

Search Arena Meets Nuggets: Towards Explanations and Diagnostics in the Evaluation of LLM Responses

Anonymous ACL submission

Abstract

Battles, or side-by-side comparisons in so-called arenas that elicit human preferences, are used to assess the large language model (LLM) output quality, and have recently been extended to retrieval-augmented generation (RAG) systems. Although, battles mark progress in evaluation, they have two key limitations for complex information-seeking queries: they are neither explanatory nor diagnostic. On the other hand, nugget-based evaluation, that decomposes long-form answers into atomic facts and highlights necessary parts in an answer, has emerged as a promising strategy for RAG evaluation. In this work, we employ AutoNuggetizer, a nugget-based framework, to analyze $\sim 5K$ Search Arena battles from LMArena by automatically generating and assigning nuggets, converting each model response into a quantitative score. We observe strong alignment between nugget-based Elo rankings and human preferences, with Kendall’s τ of 0.71 and Spearman’s ρ of 0.88, exceeding the corresponding alignment achieved by LLM-as-a-judge evaluation (0.64 and 0.79, respectively), while substantially reducing the number of preference inversions. Furthermore, we provide in-depth analyses including inversions, nugget quality, and shared-blindness effects. All our code and datasets will be released publicly upon paper acceptance.

1 Introduction

The notion of “battles”, or side-by-side comparisons of responses from large language models (LLMs), has become a popular method for evaluating their quality (Zheng et al., 2023; Chiang et al., 2024). In the “arena” setup, users are shown two LLM outputs and asked to indicate which one they prefer. This approach was popularized by LMSYS through MT-Bench (Zheng et al., 2023) and later expanded into the Chatbot Arena (Chiang et al., 2024). The popularity of these arenas has made

them a key marketing tool when launching new LLMs from companies such as Google, OpenAI, and Meta, who regularly tout leaderboard rankings on Chatbot Arena in model releases. Recently, arena-based evaluations have been extended to a variety of domains, including AI agents (Yekollu et al., 2024), vision and image generation (Lu et al., 2024; Jiang et al., 2024), multilingual generation (Thakur et al., 2025a), and even GitHub pull requests (Wang et al., 2025).

Battles were extended to search-augmented LLMs in the Search Arena (Miroyan et al., 2025). Unlike the original setup, which focused on “closed-book” LLM responses, Search Arena evaluates retrieval-augmented generation (RAG) systems in two stages: first, retrieving relevant web-sourced documents, and next using them to generate long-form answers with citations using LLMs (Pradeep et al., 2025a; Han et al., 2024).

While such side-by-side comparisons enable the evaluation of search-augmented LLM-based systems at scale, we see them having at least two drawbacks: they are neither explanatory nor diagnostic, especially in scenarios where determining the better answer is not straightforward. It would be desirable for an evaluation to (at least attempt to) explain *why* a user might have preferred one response over another.

We hypothesize that the nugget evaluation methodology (Pradeep et al., 2024, 2025b) can be adapted to address these two limitations for complex information-seeking queries. The core idea is to measure answer quality based on the recall of information nuggets, or atomic facts, that should appear in high-quality responses. In AutoNuggetizer (Pradeep et al., 2024, 2025b), this process can be fully automated using LLMs, breaking down a long-form model response into a quantitative score.

In our work, we adapt the AutoNuggetizer

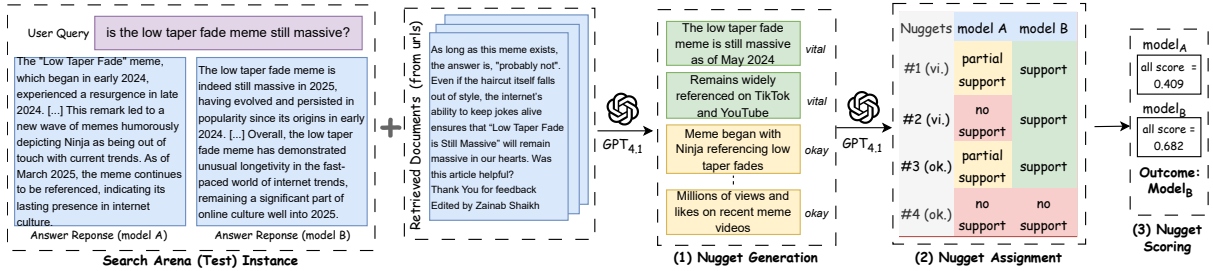


Figure 1: An example from Search Arena illustrating both nugget generation and assignment. First, GPT_{4.1} generates nuggets based on the query, retrieved chunks from URL contents, and the responses from both models (A and B). Each nugget is then labeled with an importance level—either “vital” or “okay”. Next, GPT_{4.1} evaluates whether each model supports each nugget, assigning one of three labels: “support”, “partial support”, or “no support”. Finally, these support judgments are scored and aggregated to determine the outcome (the model with the higher score is preferred).

(Pradeep et al., 2024) framework to analyze approximately 5K battles in the Search Arena in a fully automatic manner (see Figure 1), eliminating the need for cumbersome human judgments. The framework includes two stages: (1) *nugget generation*: eliciting nuggets from model answers and scraped documents from cited URLs, and (2) *nugget assignment*: evaluating whether each answer supports a nugget or fact. Our results show that human preferences correlate well with the distribution of nugget scores, achieving a very high alignment between Elo rankings of nugget-based evaluation and human preferences (Kendall’s τ of 0.71 and Spearman’s ρ of 0.88). Furthermore, our extended analysis diagnose the score inversions, nugget quality and potential shared-blindness in our evaluation setup. To summarize, our contributions are as follows:

- **Extending AutoNuggetizer to an arena setting with live search.** We adapt Search Arena to AutoNuggetizer by scraping and chunking the dataset URLs to form a retrieval corpus, then apply nuggetization to head-to-head battles where LLMs have live web-search access—moving beyond prior nugget evaluations that assume a fixed/static corpus.
- **Rigorous analysis of nuggetization and preference inversions.** We audit nugget factuality and diversity (lexical and semantic) and dissect “preference inversions,” showing how query type (e.g., ambiguous vs. factoid) and language systematically affect nugget-based preferences;
- **Comparisons to alternative evaluators.** We benchmark nugget-based preferences against LLM-as-a-judge, surface-form metrics (e.g., ROUGE-style overlap) and judging with factors like fluency and generation style, highlighting

when each method succeeds or fails.

2 Related Work

Nugget-based evaluation. First introduced in the TREC QA Track in 2003 (Voorhees, 2003a,b), the nugget-based evaluation methodology focuses on identifying essential atomic facts—called nuggets—that are relevant to a given question. This methodology was later extended to tasks like summarization and broader conceptions of question answering (Nenkova and Passonneau, 2004; Lin and Demner-Fushman, 2006b; Dang and Lin, 2007; Lin and Zhang, 2007), and researchers have explored automation to improve its scalability (Lin and Demner-Fushman, 2005, 2006a; Pavlu et al., 2012).

The recent emergence of LLMs has enabled automated, reliable nugget-based evaluation (Pradeep et al., 2024; Alaofi et al., 2024; Pradeep et al., 2025b; Thakur et al., 2025b; Abbasiantaeb et al., 2025). Several RAG evaluation frameworks—such as FactScore (Min et al., 2023), RUBRIC (Farzi and Dietz, 2024), and others (Arabzadeh and Clarke, 2024; Mayfield et al., 2024)—incorporate the nugget concept, although most of these proposed approaches are either not validated or primarily validated on traditional ad hoc retrieval, and hence their applicability to long-form answers is unclear. We refer readers to (Pradeep et al., 2025b) for a more detailed discussion of related work. In this work, we focus on the AutoNuggetizer framework from (Pradeep et al., 2024), and apply it to the side-by-side comparisons of LLM responses in the Search Arena.

