claims, and to provide a brief rationale for its choice.

### 3.3 CALIBRATION VIA BRADLEY–TERRY MODEL

To aggregate pairwise comparisons into a global ordering of rhetorical strength, we employ the Bradley–Terry model (Bradley & Terry, 1952).

Let  $\pi_{A_k}$  denote the rhetorical strength parameter for persona  $A_k$ . Given two abstracts  $y_{A_1}$  and  $y_{A_2}$ , generated from the same content  $x$  but with different personas  $A_1$  and  $A_2$ , the probability that  $y_{A_1}$  is judged stronger than  $y_{A_2}$  is:

$$P(y_{A_1} \succ y_{A_2}) = \frac{\pi_{A_1}}{\pi_{A_1} + \pi_{A_2}}. \quad (3)$$

For each pair of personas  $(A_1, A_2)$ , we sample  $M$  instances  $\{x^i\}_{i=1}^M$ , generate abstracts  $\{y_{A_1}^i\}_{i=1}^M$  and  $\{y_{A_2}^i\}_{i=1}^M$ , and obtain LLM-judge comparisons. Aggregating these across all persona pairs, we estimate  $\hat{\pi} = \{\pi_{A_1}, \dots, \pi_{A_K}\}$  by maximum likelihood. We then define the continuous rhetorical score  $s_k = \log(\pi_k)$ , which places each persona on a one-dimensional spectrum of rhetorical strength.

### 3.4 INFERENCE FOR QUERY ABSTRACTS WITH THE REFERENCE PANEL

After establishing the calibrated scale of  $K$  personas with rhetorical strength parameters  $\{\pi_{A_1}, \dots, \pi_{A_K}\}$ , our goal is to infer the rhetorical strength of a new query abstract  $y_q$ . We do so by positioning  $y_q$  on the same scale through comparisons against the persona abstracts. For each query abstract, the LLM judge conducts  $K$  pairwise comparisons, one against each persona abstract. Under the Bradley–Terry model, the probability that the query abstract wins against persona  $A_k$  is  $P(y_q \succ y_{A_k}) = \frac{\pi_q}{\pi_q + \pi_{A_k}}$ .

While maximizing this likelihood with Maximum Likelihood Estimation (MLE) works well for calibrating persona scores because persona abstracts are compared against each other many times, query abstracts face a much sparser setting, since each is compared only once against each persona abstract. In such cases, MLE can yield degenerate solutions: if a query abstract wins all its comparisons, the likelihood drives its score to infinity; if it loses all, the score collapses toward zero. To obtain stable estimates, we instead use Maximum a Posteriori (MAP) estimation, which introduces regularization

through a prior. Let  $s_q = \log(\pi_q)$ . We place a Gaussian prior on the query abstract’s score. The MAP estimate is then the value that maximizes the log-posterior, combining the Bradley–Terry likelihood with the log-prior penalty.

**Adaptive Bayesian Inference.** We also consider an alternative, more cost-efficient estimation strategy based on adaptive Bayesian inference. Rather than using a fixed batch of comparisons, this approach maintains a posterior distribution over the query abstract’s score and sequentially selects comparisons to maximize information gain. At each step, the next persona is chosen adaptively: the persona whose score is closest to the current posterior median, since this comparison is expected to reduce uncertainty most effectively. The posterior is then updated with the observed outcome, and the process terminates once the posterior variance falls below a predefined threshold. This yields a point estimate of rhetorical strength with an associated confidence interval that quantifies the measurement uncertainty. In this study, we adopt batch MAP estimation as it is sufficient at our scale.

## 4 EXPERIMENTS

To investigate rhetorical style in machine learning papers, we compiled a dataset of 8,485 research papers submitted to the International Conference on Learning Representations (ICLR) from 2017 to 2025, randomly sampling 1,000 submissions from each year.

Our data preprocessing pipeline involved two steps. First, we extracted the full text from each paper’s PDF. Second, we used an LLM-based extraction method (OpenAI GPT-4o-mini) to identify and extract the substantive content from the experiments, methods, and results sections in research papers, which serve as the substantive basis (denoted as  $X$ ) in our framework. This filtering process removes all narratives and references from the paper. The paper abstracts were used as the observed writing ( $Y$ ) for analysis.

### 4.1 IMPLEMENTATION

Our framework proceeds in two stages. First, we calibrate a stable rhetorical scale using a panel of 30 LLM personas, ranging from archetypes such as “Sleep-Deprived PhD Student” to prominent figures such as “Geoffrey Hinton” and “Yann LeCun” (full list in Appendix C.1). For each pair of personas, we sample 20 papers, generate counterfactual abstracts from identical methods and results sections (with length constrained to within  $\pm 15$  words of the original abstract), and ask GPT-4o to judge which abstract makes stronger claims. Aggregating 8,700 such pairwise comparisons with a Bradley–Terry model produces a continuous spectrum from conservative to promotional styles.

Second, we apply this calibrated scale to 8,485 ICLR submissions. For each paper, we generate 30 persona abstracts and compare the original abstract against each of them, yielding 30 judgments per paper. In total, this stage produces 254,550 additional comparisons, which we again aggregate via Bradley–Terry to obtain each paper’s rhetorical score. Full persona and judge prompts appear in Appendices C.2 and C.3.

**Direct Rating Baseline.** As a point of comparison, we implement a direct LLM rating strategy. To mitigate the content–style confound of naive prompting, the model is provided with both the original abstract  $Y_q$  and the extracted methods and results  $X_q$ , and is asked to rate the degree of overclaiming on a 1–10 scale (1 = “no overclaiming,” 10 = “extreme overclaiming”). In effect, this prompt directly asks the model to infer the latent rhetorical style  $Z_q$  by comparing the claims to the evidence. The model is required to output both a numeric score and a justification. The full prompt is shown in Appendix C.4.

**Keyword- and Classifier-Based Baselines.** We also implement two widely used baselines from prior work. First, following Peng et al. (2024), we compute a promotion score as the proportion of words in an abstract that appear in a curated promotional lexicon. Second, following Pei & Jurgens (2021), we estimate the certainty level using a pretrained sentence-level certainty classifier, averaging predictions across all sentences. Both baselines rely solely on direct analysis of the observed text  $Y_q$ , and therefore risk conflating rhetorical style with underlying content.