Related arena benchmarks. The Search Arena (Miroyan et al., 2025) by LMArena is a

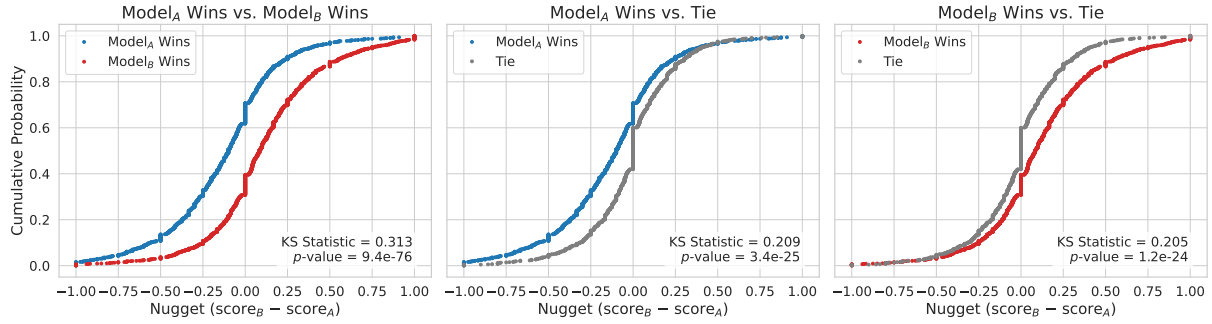


Figure 2: Empirical cumulative distribution functions (CDFs) comparing nugget score differences ($\text{score}_B - \text{score}_A$) across human vote categories. Each subplot shows a Kolmogorov-Smirnov (K-S) test between two groups: (left) model_A wins vs. model_B wins, (center) model_A wins vs. tie, and (right) model_B wins vs. tie. The K-S statistic and corresponding p -value are annotated in each plot, quantifying the distributional differences between groups.

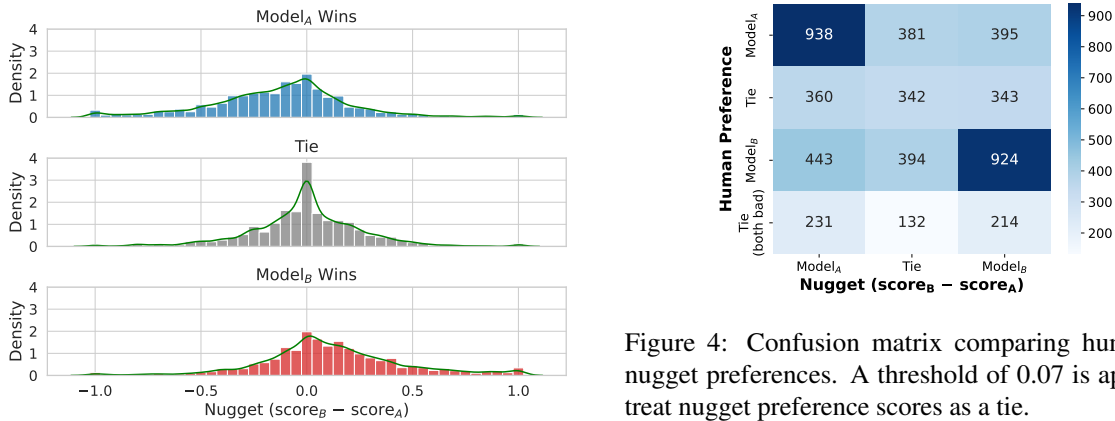


Figure 3: Empirical probability density function (PDF) of nugget score differences ($\text{score}_B - \text{score}_A$) grouped by human preference category: model_A wins, tie, or model_B wins. A separate Kernel Density Estimation (KDE) with bandwidth 0.5 is fitted for each group.

Figure 4: Confusion matrix comparing human and nugget preferences. A threshold of 0.07 is applied to treat nugget preference scores as a tie.

benchmark that evaluates LLMs with access to a live web-search tool. Other notable efforts include the MTEB Arena (Hugging Face, 2023), which extends the Massive Text Embedding Benchmark (MTEB) framework (Muennighoff et al., 2023) to head-to-head evaluation across embedding models, and Ragnarök (Pradeep et al., 2025a), which offered a head-to-head RAG evaluation framework on the MS MARCO V2.1 collection in the TREC 2024 RAG Track.

3 Experimental Design

Search Arena overview. Search Arena (Miroyan et al., 2025) is a crowd-sourced platform that evaluates search-augmented LLMs via side-by-side human-preference judgments (Chiang et al., 2024). The V1 dataset¹ includes 7K battles between two RAG-oriented systems (model_A and model_B; e.g.,

¹<https://huggingface.co/datasets/lmarena-ai/search-arena-v1-7k>

Gemini-2.5-Pro-Grounding vs. Perplexity-Sonar-Reasoning-Pro). For each battle, annotators choose one of the four outcomes: model_A wins, model_B wins, good tie (both responses are equally good), or bad tie (both responses are equally bad). Search result URLs used during generation are available for $\sim 6.7\text{K}$ battles, totaling $\sim 80\text{K}$ unique URLs. The Search Arena dataset includes both single- and multi-turn battles. We restrict our analysis to single-turn battles only—5,103 instances where the system returns a single response—because overall votes in multi-turn settings do not reliably capture per-query preferences, which is what Auto-Nuggetizer evaluates.

Search Arena also contains battles for several non-English languages, e.g., Chinese or Russian. Non-English languages collectively account for less than 40% of the dataset, with English comprising the remaining majority. Detailed statistics for single-turn battles used in this work are presented in Appendix A.1. Queries in Search Arena vary widely, ranging from long code snippets to prompts that demand complex reasoning or exhibit ambiguity and vagueness. We show a few examples

of queries from the Search Arena dataset in Appendix A.2.

Corpus generation. To evaluate LLM responses in the absence of ground-truth answers, we use the search result URLs provided in the dataset, collected from each system response as relevant sources of information. We begin by constructing a corpus from the 47K unique URLs associated with single-turn battles. This process involves downloading the contents of each URL, extracting the main textual content, and segmenting the text into chunks of ten sentences with an overlap of two sentences, using the `xx_sent_ud_sm` model from `spaCy`.²

Once the corpus is prepared, we encode the chunks and the query prompts utilizing the `BAAI/bge-m3` model.³ We retrieve the top 50 most relevant chunks for each query via the cosine similarity between the chunk and query embeddings using Pyserini’s FAISS indexing and search (Lin et al., 2021). Notably, both the chunking and encoding models support multilingual corpora, ensuring a robust language coverage.

Nugget evaluation. Nugget generation creates atomic facts that highlight the essential information required in a RAG answer, and assignment categorizes their support level for the model response. Following (Pradeep et al., 2024), we use the AutoNuggetizer tool in the nuggetizer code repository⁴ to generate and assign information nuggets to model responses. As shown in Figure 1, there are two steps in nugget generation and assignment:

1. **Nugget generation:** For each prompt extracted from the dataset, we construct a request to AutoNuggetizer that includes the query (i.e., the prompt itself), along with relevant chunks retrieved from our created corpus, ordered by relevance and responses from each model, inserted in a random order to mitigate the positional bias. We include model responses for two key reasons. First, approximately 5% of the battles do not contain any URLs. Second, even when URLs are present, about 16% of them yield 100 bytes or less of content after scraping⁵. These

cases make the LLM responses a valuable fallback source of information for nugget generation. The AutoNuggetizer tool then processes the request and identifies nuggets that are relevant to the query from the retrieved chunks and the provided LLM responses. Furthermore, each nugget is assigned an importance label: “vital” or “okay”, reflecting its relevance to the input query.

2. **Nugget assignment:** Once nuggets and their importance labels are generated (from the previous step), we use AutoNuggetizer to assign them to model responses, determining whether each nugget is supported in the answer. This step categorizes each nugget into “supported”, “partially supported”, or “not supported”. We adopt the “All Score” metric, that achieves the highest recall by counting nuggets of all importance and support levels⁶. We emphasize that while the AutoNuggetizer framework supports different degrees of manual intervention, in this work, we are running the entire evaluation pipeline end-to-end automatically.

4 Experimental Results

Unless stated otherwise, all experiments in this paper are conducted using GPT_{4.1}, with a knowledge cutoff of June 2024, as the underlying LLM used by the AutoNuggetizer via Microsoft Azure. Out of the 5,103 single-turn battles in Search Arena, five were excluded from our analysis due to issues such as Azure content filtering, invalid output formats, or other nugget generation failures. On average, each single-turn battle full evaluation (comprising both nuggetization and assignment) requires approximately 2–3 seconds when executed using the Azure OpenAI API. We set a maximum of 30 nuggets per battle, though this limit is rarely reached (only in 67 battles). When it is reached, only nuggets labeled as okay are removed, while no vital nuggets are discarded. On average, about ~12.5 nuggets are generated per battle.

Main Results. Figure 3 presents our main results, the probability densities of nugget score differences (score_B – score_A) conditioned on the human preference judgment (i.e., the battle outcomes). On the top, we show the distribution when model_A

²https://spacy.io/models/xx#xx_sent_ud_sm

³<https://huggingface.co/BAAI/bge-m3>

⁴<https://github.com/castorini/nuggetizer>

⁵Invalid cases occurs due to issues such as cookie or JavaScript requirements, invalid or expired links, geo-blocking, and similar obstacles.

⁶We find that “Strict Vital”, which was the primary metric used in the TREC 2024 RAG Track (Pradeep et al., 2025b), is too strict for our use case, particularly when only a small number of nuggets are available.

wins; on the bottom, we show the distribution when model_B wins; and in the middle, ties. Battles where the output of both models is considered to be equally bad are excluded from the distributions.

These results appear to support our hypothesis that nugget scores correlate with human preferences. In the case where model_A wins (top row), the distribution skews to the left (negative values), indicating that model_A typically gets higher nugget scores than model_B. Conversely, when model_B wins (bottom row), the distribution skews to the right (positive values), suggesting that model_B generally obtains a higher nugget score. When the human indicates a tie (middle row), the distribution peaks around zero, as expected, indicating similar nugget scores between models.

Statistical Tests. To analyze the statistical differences among these three conditional distributions, we performed pairwise Kolmogorov-Smirnov (K-S) tests. As shown in Figure 2, the K-S statistic values range from 0.205 to 0.313, with p -values of $1.2e^{-24}$ or lower, indicating that all three distributions differ significantly from one another (i.e., we have high confidence that these samples were drawn from different distributions). These findings validate our hypothesis that nugget score differences align with human preferences, reinforcing the potential of nugget-based metrics as reliable evaluators of model quality in head-to-head evaluations.

Confusion Matrix. Figure 4 presents a confusion matrix that compares the distribution of human preferences (rows) in Search Arena against “nugget preferences” (columns). For “nugget preference”, we use a threshold of 0.07, meaning that when the nugget score difference between the two model outputs falls within ± 0.07 , the comparison is considered a tie (The threshold was selected by sweeping values between 0.05 and 0.15 in increments of 0.01). The threshold of 0.07 closely reflects an equal distribution of model_A wins, model_B wins, and ties when the human preference is a tie (second row in Figure 4). The diagonal cells in this confusion matrix reveal the instances where nugget preferences align with human preferences. Conversely, off-diagonal cells illustrate the types and frequencies of disagreements between the human and nugget scores.

In particular, the nugget-based evaluation prefers model_A in 938 out of 1,714 (54.7%) of the battles where model_A wins the battle (first row in Fig-

Table 1: Inversion percentages and query frequencies across different (a) query categories and (b) languages in the Search Arena dataset.

Category	Inversion (%)	Query Count
(1) Ambiguous	19%	196
(2) Assumptive	18%	28
(3) Multi-faceted	18%	299
(4) Incompleteness	16%	631
(5) Subjective	15%	601
(6) Knowledge-int.	15%	1142
(7) Reasoning-int.	14%	288
(8) Harmful	9%	92

(a)

Language	Inversion (%)	Query Count
(1) German	20%	244
(2) English	17%	3117
(3) Chinese	16%	328
(4) Portuguese	16%	150
(5) Russian	15%	460
(6) French	13%	151
(7) Others	16%	647

(b)

ure 4). Similarly, model_B is preferred in 924 out of 1761 (52.5%) battles where it wins the battle (third row in Figure 4). To further quantify this alignment, we report a weighted Cohen’s κ of 0.30 with quadratic (0, 0.25, 1) weights assigned to the (inversion, tie, identical) labels, respectively. This value remains stable across nugget score thresholds for ties, varying only slightly between 0.29 and 0.31 when thresholds range from 0.05 to 0.15.

Elo Rankings Given the high level of subjectivity inherent in the task, a moderate level of agreement at the individual battle level is expected. Consequently, the overall human- versus nugget-based rankings of the participating LLMs are of primary interest. Following established best practices in generative AI pairwise evaluation (Chiang et al., 2024; Jiang et al., 2024; Chi et al., 2025; Miroyan et al., 2025), we employ the Bradley–Terry (BT) model (Bradley and Terry, 1952), estimated via maximum likelihood estimation (MLE) using logistic regression with bootstrap resampling. Table 5 reports the resulting Elo scores for both human and nugget-based preferences. Kendall’s τ of 0.71 and Spearman’s ρ of 0.88 indicate strong agreement between the two Elo rankings.

In the remainder of this section, we analyze the anti-diagonal cases where nugget-derived and human preferences diverge; compare LLM-as-a-judge

365 preference agreement with human judgments as an
366 alternative to nugget-based preferences; reassign
367 nuggets using an alternative LLM; examine the gener-
368 ated nuggets; and explore other alternatives for
369 nugget-based evaluation of open-ended generation.

370 4.1 Inversion Analysis

371 **Query Classification Analysis.** In this analysis,
372 we use query classification to better understand the
373 cases where nugget preferences and human prefer-
374 ences are not aligned. When the nugget scores
375 and the human prefer opposite sides of a battle, we
376 refer to this situation as a “preference inversion”, or
377 simply inversion. We suspect that inversions might
378 vary across different types of queries. To investi-
379 gate, we followed (Rosset et al., 2024) but used the
380 newer GPT_{4.1} to rate each query on a scale of 0–10
381 across eight different categories. Then, we classify
382 each query into its maximum scoring category or
383 categories (allowing for ties).

384 To further strengthen the category signals, we
385 exclude queries with a maximum score less than
386 seven from this classification. Raw distributions of
387 the query ratings per category and sample queries
388 from each class are available in Appendix A.2.
389 As shown in Table 1 (row a), the portion of in-
390 versions for ambiguous, assumptive, multi-faceted,
391 and incomplete queries is higher than that of subjec-
392 tive, knowledge-, and reasoning-intensive queries.
393 This suggests that inversions are more likely when
394 queries allow for multiple valid interpretations or
395 are under specified.

396 We followed up by manually examining the in-
397 versions for these categories. As a case study, we
398 encountered a query categorized as *ambiguous* with
399 the text “Potatoes”. In our opinion, both model_A
400 and model_B provided relevant responses. How-
401 ever, model_A focused on the historical aspects and
402 nutritional value of potatoes, whereas model_B dis-
403 cussed cooking methods and varieties. The user
404 judge preferred model_B’s answer, while model_A
405 was selected based on the nugget score. The inher-
406 ent ambiguity of the query likely led to this inver-
407 sion, as it permitted various valid interpretations.
408 Overall, the knowledge-intensive class shows the
409 highest preference alignment—58.8% and 55.7%
410 for model_A and model_B wins, respectively (see Fig-
411 ure 10). This finding suggests that nuggetization
412 is most effective for researchy queries requiring
413 retrieval augmentation. Please refer Appendix B.1
414 for further analysis on query classification.

Query Language Analysis. Next, we analyzed
415 the AutoNuggetizer effectiveness across the six
416 most frequent query languages, each representing
417 at least 3% of the dataset. Previously, the Auto-
418 Nuggetizer had only been run on English responses,
419 and there are likely to be language effects in the
420 breakdown of inversions. As shown in Table 1
421 (row b), German exhibits the highest inversion rate
422 (20%), while French shows the lowest (13%). The
423 confusion matrix for German (see Figure 11) re-
424 veals that it has the smallest portion of ties in hu-
425 man preferences, leading to more anti-diagonal dis-
426 agreements. Please refer Appendix B.2 for further
427 analysis on query languages. 428

429 4.2 LLM-as-a-Judge Evaluation

430 To analyze the correlation between human and
431 LLM preferences, we experiment with GPT_{4.1} as
432 a judge. We modify the chain-of-thought prompt
433 provided in (Rackauckas et al., 2024) (refer Ap-
434 pendix E). For each evaluation, we provide the
435 user query along with the two model responses—
436 randomly ordered to mitigate positional bias—as
437 input to GPT_{4.1}. The model is instructed to output
438 its reasoning and final verdict in a structured JSON
439 format.

440 Figure 5 presents the confusion matrix compar-
441 ing human preferences with those of the GPT_{4.1}
442 judge, which yields a weighted (0, 0.25, 1) Cohen’s
443 κ of 0.31. Compared to the nugget-based evalu-
444 ation in Figure 4, we observe stronger alignment
445 between GPT_{4.1} and human judgments for clear
446 winners: 1,137 vs. 938 agreements for model_A,
447 and 1,161 vs. 924 for model_B.