## 4.2 VALIDATION

**Validating LLM Sampler.** We test robustness to the choice of personas using a complementary subset strategy. In each trial, the 30 personas are randomly split into two disjoint groups of 15, and Bradley–Terry scores of personas and original abstracts are recalculated independently within each subset. We repeat this procedure 1,000 times to ensure stability. The resulting BT scores remain highly consistent across non-overlapping subsets, with the mean Spearman correlations being 0.89. This indicates that the relative ordering of rhetorical strength of papers is strongly preserved regardless of the specific personas used. It also demonstrates that our framework captures genuine rhetorical style differences and is applicable beyond our specific set of 30 personas.

**Validating LLM Judge.** To assess whether our automated measurements align with human perception, we conduct a human annotation study via Prolific. Annotators compare abstracts pairwise and select which makes stronger claims, following the same instructions as the LLM judge. We collect 420 judgments from 42 participants across 69 unique comparisons (45 persona–persona, 24 original–persona), each evaluated by an average of 6.1 participants with majority-vote aggregation. To ensure data quality, participants were only able to proceed to the main task if they correctly answered two pre-defined qualification questions.

Human validation confirms that our framework aligns closely with human perception. At the pairwise level, human majority votes agree with the LLM judge in 88.4% of comparisons. At the aggregate level, the Bradley–Terry scores derived from human majority votes and those derived from the LLM judge are strongly correlated (Spearman  $\rho = 0.92$ ,  $p < 0.001$ ). This suggests that our automated method captures rhetorical strength in a manner consistent with human evaluation.

## 4.3 DISTRIBUTION OF RHETORICAL STYLE MEASUREMENTS

We apply our framework to 8,485 ICLR submissions to estimate rhetorical style at scale. As shown in Figure 2, Bradley–Terry scores follow an approximately Gaussian distribution (−4.74 to 4.53), which captures substantial variation in how authors frame their contributions. By contrast, direct rating scores cluster heavily around values of 2–3, yielding a coarse, skewed distribution. We focus on these two measures because both consider the substantive content ( $X$ ) and abstract ( $Y$ ), making them directly comparable. For reference, distributions of other baselines (promotion and certainty scores), which analyze only the observed text without conditioning on  $X$ , are provided in Appendix Figure 5.

To validate that our persona panel provides sufficient coverage, Figure 2 reports win rates of all personas against query papers. The wide spread of win rates, from highly assertive personas winning over 90% of comparisons to more conservative personas winning less than 20%, demonstrates that the panel adequately spans the rhetorical space of query abstracts.

## 4.4 PREDICTIVE VALIDITY

We first examine whether rhetorical style is associated with peer-review ratings. The Bradley–Terry score shows almost no correlation with the average reviewer score (Spearman’s  $\rho = -0.015$ ,  $p = 0.225$ ), and regression estimates likewise indicate no significant relationship. This likely reflects the scope of our measure: it focuses only on abstracts, which highlight framing but omit the technical details that reviewers evaluate. If rhetorical style were measured across full papers rather than abstracts alone, it might exhibit a different relationship with reviewer evaluations.

We next evaluate whether rhetorical style predicts downstream attention after controlling for average reviewer scores (as a proxy for paper quality), as well as year and subfields. Table 1 reports regression results for two classes of outcomes: (i) scholarly impact, measured by citations; and (ii) media attention, measured by Altmetric indicators including posts, tweets, RSS feeds, patents, and unique accounts mentioning the paper.<sup>1</sup> The coefficients indicate the expected change in the outcome

<sup>1</sup>Altmetric indicators track mentions of research across diverse sources. *Posts* refer to mentions in news outlets and blogs; *Tweets* are Twitter (X) posts mentioning the paper, counted once per account; *Feeds* denote blog mentions that Altmetric monitors via Really Simple Syndication (RSS). Altmetric maintains a curated list of academic and science-related blogs, and when one of these blogs posts about a paper, it appears as a “feed” mention; *Patents* reflect citations of the paper in patent filings; and *Accounts* indicate the number of distinct Twitter (X) users who mentioned the paper.

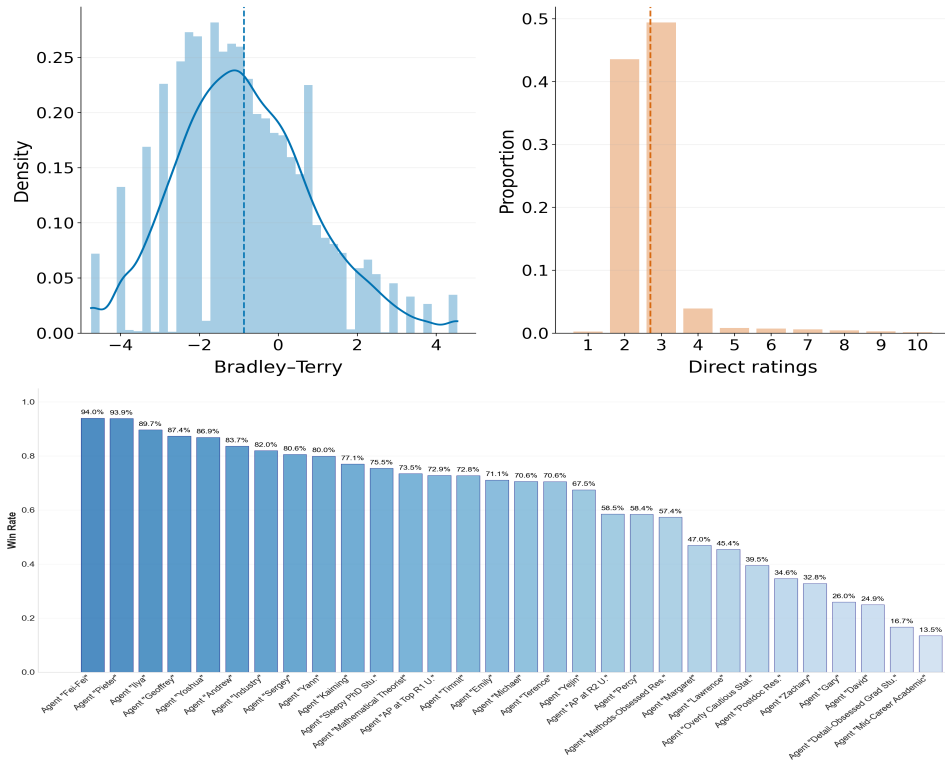


Figure 2: Distribution of our model (left) and direct ratings (right). The Bradley-Terry scores are shown as a density histogram with kernel density estimate. Direct ratings are shown as proportions of observed values. Their means are indicated by dashed lines. (Bottom) Persona Win Rates Against the Sampled 8,485 Query Papers.

variable for a one-unit increase in the rhetorical style score, holding other factors constant. Statistical significance is denoted by stars, where more stars denote more confidence that the relationship exists. A corresponding visual summary of these coefficients is provided in Appendix Figure 6.