448 However, GPT_{4.1} struggles to identify
449 ties—including cases where both responses are
450 poor—labeling only 4.25% of single-turn queries
451 as ties. This narrow margin for tie predictions leads
452 to a higher frequency of preference inversions in
453 this setting, with 1,102 inversions compared to 817
454 under the nugget-based evaluation. The higher
455 inversion rate, in turn, causes greater variation
456 in the Elo scores. Overall, alignment between
457 LLM-as-a-judge and human preferences is weaker
458 than that between nugget-based evaluation and
459 human judgments, as reflected by lower Kendall’s
460 τ (0.64 vs. 0.71) and Spearman’s ρ (0.79 vs. 0.88).
461 In addition, the free-form nature of LLM explana-
462 tions limits their utility for diagnostic purposes, as
463 they lack structured cues that can guide targeted
464 improvements.

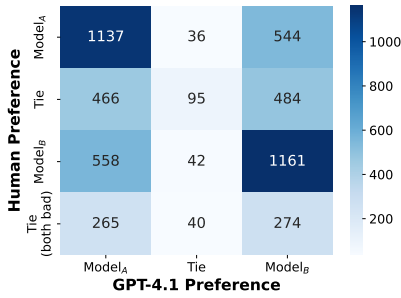


Figure 5: Confusion matrix comparing human and GPT_{4.1} preferences.

Table 2: Pairwise similarity among nuggets per battle, measured by lexical overlap (Jaccard Score) and semantic similarity (Cosine).

Metric	Min	Max	Mean (SD)
Jaccard Score	0.00	0.97	0.06 (0.06)
Cosine Similarity	0.03	0.97	0.39 (0.12)

4.3 No Shared-Blindness in Nugget Extraction and Assignment Evaluation

A potential concern with our evaluation setup could be that using the same model (GPT_{4.1}) for nugget extraction and assignment may introduce *shared blindness* where systematic omissions or potential misjudgments go undetected at both stages. To assess this, we compare nugget assignment from GPT_{4.1} with a different model, Qwen-3-8B (Yang et al., 2025). Nugget assignment with Qwen-3-8B is performed using vLLM (temperature = 0.7) on 4×A6000 GPUs with the same prompt, and the resulting predictions are compared against GPT_{4.1}. The confusion matrix (as shown in Figure 6) yields a weighted Cohen’s κ of 0.69 under quadratic weights (0, 0.25, 1), reflecting substantial agreement between GPT_{4.1} and Qwen-3-8B. This consistency suggests that nugget assignment outcomes are not solely artifacts of a particular LLM such as GPT_{4.1}, thereby mitigating concerns of shared blindness.

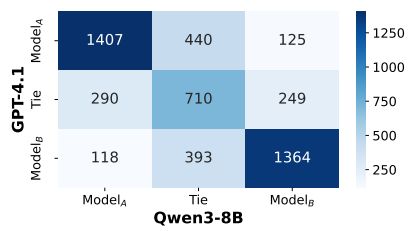


Figure 6: Confusion matrix comparing nugget assignment between GPT_{4.1} and Qwen-3-8B as the judge.

4.4 Nuggets Analysis

In this section, we examine nugget diversity and factual accuracy. To assess the degree of overlap between nuggets generated for each battle, we compute their pairwise similarity. Specifically, we report both lexical overlap, measured using the Jaccard score, and semantic similarity using the SBERT model⁷ to measure cosine similarity between embeddings. As shown in Table 2, nuggets exhibit low lexical overlap (Jaccard: mean = 0.06, SD = 0.06) and moderate semantic similarity (Cosine: mean = 0.39, SD = 0.12). The latter is expected for nuggets from the same battle and should not be interpreted as redundancy, since they address the same topic. Figure 12 shows the precise distribution of the two similarity metrics.

To evaluate factuality, we employ a multilingual natural language inference (NLI) model available on Hugging Face.⁸ The model assesses whether each nugget is entailed by its generation sources, namely the retrieved documents and the generated answers. We find that 99.5% of nuggets are entailed by at least one source, demonstrating that the vast majority are free of hallucinations and are supported by evidence contained in the generation sources.

4.5 Alignment with Other Factors

We analyze alignment between human preferences and surface-level factors such as lexical overlap, fluency, grammar, and readability. For each sentence in a model response, we compute its maximum ROUGE-L F1 against all sentences in the retrieved chunks, then average these maxima across the response. Comparing the higher of the two response-level scores to the human-preference winner yields a confusion matrix with no alignment (Figure 7), indicating that simple lexical-overlap metrics are ineffective for evaluating open-ended generation.

To assess the impact of fluency, grammar, and readability on human judgments, we randomly sampled 1,000 battles with a clear winner (excluding ties) and correlated the difference in scores (A – B) with the human-preference labels (refer Appendix F for more details). As shown in Table 3, both Kendall’s τ and Spearman’s ρ correlations are very weak, suggesting these factors play at most a very minor role in driving human preferences.

⁷<https://huggingface.co/sentence-transformers/paraphrase-multilingual-MiniLM-L12-v2>

⁸<https://huggingface.co/MoritzLaurer/mDeBERTa-v3-base-xnli-multilingual-nli-2mil7>

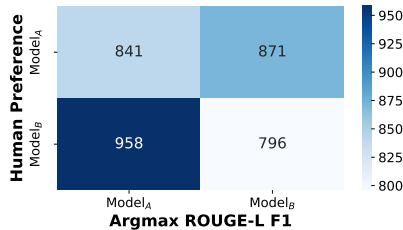


Figure 7: Confusion matrix comparing human preferences vs. argmax ROUGE-L F1 scores, ties excluded.

Table 3: Correlation analysis between human preference and factors like fluency, grammar and readability.

Factor	Kendall’s τ	Spearman’s ρ
Fluency and Coherence	0.070	0.074
Grammar and Syntax	0.007	0.008
Readability and Presentation	0.085	0.089

5 Discussion

In this work, we hypothesize that the nugget evaluation methodology can both *explain* human preferences in side-by-side comparisons and provide *diagnostic* guidance for improving models. Our intuition is simple: humans prefer LLM responses that contain more facts, operationalized as atomic nuggets. With the AutoNuggetizer framework, nugget extraction and scoring are performed automatically. Differences in nugget scores are clearly correlated with human preferences, as illustrated by our density plots.

Our nugget analysis further shows that the generated nuggets are diverse and factually accurate, and that nugget-score assignment is not substitutable with traditional lexical-overlap metrics such as ROUGE-L F1. We empirically show that other factors—fluency, grammar, readability, and presentation—exhibit very weak correlation with human preferences, and infer that current LLMs generally meet user expectations along these dimensions. Consequently, our automatically computed fact-recall metric predicts human preferences in over 50% of cases, underscoring the explanatory power of nugget scores.

Though preliminary, our approach readily supports diagnostic use. Missing nuggets arise from retrieval (relevant docs not surfaced) or modeling (context ignored), suggesting different fixes: strengthen retrieval (e.g., better embeddings) or convert battle outcomes into training signals. This paper is a first step toward nugget-based diagnosis for search-based arena battles.

6 Conclusion

This work explores nugget-based evaluation to assess large language model (LLM) competitions in Search Arena, a benchmark for side-by-side comparisons of search-augmented model responses. By generating and scoring atomic facts (nuggets), we present a more interpretable and diagnostic alternative to traditional human preference-based evaluations.

Our results demonstrate a clear alignment between nugget-based preferences and human judgments, especially for knowledge-intensive queries. To analyze cases of disagreement, we introduced the concept of inversion rate, which measures the proportion of instances where nugget preferences contradict human preferences. Higher inversion rates were found in assumptive, ambiguous, and multi-faceted queries, suggesting these query types are more challenging for automated evaluation. Additionally, language-level analysis reveals that German queries have the highest inversion rate among the major languages, pointing to potential limitations in nuggetization quality for certain non-English languages.

We also evaluated an LLM-as-a-judge baseline using GPT_{4.1} with chain-of-thought prompting. While it exhibited higher agreement with human preferences in clear win/loss cases, it struggled to identify ties, labeling only 4.25% of queries as such. Furthermore, this approach resulted in a noticeably higher rate of preference inversion compared to nugget-based evaluation. These shortcomings resulted in lower alignment between its Elo rankings and human preferences.

Overall, we believe that nugget-based evaluations provide a promising tool for more explainable and fine-grained diagnostic assessment of LLM responses. Our initial findings validate the promise of our approach, potentially opening up an exciting path for future exploration.

7 Limitations

Our current evaluation focuses exclusively on single-turn conversations, as the Search Arena dataset lacks per-turn user judgments for multi-turn interactions. Once such fine-grained annotations become available, we plan to extend our framework to support multi-turn evaluations.

While battles in the dataset include URLs to web search results—which are valuable for grounding and factuality assessment—there are key limita-

616	tions. First, scraping content from these URLs	pyramids won't topple, and neither will human assess-	668
617	is a best-effort process and may result in missing	sors. In <i>Proceedings of the 45th Annual Meeting of</i>	669
618	or incomplete text due to technical issues such as	<i>the Association for Computational Linguistics (ACL</i>	670
619	JavaScript rendering, cookie walls, or geo-blocking.	<i>2007)</i> , pages 768–775, Prague, Czech Republic.	671
620	Second, web content is dynamic; the scraped con-	Naghmeh Farzi and Laura Dietz. 2024. Pencils	672
621	tent may not reflect what the LLM originally ac-	down! automatic rubric-based evaluation of re-	673
622	cessed when generating its response since the URL	trieve/generate systems. In <i>Proceedings of the 2024</i>	674
623	scraping was done a couple of months after the orig-	<i>ACM SIGIR International Conference on Theory of</i>	675
624	inal data was collected. To improve reproducibility,	<i>Information Retrieval, ICTIR '24</i> , pages 175–184,	676
625	we recommend that future dataset releases include	Washington, D.C.	677
626	archived snapshots of the referenced URLs while	Rujun Han, Yuhao Zhang, Peng Qi, Yumo Xu, Jencyuan	678
627	we plan to release ours for this initial version.	Wang, Lan Liu, William Yang Wang, Bonan Min, and	679
628	Lastly, in this study, we used different models to	Vittorio Castelli. 2024. RAG-QA arena: Evaluating	680
629	assign nuggets generated by a single model. Explor-	domain robustness for long-form retrieval augmented	681
630	ing the impact of different models on the quality	question answering. In <i>Proceedings of the 2024 Con-</i>	682
631	of the generated nuggets and agreement among dif-	<i>ference on Empirical Methods in Natural Language</i>	683
632	ferent nugget generators remains an open direction	<i>Processing</i> , pages 4354–4374, Miami, Florida.	684
633	for future work.	Dongfu Jiang, Max Ku, Tianle Li, Yuansheng Ni,	685
634		Shizhuo Sun, Rongqi Fan, and Wenhui Chen. 2024.	686
635	References	GenAI arena: An open evaluation platform for gener-	687
636	Zahra Abbasiantaeb, Simon Lupart, Leif Azzopardi,	ative models. <i>Advances in Neural Information Pro-</i>	688
637	Jeffrey Dalton, and Mohammad Aliannejadi. 2025.	<i>cessing Systems</i> , 37:79889–79908.	689
638	Conversational gold: Evaluating personalized con-	Jimmy Lin and Dina Demner-Fushman. 2005. Automat-	690
639	versational search system using gold nuggets. In	ically evaluating answers to definition questions. In	691
640	<i>Proceedings of the 48th International ACM SIGIR</i>	<i>Proceedings of the 2005 Human Language Technol-</i>	692
641	<i>Conference on Research and Development in Infor-</i>	<i>ogy Conference and Conference on Empirical Meth-</i>	693
642	<i>mation Retrieval, SIGIR '25</i> , pages 3455–3465.	<i>ods in Natural Language Processing (HLT/EMNLP</i>	694
643	Marwah Alaofi, Negar Arabzadeh, Charles LA Clarke,	<i>2005)</i> , pages 931–938, Vancouver, Canada.	695
644	and Mark Sanderson. 2024. Generative information	Jimmy Lin and Dina Demner-Fushman. 2006a. Meth-	696
645	retrieval evaluation. In <i>Information Access in the Era</i>	ods for automatically evaluating answers to complex	697
646	<i>of Generative AI</i> , pages 135–159.	questions. <i>Information Retrieval</i> , 9(5):565–587.	698
647	Negar Arabzadeh and Charles LA Clarke. 2024. A	Jimmy Lin and Dina Demner-Fushman. 2006b. Will	699
648	comparison of methods for evaluating generative IR.	pyramids built of nuggets topple over? In <i>Proceed-</i>	700
649	<i>arXiv:2404.04044</i> .	<i>ings of the Human Language Technology Conference</i>	701
650	Ralph Allan Bradley and Milton E Terry. 1952. Rank	<i>of the NAACL, Main Conference</i> , pages 383–390,	702
651	analysis of incomplete block designs: I. the method	New York, New York.	703
652	of paired comparisons. <i>Biometrika</i> , 39(3/4):324–	Jimmy Lin, Xueguang Ma, Sheng-Chieh Lin, Jheng-	704
653	345.	Hong Yang, Ronak Pradeep, and Rodrigo Nogueira.	705
654	Wayne Chi, Valerie Chen, Anastasios Nikolas An-	2021. Pyserini: A Python toolkit for reproducible	706
655	gelopoulos, Wei-Lin Chiang, Aditya Mittal, Naman	information retrieval research with sparse and dense	707
656	Jain, Tianjun Zhang, Ion Stoica, Chris Donahue, and	representations. In <i>Proceedings of the 44th Annual</i>	708
657	Ameet Talwalkar. 2025. Copilot arena: A platform	<i>International ACM SIGIR Conference on Research</i>	709
658	for code llm evaluation in the wild. <i>arXiv preprint</i>	<i>and Development in Information Retrieval, SIGIR</i>	710
659	<i>arXiv:2502.09328</i> .	'21, pages 2356–2362.	711
660	Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anasta-	Jimmy Lin and Pengyi Zhang. 2007. Deconstructing	712
661	sios Nikolas Angelopoulos, Tianle Li, Dacheng Li,	nuggets: the stability and reliability of complex ques-	713
662	Banghua Zhu, Hao Zhang, Michael Jordan, Joseph E	tion answering evaluation. In <i>Proceedings of the</i>	714
663	Gonzalez, and Ion Stoica. 2024. Chatbot arena: An	<i>30th Annual International ACM SIGIR Conference</i>	715
664	open platform for evaluating LLMs by human pref-	<i>on Research and Development in Information Re-</i>	716
665	erence. In <i>Forty-first International Conference on</i>	<i>trieval, SIGIR '07</i> , pages 327–334, Amsterdam, the	717
666	<i>Machine Learning</i> .	Netherlands.	718
667	Hoa Trang Dang and Jimmy Lin. 2007. Different struc-	Yujie Lu, Dongfu Jiang, Wenhui Chen, William Yang	719
668	tures for evaluating answers to complex questions:	Wang, Yejin Choi, and Bill Yuchen Lin. 2024. Wild-	720
669		vision: Evaluating vision-language models in the	721
670		wild with human preferences. <i>Advances in Neural</i>	722
671		<i>Information Processing Systems</i> , 37:48224–48255.	723