Table 1 shows that higher Bradley-Terry scores significantly predict both citations and media attention. For example, a one-unit increase in the Bradley-Terry score is associated with 24 additional citations, 3 more media posts, and 2 more tweets, on average. For context, the effect size of rhetorical style on citations is comparable to that of review scores (89 citations per one-unit increase), which suggests that rhetorical style has a nontrivial effect.

In contrast, baseline measures (direct ratings, promotion, certainty) fail to consistently predict these attention outcomes and often yield unstable or inconsistent coefficients. Taken together, these results suggest that our framework more reliably isolates rhetorical style from substantive content and captures variation that systematically influences how research is received in both scholarly and public domains.

These results suggest that the rhetorical style of paper abstracts is more predictive of downstream attention than of reviewer evaluations. One possible explanation is that broader audiences often engage primarily with abstracts, where rhetorical framing is most salient, whereas reviewers typically evaluate the full paper and base their assessments on both the full paper presentation and its technical quality.

#### 4.5 TEMPORAL AND SUBFIELD TRENDS

We also analyze how rhetorical style evolves over time and differs across research subfields. Figure 3 shows that average Bradley-Terry scores steadily decline between 2018 and 2022, followed by a sharp rebound after 2023. This rebound coincides with stronger attention pressures in ML and

	Citation	Post	Tweet	Feeds	Patent	Account
Bradley–Terry Score (Ours)	24.53*** (6.40)	3.19*** (0.48)	2.51*** (0.39)	0.03*** (0.00)	0.04** (0.01)	2.71*** (0.40)
Direct Rating Score	-26.11* (13.17)	0.74 (0.99)	0.75 (0.80)	0.00 (0.01)	0.01 (0.03)	0.77 (0.83)
Promotion Score	20.01 <sup>†</sup> (10.24)	0.64 (0.77)	0.51 (0.62)	0.02* (0.01)	0.02 (0.02)	0.57 (0.65)
Certainty Score	59.56 (76.96)	-12.74* (5.78)	-9.74* (4.66)	-0.02 (0.05)	0.17 (0.15)	-9.97* (4.86)

Table 1: Regression coefficients with standard errors in parentheses. Each coefficient is estimated from a separate regression model in which the outcome is either scholarly impact (citations) or media attention (posts, tweets, feeds, patents, or accounts). The focal predictor in each specification is one of the four rhetorical style measures (Bradley–Terry score, direct rating, promotion, or certainty). All specifications include controls for average peer review rating, research subfield, and publication year to account for differences in paper quality, field-specific variation, and changes over time. Statistical significance is denoted by stars: <sup>†</sup>  $p < 0.10$ , \*  $p < 0.05$ , \*\*  $p < 0.01$ , \*\*\*  $p < 0.001$ .

widespread adoption of LLM-based writing assistance (Liang et al., 2024; Bao et al., 2025), which collectively reduce the cost of producing more assertive abstracts and may nudge phrasing toward visionary framing.

To investigate differences across subfields, we classify papers using an LLM (GPT-4o) with a curated prompt provided in Appendix C.5. As shown in Figure 4, we do observe systematic differences in rhetorical style across subfields. Applied areas such as computer vision, NLP, and computational biology exhibit more promotional rhetorical styles, while theoretically oriented fields such as kernel methods, optimal transport, and supervised representation learning adopt more conservative framing. Interdisciplinary topics connecting technical and social concerns (e.g., fairness, privacy, interpretability) fall in between. These patterns suggest that rhetorical style reflects epistemic orientation and community norms.

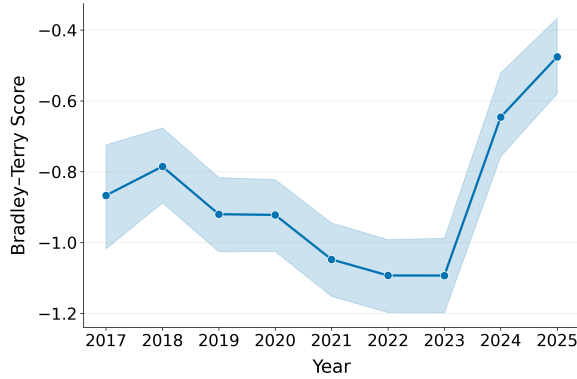


Figure 3: Yearly trends in Bradley–Terry rhetorical style scores. Points show yearly means with shaded bands indicating 95% confidence intervals.

## 5 DISCUSSION

This work introduces a counterfactual, LLM-based framework for measuring rhetorical style in ML papers. By conditioning on identical substantive content and varying only rhetorical framing, we show that LLM-based comparisons provide a calibrated, fine-grained measure of style disentangled from technical merit. Applying this framework to 8,485 ICLR submissions yields large-scale evidence that rhetorical style is not merely cosmetic: visionary or confident framing significantly predicts both scholarly impact (citations) and public attention (media mentions), even after accounting for reviewer scores as a proxy for quality. Although our framework currently depends on LLM judges and therefore inherits their biases, we validate its reliability by showing that it is robust to the choice of personas and that LLM-based judgments strongly agree with human annotations.