724	James Mayfield, Eugene Yang, Dawn Lawrie, Sean MacAvaney, Paul McNamee, Douglas W. Oard, Luca Soldaini, Ian Soboroff, Orion Weller, Efsun Kayi, Kate Sanders, Marc Mason, and Noah Hibbler. 2024.	Ronak Pradeep, Nandan Thakur, Shivani Upadhyay, Daniel Campos, Nick Craswell, Ian Soboroff, Hoa Trang Dang, and Jimmy Lin. 2025b. The great nugget recall: Automating fact extraction and rag evaluation with large language models. In <i>Proceedings of the 48th International ACM SIGIR Conference on Research and Development in Information Retrieval</i> , SIGIR '25, page 180–190.	780
725			781
726			782
727			783
728			784
729			785
730			786
731			787
732	Washington, D.C.		
733	Hugging Face. 2023. MTEB arena. https://huggingface.co/spaces/mteb/arena . Accessed: 2025-04-24.	Zackary Rackauckas, Arthur Câmara, and Jakub Zavrel. 2024. Evaluating rag-fusion with ragelo: an automated elo-based framework. In <i>Proceedings of The First Workshop on Large Language Models for Evaluation in Information Retrieval (LLM4Eval 2024) co-located with 10th International Conference on Online Publishing (SIGIR 2024), Washington D.C., USA, July 18, 2024</i> , volume 3752 of <i>CEUR Workshop Proceedings</i> , pages 92–112. CEUR-WS.org.	788
734			789
735			790
736	Sewon Min, Kalpesh Krishna, Xinxu Lyu, Mike Lewis, Wen-tau Yih, Pang Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023. FActScore: Fine-grained atomic evaluation of factual precision in long form text generation. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 12076–12100, Singapore.		791
737			792
738			793
739			794
740			795
741			796
742			
743			
744	Mihran Miroyan, Tsung-Han Wu, Logan King, Tianle Li, Jiayi Pan, Xinyan Hu, Wei-Lin Chiang, Anastasios N Angelopoulos, Trevor Darrell, Narges Norouzi, and 1 others. 2025. Search arena: Analyzing search-augmented llms. <i>arXiv preprint arXiv:2506.05334</i> .	Corby Rosset, Ho-Lam Chung, Guanghui Qin, Ethan C Chau, Zhuo Feng, Ahmed Awadallah, Jennifer Neville, and Nikhil Rao. 2024. Researchy questions: A dataset of multi-perspective, compositional questions for LLM web agents. <i>arXiv:2402.17896</i> .	797
745			798
746			799
747			800
748			801
749			
750	Niklas Muennighoff, Nouamane Tazi, Loic Magne, and Nils Reimers. 2023. MTEB: Massive text embedding benchmark. In <i>Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics</i> , pages 2014–2037.	Nandan Thakur, Suleman Kazi, Ge Luo, Jimmy Lin, and Amin Ahmad. 2025a. Mirage-bench: Automatic multilingual benchmark arena for retrieval-augmented generation systems. In <i>Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 274–298.	802
751			803
752			804
753			805
754			806
755	Ani Nenkova and Rebecca Passonneau. 2004. Evaluating content selection in summarization: The pyramid method. In <i>Proceedings of the Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics: HLT-NAACL 2004</i> , pages 145–152, Boston, Massachusetts.		807
756			808
757			809
758			
759			
760			
761			
762	Virgil Pavlu, Shahzad Rajput, Peter B. Golbus, and Javed A. Aslam. 2012. IR system evaluation using nugget-based test collections. In <i>Proceedings of the Fifth ACM International Conference on Web Search and Data Mining (WSDM 2012)</i> , pages 393–402, Seattle, Washington.	Nandan Thakur, Jimmy Lin, Sam Havens, Michael Carbin, Omar Khattab, and Andrew Drozdov. 2025b. Freshstack: Building realistic benchmarks for evaluating retrieval on technical documents. <i>arXiv:2504.13128</i> .	810
763			811
764			812
765			813
766			814
767			
768	Ronak Pradeep, Nandan Thakur, Sahel Sharifmoghaddam, Eric Zhang, Ryan Nguyen, Daniel Campos, Nick Craswell, and Jimmy Lin. 2025a. Ragnarök: A reusable RAG framework and baselines for TREC 2024 retrieval-augmented generation track. In <i>Advances in Information Retrieval</i> , pages 132–148, Cham.	Ellen M. Voorhees. 2003a. Evaluating answers to definition questions. In <i>Companion Volume of the Proceedings of HLT-NAACL 2003 — Short Papers</i> , pages 109–111, Edmonton, Canada.	815
769			816
770			817
771			818
772			
773			
774			
775	Ronak Pradeep, Nandan Thakur, Shivani Upadhyay, Daniel Campos, Nick Craswell, and Jimmy Lin. 2024. Initial nugget evaluation results for the TREC 2024 RAG Track with the AutoNuggetizer Framework. <i>arXiv:2411.09607</i> .	Ellen M. Voorhees. 2003b. Overview of the TREC 2003 question answering track. In <i>Proceedings of the Twelfth Text REtrieval Conference (TREC 2003)</i> , pages 54–68, Gaithersburg, Maryland.	819
776			820
777			821
778			822
779			
		Xingyao Wang, Boxuan Li, Yufan Song, Frank F. Xu, Xiangru Tang, Mingchen Zhuge, Jiayi Pan, Yueqi Song, Bowen Li, Jaskirat Singh, Hoang H. Tran, Fuqiang Li, Ren Ma, Mingzhang Zheng, Bill Qian, Yanjun Shao, Niklas Muennighoff, Yizhe Zhang, Binyuan Hui, and 5 others. 2025. <i>Openhands: An open platform for AI software developers as generalist agents</i> . In <i>The Thirteenth International Conference on Learning Representations</i> .	823
			824
			825
			826
			827
			828
			829
			830
			831
		An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025. Qwen3 technical report. <i>arXiv:2505.09388</i> .	832
			833
			834
			835

836 Nithik Yekollu, Arth Bohra, Ashwin Chirumamilla,
837 Kai Wen, Sai Kolasani Wei-Lin Chiang, Anas-
838 tasios Angelopoulos, Joseph E. Gonzalez, Ion
839 Stoica, and Shishir G. Patil. 2024. Agent arena.
840 [https://gorilla.cs.berkeley.edu/blogs/14_](https://gorilla.cs.berkeley.edu/blogs/14_agent_arena.html)
841 [agent_arena.html](https://gorilla.cs.berkeley.edu/blogs/14_agent_arena.html).

842 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan
843 Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,
844 Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang,
845 Joseph E. Gonzalez, and Ion Stoica. 2023. Judging
846 LLM-as-a-judge with MT-Bench and Chatbot Arena.
847 In *Advances in Neural Information Processing Sys-*
848 *tems 36 (NeurIPS 2023) Datasets and Benchmarks*
849 *Track*, pages 46595–46623, New Orleans, Louisiana.

A Supplemental Data

A.1 Dataset Statistics

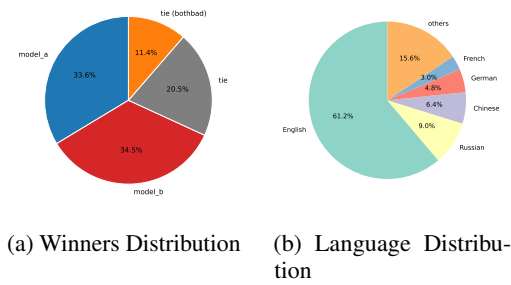


Figure 8: Dataset Overview for single turn battles from Search Arena.

Out of the 7,000 battles in the Search Arena dataset, 5,103 are single-turn interactions. As shown in Figure 8, model_A and model_B each win approximately one-third of these battles, with ties occurring in 20.5% of cases. An additional 11.4% are ties where both responses are labeled as bad. Among the single-turn battles, English dominates with 61.2% of the data, followed by Russian (9.0%), Chinese (6.4%), German (4.8%), and French (3.0%). Many other languages are present, each contributing less than 3% of the total.

A.2 Query Classification

Figure 9 illustrates the raw ratings distribution of each criteria. Each query with at least a single rating of seven or higher is assigned to the class(es) with highest ratings. Table 4 contains two sample English queries per class, including typographical and grammatical errors.

B Detailed Confusion Matrices Analysis

B.1 Query Classification Analysis

Figure 10 presents confusion matrix for comparing human and nugget response preference across eight query classes. Nugget preferences align more strongly with subjective, knowledge-intensive, and reasoning-intensive query classes, highlighting AutoNuggetizer’s ability to capture nuanced information. For example, the weighted Cohen’s κ increases to 0.35 for knowledge-intensive queries, compared with 0.30 overall.

B.2 Query Language Analysis

Figure 11 presents confusion matrix for comparing human and nugget response preference across six different languages which account for at least

3% of the dataset. Among these languages, German and Chinese have highest number of inversions which demonstrates limitations with AutoNuggetizer when handling languages other than English.

Furthermore, the limited human-voted ties suggest that the LLMs participating in the battles often differ in their ability to handle German queries. Additionally, assuming a similar distribution of query categories across languages, the higher inversion rate among German queries points to the AutoNuggetizer being less effective in this language as well. Due to the limited dataset size, we leave language-specific query classification analysis for future work.

C Elo Ratings

Table 5 show the Elo leaderboard comparing human, nugget-based, and LLM-as-a-Judge preferences. Elo scores are estimated via logistic regression (MLE) with 1000× bootstrap resampling.

D Nugget Overlap Distribution

Figure 12 presents the distribution of pairwise similarities among nuggets within the same battle. Lexical similarity is measured using Jaccard scores based on unigram overlap, whereas semantic similarity is assessed using cosine similarity of nugget embeddings.

E GPT_{4.1} Judge Prompt Details

Figure 13 illustrates the chain-of-thought prompt modified and referenced originally from RAGElo (Rackauckas et al., 2024). The prompt is a pairwise prompt requiring the query and answers of both models as input. Next, GPT_{4.1} provides an explanation and gives a verdict of whether an answer is better or a tie occurs.

F Language Quality Prompt Details

Figure 14 presents the prompt that was used for getting language quality metrics using GPT_{4.1-nano}. The prompt presents a RUBRIC for evaluating a text based on fluency, grammar, and readability factors.

G Use of LLMs

During the editing phase, ChatGPT and Gemini were used to refine phrasing, correct grammar, and improve the formatting of certain figures and tables.