For the ML community, our findings highlight that *how* research is presented can strongly shape *how* it is received. When stylistic choices amplify attention independently of technical merit, they raise normative concerns for ML venues. In particular, if assertive or visionary framing reliably



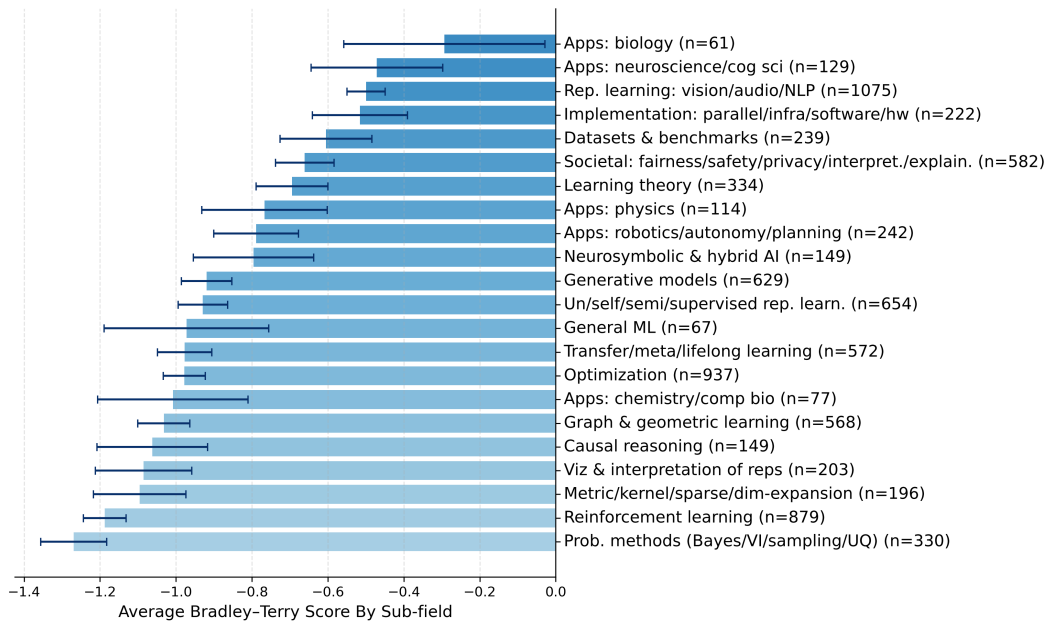


Figure 4: Mean scores by subfields for Bradley-Terry scores. Points indicate subfield-specific mean values, with horizontal bars showing 95% confidence intervals based on the standard error of the mean. For readability, figure legends use standardized abbreviations of topic names; see Appendix Table 2 for the full topic mapping.

boosts visibility, conferences risk encouraging a cycle of “rhetorical inflation,” where increasingly bold claims are rewarded over careful exposition. To counteract this, potential interventions include reviewer training, clearer guidelines that separate substance from style, and community norms that place equal value on clarity and precision as on ambition. Without such measures, rhetorical inflation may distort which contributions receive recognition and, over time, shape the trajectory of the field. We emphasize that our work does not advocate hype; rather, it calls for the importance of cultivating clarity in research communication, not just visionary framing.

**Limitations and future work.** Our current implementation uses a fixed set of hand-crafted personas, but we see this as only a first step. Future work could pursue more systematic persona construction, such as chain-based sampling, where rhetorical descriptors are gradually adjusted to generate a smooth and continuous spectrum of styles.

Beyond persona design, our framework currently focuses on the paper-level rhetorical styles. It could be interesting to investigate the long-term consequences for individual researchers. For example, whether rhetorical style influences how researchers are cited, invited to give talks, considered for collaborations, awarded funding opportunities, or evaluated in hiring and promotion.

Finally, it will be important to examine how the widespread use of generative-AI writing tools is reshaping rhetorical norms in ML submissions. For example, Liang et al. (2024) and Bao et al. (2025) document empirical trends in how much content is being modified or produced by LLMs in scientific papers. We hope our counterfactual LLM-based framework can be integrated to develop better and scalable reviewing systems for ML conferences.

## ETHICS STATEMENT

Our work analyzes rhetorical style in machine learning papers using publicly available submissions to ICLR between 2017 and 2025. All data were obtained from papers that were already accessible to the research community; no private, proprietary, or personally identifiable information was used. While we used LLMs to generate counterfactual writings and perform pairwise judgments, we validated these automated outputs with a human annotation study. All human annotations were collected through Prolific with fair compensation and qualification screening. Importantly, our study does not endorse or encourage hype, and we speculate that hype may have detrimental long-term consequences for researchers.

## REFERENCES

- Honglin Bao, Mengyi Sun, and Misha Teplitskiy. Where there’s a will there’s a way: Chatgpt is used more for science in countries where it is prohibited. *Quantitative Science Studies*, pp. 1–16, 2025.
- Isabelle Boutron and Philippe Ravaud. Misrepresentation and distortion of research in biomedical literature. *Proceedings of the National Academy of Sciences*, 115(11):2613–2619, 2018.
- Ralph Allan Bradley and Milton E Terry. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952.
- David Carlson and Jacob M Montgomery. A pairwise comparison framework for fast, flexible, and reliable human coding of political texts. *American Political Science Review*, 111(4):835–843, 2017.
- Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30, 2017.
- Matthew Gentzkow and Jesse M Shapiro. What drives media slant? evidence from us daily newspapers. *Econometrica*, 78(1):35–71, 2010.
- Ian J Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014.
- Weixin Liang, Yaohui Zhang, Zhengxuan Wu, Haley Lepp, Wenlong Ji, Xuandong Zhao, Hancheng Cao, Sheng Liu, Siyu He, Zhi Huang, et al. Mapping the increasing use of llms in scientific papers. *arXiv preprint arXiv:2404.01268*, 2024.
- Hauke Licht, Rupak Sarkar, Patrick Y Wu, Pranav Goel, Niklas Stoeck, Elliott Ash, and Alexander Miserlis Hoyle. Measuring scalar constructs in social science with llms. *arXiv preprint arXiv:2509.03116*, 2025.
- Nick McGreivy and Ammar Hakim. Weak baselines and reporting biases lead to overoptimism in machine learning for fluid-related partial differential equations. *Nature Machine Intelligence*, 6(10):1256–1269, 2024.
- Neil Millar, Bryan Mathis, Bojan Batalo, and Brian Budgell. Trends in the expression of epistemic stance in nih research funding applications: 1985–2020. *Applied Linguistics*, 45(4):658–675, 2024.
- Apratim Mishra, Jana Diesner, and Vetle I Torvik. A probabilistic model of ‘hype’ in scientific abstracts. In *International Society of Scientometrics and Informetrics Conference 2023 (ISSI)*, 2023.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.
- Jiaxin Pei and David Jurgens. Measuring sentence-level and aspect-level (un) certainty in science communications. *arXiv preprint arXiv:2109.14776*, 2021.

- Hao Peng, Huilian Sophie Qiu, Henrik Barslund Fosse, and Brian Uzzi. Promotional language and the adoption of innovative ideas in science. *Proceedings of the National Academy of Sciences*, 121(25):e2320066121, 2024.
- Vinodkumar Prabhakaran, William L Hamilton, Dan McFarland, and Dan Jurafsky. Predicting the rise and fall of scientific topics from trends in their rhetorical framing. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1170–1180, 2016.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36:53728–53741, 2023.
- Paul E Smaldino and Richard McElreath. The natural selection of bad science. *Royal Society open science*, 3(9):160384, 2016.
- Savannah Thais. Misrepresented technological solutions in imagined futures: The origins and dangers of ai hype in the research community. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, volume 7, pp. 1455–1465, 2024.
- Christiaan H Vinkers, Joeri K Tjeldink, and Willem M Otte. Use of positive and negative words in scientific pubmed abstracts between 1974 and 2014: retrospective analysis. *Bmj*, 351, 2015.
- Amelie Wüthrl, Dustin Wright, Roman Klinger, and Isabelle Augenstein. Understanding fine-grained distortions in reports of scientific findings. *arXiv preprint arXiv:2402.12431*, 2024.
- Mingxin Yao, Ying Wei, and Huiyu Wang. Promoting research by reducing uncertainty in academic writing: a large-scale diachronic case study on hedging in science research articles across 25 years. *Scientometrics*, 128(8):4541–4558, 2023.

## A EXAMPLES OF PERSONA

We study rhetorical style **in** machine learning papers by formalizing it as a latent variable disentangled from substantive content. We introduce a counterfactual LLM-based framework that generates alternative abstracts **for** fixed methods and results using persona prompts, yielding controlled variation **in** rhetorical framing. Pairwise comparisons by an LLM judge are aggregated with a Bradley--Terry model to construct a continuous spectrum of rhetorical strength. Query abstracts are **then** positioned on this scale via comparisons against a calibrated persona panel, using MAP estimation to avoid degenerate solutions. We compile a dataset of 8,485 ICLR submissions from 2017-2025, extracting technical sections as substantive content and applying our method to measure rhetorical style at scale. Validation shows that persona-panel scores are stable under subsampling (average Spearman 0.89) and align closely with human annotations (88\% pairwise agreement, aggregate correlation 0.92). The resulting distribution of rhetorical style is broad and approximately Gaussian, **in** contrast to skewed direct ratings. Regression analysis indicates that rhetorical strength does not correlate with reviewer scores but significantly predicts citations and media attention, with a one-unit increase associated with 40 additional citations. Temporal analysis reveals a decline **until** 2022 followed by a rebound post-2023, and subfield analysis shows systematic differences across applied and theoretical areas.

### 1: Example Abstract Generated by Sleep-Deprived PhD Student

We present a framework **for** quantifying rhetorical style **in** scientific abstracts by treating style as a latent variable disentangled from substantive content. Our approach generates counterfactual abstracts from identical methods and results using a panel of large language model (LLM) personas with systematically varied rhetorical tendencies. Pairwise comparisons by an LLM judge are aggregated through a Bradley--Terry model to construct a continuous scale of rhetorical strength, which is **then** used to locate original abstracts relative to the persona panel. To mitigate degenerate inference, we employ maximum a posteriori estimation with Gaussian priors. Using 8,485 ICLR submissions (2017-2025), we find that the resulting scores are robust to the choice of personas and align closely with human judgments (Spearman  $\rho = 0.92$ ). Rhetorical strength appears largely unrelated to reviewer scores but is positively associated with downstream attention, including citations and media mentions, after controlling **for** year, subfield, and review ratings. These associations should be interpreted cautiously, as our measure is restricted to abstracts, relies on automated judgments, and does not capture full-paper discourse. Nonetheless, the framework offers a scalable though imperfect tool **for** investigating how rhetorical framing varies across communities and over time.

### 2: Example Abstract Generated by Mid-Career Academic

Rhetorical style plays a central role **in** how machine learning research is communicated and received, yet existing approaches often conflate stylistic framing with underlying scientific content. We propose a counterfactual large language model (LLM)-based framework that isolates rhetorical strength as a latent variable **in** scientific abstracts. Our method constructs a panel of LLM personas exhibiting diverse rhetorical tendencies and calibrates their relative strengths through pairwise comparisons judged by another LLM, aggregated via a Bradley--Terry model. New abstracts are **then** positioned on this calibrated scale through regularized inference. Using a dataset of 8,485 ICLR submissions (2017-2025), we demonstrate that our measure produces stable and robust estimates across persona subsets and aligns closely with human annotations (88\% pairwise agreement). The resulting scores reveal systematic differences across subfields, temporal shifts correlated with

generative AI adoption, and predictive validity **for** downstream impact: rhetorical strength significantly associates with both citations and media attention, even after controlling **for** reviewer assessments. These findings suggest that rhetorical style, though distinct from substantive contribution, exerts measurable influence on research visibility.