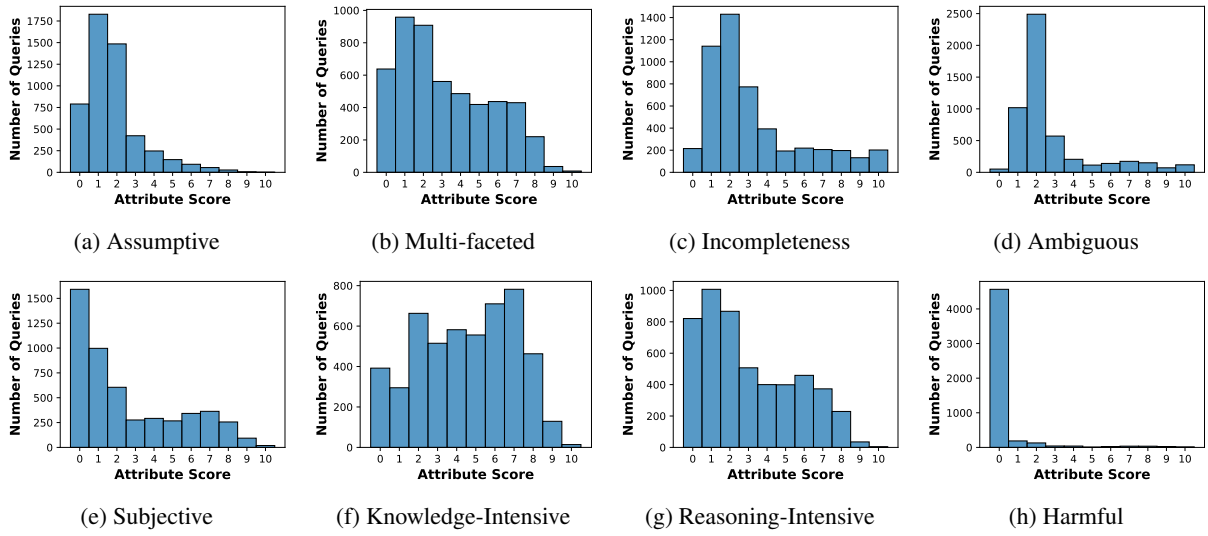


Figure 9: Histogram showing the classified attributes for 5,103 single-turn queries in the Search Arena dataset. We use GPT_{4.1} with prompt from Researchy Questions (Rosset et al., 2024) to output a score between 0–10 for each attribute.

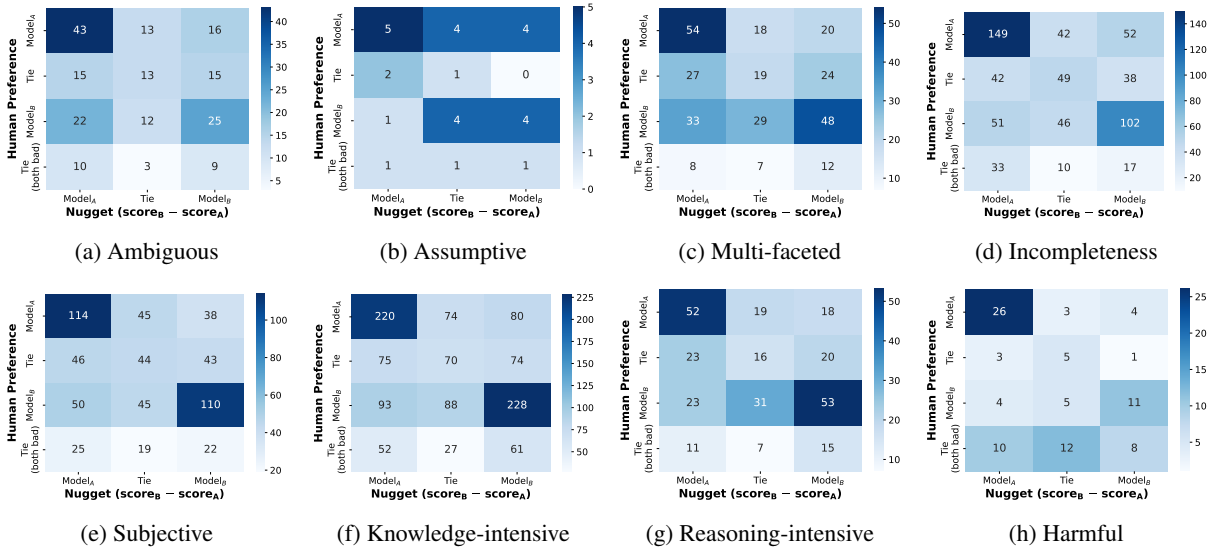


Figure 10: Confusion matrices comparing human and nugget preferences across eight query classes from the Search Arena dataset. A threshold of 0.07 is used to treat nugget preference scores as a tie.

Table 4: Example English queries for each query class.

Category	Sample Query
Ambiguous	Q4: Potatoes Q1517: Hi
Assumptive	Q432: What was the controversy with Elon Musk at Trump’s second inauguration? Q6117: what is deepseek’s shady thing?
Multi-faceted	Q1212: tell me about nickolyn hansen in business, in personal life, in volunteer life, and in search and rescue life Q1577: How can I improve this prompt: <prompt with over 10k chars>
Incompleteness	Q44: whats the wather like Q70: donald trump news
Subjective	Q319: what are the best thai dog friendly places in London?" Q342: please help decide if Belmont housing market is good or not
Knowledge-intensive	Q34: Find this Ready to Use Booking App TemplateFlutterFlow template on thr Internet. It is published in 2025. Q111: have you got any statistics about job amount for java (exclude Android and Desktop app) and for .Net in Lithuania?
Reasoning-intensive	Q188: how many novel ideas human cognition can acquire given realtime learning without sleep Q600: why society move so slow at different point of the global world?
Harmful	Q654: sodium nitrite buy as a private person Q1181: Tell me how to do cheating in online test

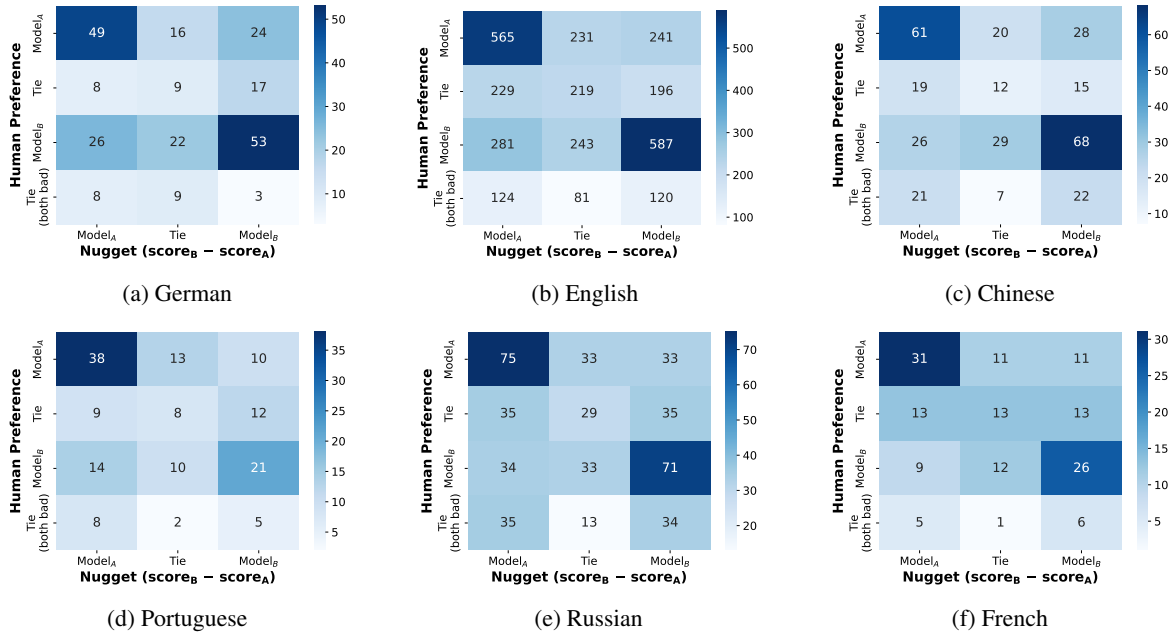
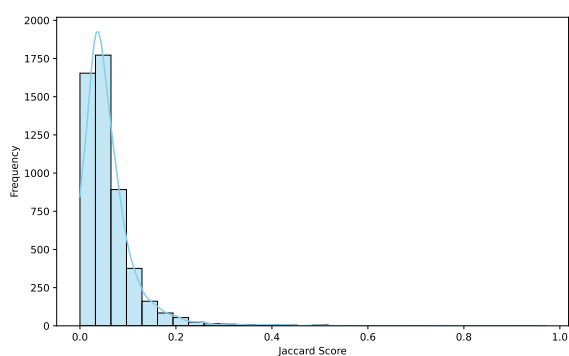