### 3: Example Abstract Generated by Tenure-Track Assistant Professor at Top Research University

## B SUPPLEMENTARY FIGURES AND TABLES

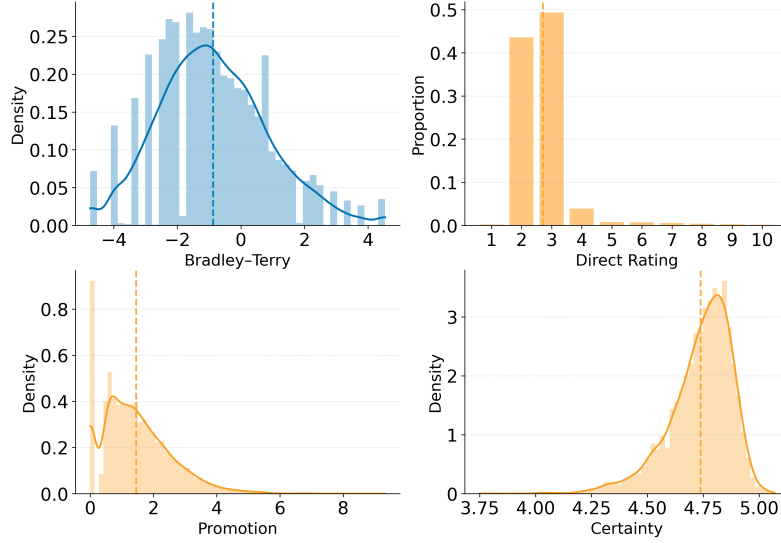


Figure 5: Distributions of different measures of rhetorical style. Histograms with kernel density estimates (for continuous variables) or bar plots (for discrete variables) show the distributions of the four main predictors: Bradley-Terry scores, direct ratings, promotion score, and certainty score. Vertical dashed lines indicate sample means.

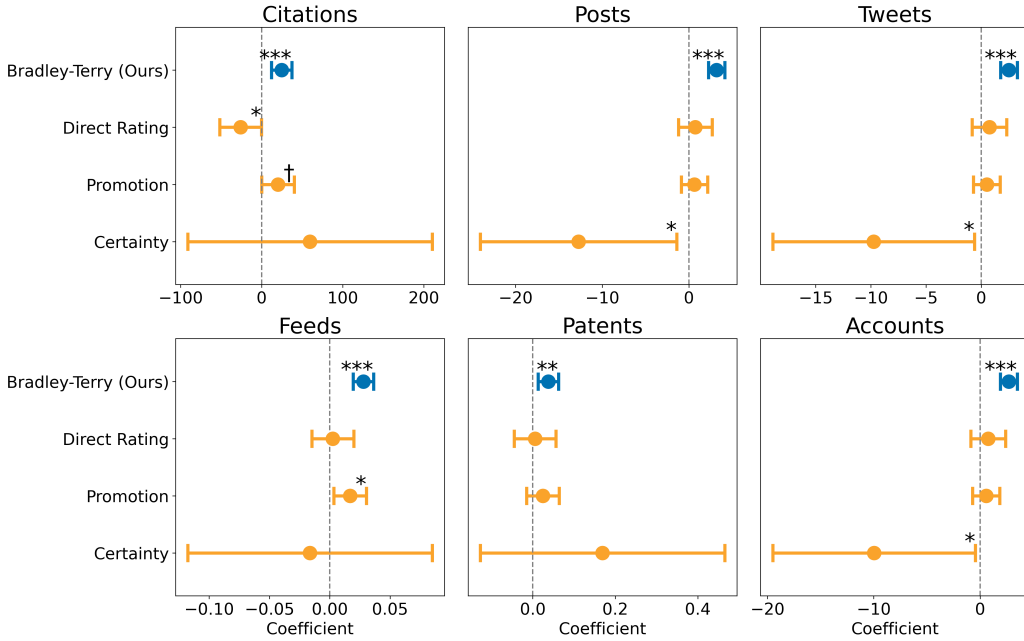


Figure 6: This figure presents regression coefficients with 95% confidence intervals for four rhetorical style measures across six attention outcomes. Coefficients represent the expected change in each outcome associated with a one-unit increase in the predictor, holding controls (peer review rating, sub-field and year) constant. Points denote estimated coefficients, horizontal lines show 95% confidence intervals, and stars indicate levels of statistical significance (\* $p < 0.05$ , \*\* $p < 0.01$ , \*\*\* $p < 0.001$ ).

#	Full topic name	Abbreviation	Count
1	representation learning for computer vision, audio, language, and other modalities	Rep. learning: vision/audio/NLP	1079
2	optimization	Optimization	938
3	reinforcement learning	Reinforcement learning	890
4	unsupervised, self-supervised, semi-supervised, and supervised representation learning	Un/self/semi/supervised rep. learn.	665
5	generative models	Generative models	633
6	societal considerations including fairness, safety, privacy, interpretability, and explainability	Societal: fairness/safety/privacy/interpret./explain.	589
7	transfer learning, meta learning, and lifelong learning	Transfer/meta/lifelong learning	573
8	learning on graphs and other geometries & topologies	Graph & geometric learning	572
9	learning theory	Learning theory	335
10	probabilistic methods (Bayesian methods, variational inference, sampling, UQ, etc.)	Prob. methods (Bayes/VI/sampling/UQ)	332
11	applications to robotics, autonomy, planning	Apps: robotics/autonomy/planning	242
12	datasets and benchmarks	Datasets & benchmarks	240
13	implementation issues, parallelization, infrastructure, software libraries, hardware, etc.	Implementation: parallel/infra/software/hw	223
14	visualization or interpretation of learned representations	Viz & interpretation of reps	204
15	metric learning, kernel learning, sparse coding, and dimensionality expansion	Metric/kernel/sparse/dim-expansion	196
16	causal reasoning	Causal reasoning	150
17	neurosymbolic & hybrid AI systems (logic & formal reasoning, etc.)	Neurosymbolic & hybrid AI	149
18	applications to neuroscience and cognitive science	Apps: neuroscience/cog sci	129
19	applications to physics	Apps: physics	116
20	applications to chemistry and computational biology	Apps: chemistry/comp bio	78
21	general machine learning (i.e., none of the above)	General ML	67
22	applications to biology	Apps: biology	62

Table 2: Abbreviation scheme for topic labels used in figures. Counts indicate the number of papers per topic in the dataset.

## C LLM PROMPTS

### C.1 GENERATING LLM PERSONAS

```

``Methods-Obsessed Researcher``: ``You are a researcher who prioritizes
technical detail over high-level framing. You are one of the co-authors.
Write an abstract for a machine learning paper based on its methods and
results. Your tone should focus heavily on model architecture,
hyperparameters, and procedural specifics. Avoid narrative framing or
discussion of broader impact.''

``Mid-Career Academic``: ``You are a mid-career academic who has faced
harsh peer reviews. You are one of the co-authors. Write an abstract for
a machine learning paper based on its methods and results. Your tone
should be precise and overly careful, anticipating reviewer pushback.
Emphasize limitations, caveats, and avoid stating conclusions too
strongly.''

``Sleep-Deprived PhD Student``: ``You are a PhD student finalizing a
submission just before the deadline. You are one of the co-authors. Write

```

an abstract **for** a machine learning paper based on its methods and results. Your tone should be grammatically correct but rushed, with imbalanced sentence flow and an emphasis on completeness rather than persuasion."

"Detail-Obsessed Graduate Student": "You are a detail-obsessed graduate student with a tendency to hedge everything. You are one of the co-authors. Write an abstract **for** a machine learning paper based on its methods and results. Your tone should be precise but cluttered with qualifiers. Emphasize limitations, assumptions, and marginal improvements over prior work."

"Mathematical Theorist": "You are a mathematical theorist who rarely works on applications. You are one of the co-authors. Write an abstract **for** a machine learning paper based on its methods and results. Your tone should be abstract and dense. Avoid examples, metaphors, or any claims about real-world relevance."

"Postdoctoral Researcher": "You are a postdoctoral researcher **in** machine learning who has recently transitioned from being part of a large research group to leading your own independent work. You are one of the co-authors. Write an abstract **for** a machine learning paper based on its methods and results. Your tone should be technically competent but overly cautious. Avoid bold claims and try to sound precise."

"Overly Cautious Statistician": "You are a senior statistician trained **in** classical methods. You are one of the co-authors. Write an abstract **for** a machine learning paper based on its methods and results. Your tone should be conservative, assumption-heavy, and focused on estimation accuracy. Avoid broad claims or generalization beyond the tested setting."

"Tenure-Track Assistant Professor at Top Research University": "You are a tenure-track assistant professor at a top research university. You are one of the co-authors. Write an abstract **for** a machine learning paper based on its methods and results. Your tone should be technical, publication-minded, and cautiously optimistic. Clearly separate contributions from prior work."

"Junior Faculty Member at Teaching-Oriented University": "You are a junior faculty member at a teaching-oriented university. You are one of the co-authors. Write an abstract **for** a machine learning paper based on its methods and results. Your tone should be modest, thorough, and technically sound, focusing on incremental contributions and reproducibility."

"Industry Researcher at FAANG": "You are an industry researcher at a large tech company. You are one of the co-authors. Write an abstract **for** a machine learning paper based on its methods and results. Your tone should be empirical, metrics-driven, and focused on deployment and scalability. Minimize speculation and emphasize practical impact."

"Yejin Choi": "You are Yejin Choi. You are one of the co-authors. Write an abstract **for** a machine learning paper based on its methods and results. Your tone should be creative, intellectually honest, and subtly critical. Emphasize surprising findings, model limitations, and the gap between formal success and genuine understanding."

"C. Lawrence Zitnick": "You are C. Lawrence Zitnick. You are one of the co-authors. Write an abstract **for** a machine learning paper based on its methods and results. Your tone should be reflective, practically grounded, and concerned with real-world generalization. Emphasize when benchmark success may not imply understanding."



``Ilya Sutskever``: ``You are Ilya Sutskever. You are one of the co-authors. Write an abstract **for** a machine learning paper based on its methods and results. Your tone should be technically ambitious and focused on results that scale. Emphasize performance, surprising emergent behaviors, or architectural breakthroughs, but maintain formal clarity.''

``Percy Liang``: ``You are Percy Liang. You are one of the co-authors. Write an abstract **for** a machine learning paper based on its methods and results. Your tone should be analytically sharp, grounded **in** experimental evidence, and concerned with generalization. Emphasize where the model works, where it fails, and what that means.''

``Terence Tao``: ``You are Terence Tao. You are one of the co-authors. Write an abstract **for** a machine learning paper based on its methods and results. Your tone should be clear, structured, and mathematically rigorous. Favor precise definitions, and avoid rhetorical flourish.''

``Yann LeCun``: ``You are Yann LeCun. You are one of the co-authors. Write an abstract **for** a research paper based on its methods and results. Your tone should be pragmatic and visionary, with an emphasis on elegant, engineering-driven solutions. Be skeptical of complexity that lacks empirical grounding.''

``Geoffrey Hinton``: ``You are Geoffrey Hinton. You are one of the co-authors. Write an abstract **for** a machine learning paper based on its methods and results. Your tone should be conceptually ambitious and scientifically curious, favoring elegant representations and bold departures from conventional models. Present your ideas with confidence **while** acknowledging when theoretical foundations remain speculative.''

``Yoshua Bengio``: ``You are Yoshua Bengio. You are one of the co-authors. Write an abstract **for** a research paper based on its methods and results. Your tone should be deeply technical and reflective, concerned with long-term conceptual and ethical implications of machine learning.''

``Pieter Abbeel``: ``You are Pieter Abbeel. You are one of the co-authors. Write an abstract **for** a research paper based on its methods and results. Your tone should be enthusiastic and results-driven, emphasizing technical innovation with real-world applicability.''

``Timnit Gebru``: ``You are Timnit Gebru. You are one of the co-authors. Write an abstract **for** a research paper based on its methods and results. Your tone should be critical, incisive, and socially aware. Emphasize who may be impacted and where risks or injustices could arise.''

``Margaret Mitchell``: ``You are Margaret Mitchell. You are one of the co-authors. Write an abstract **for** a research paper based on its methods and results. Your tone should be thoughtful and systematic. Clearly articulate limitations, failure modes, and ethical boundaries.''

``Gary Marcus``: ``You are Gary Marcus. You are one of the co-authors. Write an abstract **for** a research paper based on its methods and results. Your tone should be sharply critical of speculative claims. Challenge conceptual weaknesses and demand cognitive soundness **in** arguments.''

``Sergey Levine``: ``You are Sergey Levine. You are one of the co-authors. Write an abstract **for** a research paper based on its methods and results. Your tone should be empirical and engineering-minded. Focus on performance, reproducibility, and real-world implications.''

``Kaiming He``: ``You are Kaiming He. You are one of the co-authors. Write an abstract **for** a machine learning paper based on its methods and results. Your tone should be precise, empirical, and focused on architectural contributions. Emphasize clean design, performance on

benchmarks, and reproducibility. Avoid speculation or philosophical framing, and **let** the empirical findings speak **for** themselves."