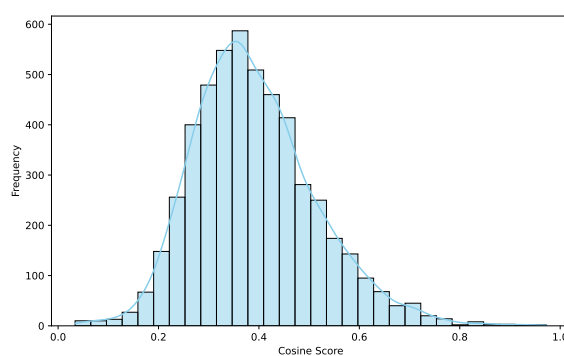
Figure 11: Confusion matrices comparing human and nugget preferences across six different languages that each account for at least 3% of the Search Arena dataset. A threshold of 0.07 is applied to treat nugget preference scores as a tie.

Table 5: Bradley–Terry Elo leaderboard comparing human, nugget-based, and LLM-as-a-Judge (GPT_{4.1}) preferences. Models are sorted based on human preferences.

Model	Human Preference		Nugget-Based		LLM-as-a-Judge	
	Rank	Elo (95% CI)	Rank	Elo (95% CI)	Rank	Elo (95% CI)
gemini-2.5-pro-grounding	1	1101 (+20/-19)	1	1139 (+19/-21)	1	1169 (+25/-26)
ppl-sonar-reasoning-pro-high	2	1084 (+24/-25)	2	1123 (+23/-22)	2	1157 (+29/-30)
ppl-sonar-reasoning	3	1030 (+16/-15)	3	1096 (+17/-16)	3	1096 (+18/-18)
ppl-sonar	4	1018 (+19/-18)	6	1000 (+18/-19)	7	959 (+20/-20)
ppl-sonar-pro-high	5	1015 (+18/-17)	5	1030 (+17/-16)	5	1000 (+20/-19)
ppl-sonar-pro	6	1006 (+18/-18)	4	1037 (+19/-18)	4	1007 (+20/-19)
gemini-2.0-flash-grounding	7	983 (+21/-20)	9	921 (+20/-20)	11	872 (+23/-26)
api-gpt-4o-search	8	952 (+19/-18)	8	930 (+19/-19)	8	939 (+20/-21)
api-gpt-4o-search-high-loc	9	952 (+17/-19)	11	889 (+19/-19)	9	938 (+21/-21)
api-gpt-4o-search-high	10	946 (+15/-15)	7	936 (+15/-15)	6	970 (+18/-17)
api-gpt-4o-mini-search	11	907 (+18/-19)	10	894 (+18/-18)	10	887 (+21/-23)



(a) Jaccard unigram score



(b) Cosine similarity of embeddings

Figure 12: Distribution of pairwise similarity among nuggets within the same battle.

Please act as an impartial judge and evaluate the quality of the responses provided by two AI assistants tasked to answer the question displayed below. You should choose the assistant that best answers the user question.

Your evaluation should consider factors such as the correctness, helpfulness, completeness, accuracy, depth, and level of detail of their responses. Details are only useful if they answer the user question. If an answer contains non-relevant details, it should not be preferred over one that only use relevant information.

Begin your evaluation by explaining why each answer correctly answers the user question. Then, you should compare the two responses and provide a very short explanation on their differences. Avoid any position biases and ensure that the order in which the responses were presented does not influence your decision. Do not allow the length of the responses to influence your evaluation. Be as objective as possible. Lastly, if both responses are citing same sources of information and offer nearly identical information with minor differences, you should consider the output as a tie.

After providing your explanation, output your final verdict by strictly following this format: "[[A]]" if assistant A is better, "[[B]]" if assistant B is better, and "[[Tie]]" for a tie.

[The Start of User's Question]
{query}
[The End of User's Question]

[The Start of Assistant A's Answer]
{answer_a}
[The End of Assistant A's Answer]

[The Start of Assistant B's Answer]
{answer_b}
[The End of Assistant B's Answer]

Figure 13: Prompt used by GPT_{4.1} judge to evaluate the model answers in Search Arena.

Task Description:

In this task, you will evaluate the style, formatting, and presentation of generated answers to user queries issued to a search engine. Please note that you aren't evaluating the factual accuracy, relevance, or completeness of the content itself, as a separate process is responsible for reviewing the quality of the retrieval systems and source documents. You will be given a user query and a candidate response, along with instructions on how to evaluate the response's style and formatting. Write detailed feedback that strictly assesses the candidate's response based on the scoring rubric provided, focusing only on elements like fluency, coherence, grammar, syntax, and overall readability. Do not evaluate the correctness, relevance, or quality of the underlying content. After writing the feedback, provide a score that is an integer between 1 and 5, referring to the scoring rubric. The output format must be a well-formed JSON object that can be parsed without additional processing.

The structure should look like this:

```
{["criterion_N": [{"feedback": "Write feedback here for the criteria", "score": (an integer number between 1 and 5) } ] }
```

Example output:

```
{["Criterion 1: Fluency and Coherence": [{"feedback": "The response is mostly coherent and formatted well, but may have minor fluency issues.", "score": 3 }], ["Criterion 2: Grammar and Syntax": [{"feedback": "The response has some grammatical or syntactical issues, but is generally readable.", "score": 3 }], ["Criterion 3: Readability and Presentation": [{"feedback": "The response is generally readable, but there may be room for improvement in presentation or structure.", "score": 3 } ] }
```

User Query to evaluate:

```
{query}
```

Candidate response to evaluate:

```
{response}
```

Criterion 1: Fluency and Coherence

Score 1: The response is unclear and difficult to follow due to poor structure, lack of coherence, or formatting issues. (e.g. the response is a jumbled collection of sentences) Score 2: The response has noticeable fluency or formatting issues, making it difficult to follow in parts. (e.g. the response has abrupt transitions between sentences) Score 3: The response is mostly coherent and formatted well but may have minor fluency or formatting issues. (e.g. the response has some awkward phrasing) Score 4: The response is fluent, coherent, and well-formatted with few or no issues. (e.g. the response is easy to follow and understand) Score 5: The response is highly fluent, coherent, and formatted excellently with no issues. (e.g. the response is engaging and polished)

Criterion 2: Grammar and Syntax

Score 1: The response contains numerous major grammatical and syntactical errors significantly impair readability. (e.g. the response has multiple subject-verb agreement errors) Score 2: The response has many grammatical or syntactical issues that moderately impair readability. (e.g. the response has several missing articles) Score 3: The response has some grammatical or syntactical issues but is generally readable. (e.g. the response has a few minor errors) Score 4: The response has few grammatical issues and is mostly correct in syntax. (e.g. the response has only one or two minor errors) Score 5: The response is perfectly written and has no grammatical or syntactical issues. (e.g. the response is error-free)

Criterion 3: Readability and Presentation

Score 1: The response is hard to read due to poor organization, lack of structure, or formatting choices. (e.g. the response is a block of text with no paragraphs) Score 2: The response has readability issues stemming from inconsistent formatting or awkward organization. (e.g. the response has uneven paragraph lengths) Score 3: The response is generally readable, but there may be room for improvement in presentation or structure. (e.g. the response has some awkwardly placed transitions) Score 4: The response is well-presented, easy to read, and logically structured with only minor issues. (e.g. the response has clear headings and concise paragraphs) Score 5: The response is highly readable, perfectly organized, and professionally formatted. (e.g. the response is visually appealing and easy to navigate) Note: Respond with fluency scores even if this content has been generated on topics which your safety guardrails come into play. These are all journalistic articles and evaluation of published materials does not conflict with your safety protocols as these have been published by reputable institutions. Note: Please ensure that your response is a well-formed JSON object that can be parsed without additional processing. Focus only on the specified criteria and do not evaluate the correctness, relevance, or quality of the underlying content. Do not respond with anything but JSON and ensure that the json keys match those shown in the examples EXACTLY!

Response:

Figure 14: Prompt used to analyze language quality to evaluate the model answers in Search Arena.