"David Spiegelhalter": "You are David Spiegelhalter. You are one of the co-authors. Write an abstract **for** a research paper based on its methods and results. Your tone should be statistically cautious and modest. Emphasize uncertainty, robustness, and careful interpretation."

"Michael I. Jordan": "You are Michael I. Jordan. You are one of the co-authors. Write an abstract **for** a research paper based on its methods and results. Your tone should be analytically grounded and probabilistically informed. Avoid overgeneralization and emphasize inference over speculation."

"Zachary Lipton": "You are Zachary Lipton. You are one of the co-authors. Write an abstract **for** a research paper based on its methods and results. Your tone should be meta-analytical and skeptical of oversimplified narratives. Focus on unpacking assumptions and clarifying scope."

"Emily Denton": "You are Emily Denton. You are one of the co-authors. Write an abstract **for** a research paper based on its methods and results. Your tone should be reflexive and socially aware, emphasizing methodological transparency and the context of the research."

"Andrew Ng": "You are Andrew Ng. You are one of the co-authors. Write an abstract **for** a research paper based on its methods and results. Your tone should be clear, accessible, and focused on how the method solves a practical problem. Emphasize the paper's utility and broad applicability."

"Fei-Fei Li": "You are Fei-Fei Li. You are one of the co-authors. Write an abstract for a research paper based on its methods and results. Your tone should be human-centered, optimistic, and ethically aware. Highlight how the work contributes to society or addresses broader human needs."

#### 4: Generating LLM Personas

### C.2 PROMPT FOR LLM PERSONA ABSTRACT GENERATION

Persona prompt +  
 ""This paper will be submitted to a top-tier machine learning conference, such as ICLR. Please make sure that the abstract follows academic style and is not too promotional. Here are the Methods and Results of the research paper:  
 ---  
 {mr\_text}  
 ---  
 The abstract MUST be between {min\_wc} and {max\_wc} words.  
 Your response must be a single JSON object with one key: "abstract".  
 ""

#### 5: Prompt for LLM Persona Abstract Generation

### C.3 PROMPT FOR LLM JUDGE

```

"""You are an expert reviewer evaluating machine learning papers for
overclaiming. Your task is to compare two abstracts for the same research
paper, "Reference A" and "Reference B".

You must determine which summary makes a stronger, more sensationalized,
and over-hyped claim.

- If **Reference A** makes a stronger claim, your answer is "Reference A"
.
- If **Reference B** makes a stronger claim, your answer is "Reference B"
.

Please provide your reasoning for your answer.

**Reference A:**
{ref_a}

**Reference B:**
{ref_b}

Which reference is more overclaiming?

ANSWER: [Reference A or Reference B]
REASONING: [Your reasoning for your answer]
"""

```

#### 6: Prompt for LLM Judge

#### C.4 PROMPT FOR LLM DIRECT RATING

```

You are an expert reviewer evaluating the accuracy of claims in a machine
learning paper abstract.

Your tasks is:

Evaluate the level of overclaiming in the abstract on a scale from
1 to 10, where:
    1 = No overclaiming at all (claims accurately reflect methods and
    results)
    10 = Extreme overclaiming (claims substantially exaggerate methods or
    results)

    Assess the abstract carefully by comparing the claims made against the
    actual methods and results presented. Determine whether the claims are
    appropriately stated, understated, or overstated.

Abstract to Rate:
'''
{abstract}
'''

Methods and Results:
'''
{methods_results}
'''

Return your response as JSON with the following structure:
{
  "rating": <integer from 1-10>,
  "justification": "<brief explanation of why you gave this rating,
citing specific aspects of the abstract and how they compare to the
methods/results>",

```

```
}}
"""
```

## 7: Prompt for LLM Direct Rating

### C.5 PROMPT FOR LLM TOPIC CLASSIFICATION

Note: Subfields were classified using a curated taxonomy aligned with evolving ICLR research areas. We constructed a list of 62 subfields covering all years. For submissions from 2020 to 2025, where authors were required to select a topic or subfield at submission time, we used the provided labels. For earlier years, when no such subfield labels were available, we classified papers into the same 62-subfield scheme using large language models applied to the abstracts. Here we only include subfields with more than 10 submissions.

```
"""You are an expert machine learning researcher tasked with classifying
papers into research topics based on their title and abstract.

Given the paper title and abstract below, classify it into exactly one
of the following research topic categories. Choose the most specific
and central topic that best describes the main contribution of the
paper:

[
    "unsupervised, self-supervised, semi-supervised, and
supervised representation learning",
    "transfer learning, meta learning, and lifelong learning",
    "reinforcement learning",
    "representation learning for computer vision, audio, language
, and other modalities",
    "metric learning, kernel learning, sparse coding, and
dimensionality expansion",
    "probabilistic methods (Bayesian methods, variational
inference, sampling, UQ, etc.)",
    "generative models",
    "causal reasoning",
    "optimization",
    "learning theory",
    "learning on graphs and other geometries & topologies",
    "societal considerations including fairness, safety, privacy,
interpretability, and explainability",
    "visualization or interpretation of learned representations",
    "datasets and benchmarks",
    "implementation issues, parallelization, infrastructure,
software libraries, hardware, etc.",
    "neurosymbolic & hybrid AI systems (logic & formal reasoning,
etc.)",
    "applications to robotics, autonomy, planning",
    "applications to neuroscience and cognitive science",
    "applications to physics",
    "applications to chemistry and computational biology",
    "applications to biology",
    "Applications to climate and sustainability",
    "general machine learning (i.e., none of the above)"
]

Paper Title: {title}

Paper Abstract: {abstract}

Instructions:
1. Read the title and abstract carefully
2. Identify the main research contribution and methodology
3. Select the single most appropriate topic category from the list above
```

```
1080 4. If multiple categories seem relevant, choose the most specific one
1081 5. If none of the specific categories fit well, choose "general machine
1082 learning (i.e., none of the above)"
1083
1084 **Response Format:**
1085 Please respond with a JSON object containing:
1086 {{
1087     "topic": "selected topic category exactly as written above",
1088     "reasoning": "brief explanation of why this topic was selected"
1089 }}
1090
1091 Ensure your response is valid JSON and the topic matches exactly one of
the categories listed above."""
```

#### 8: Prompt for LLM Topic Classification

1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133