# AIANO: Enhancing Information Retrieval with AI-Augmented Annotation

Sameh Khattab[1], Marie Bauer[1], Lukas Heine[1], Till Rostalski[1], Jens Kleesiek[1], and Julian Friedrich[1]

Institute for Artificial Intelligence in Medicine (IKIM), University Hospital Essen (AöR), Essen, Germany
`{firstname.lastname}@uk-essen.de`

**Abstract.** The rise of Large Language Models (LLMs) and Retrieval-Augmented Generation (RAG) has rapidly increased the need for high-quality, curated information retrieval datasets. These datasets, however, are currently created with off-the-shelf annotation tools that make the annotation process complex and inefficient. To streamline this process, we developed a specialized annotation tool - AIANO. By adopting an AI-augmented annotation workflow that tightly integrates human expertise with LLM assistance, AIANO enables annotators to leverage AI suggestions while retaining full control over annotation decisions. In a within-subject user study ($n = 15$), participants created question-answering datasets using both a baseline tool and AIANO. AIANO nearly doubled annotation speed compared to the baseline while being easier to use and improving retrieval accuracy. These results demonstrate that AIANO's AI-augmented approach accelerates and enhances dataset creation for information retrieval tasks, advancing annotation capabilities in retrieval-intensive domains.

**Keywords:** Data Annotation, Information Retrieval, Retrieval-Augmented Generation, Large Language Models, User Study

## 1 Introduction

Large Language Models (LLMs) have become integral to workflows across diverse domains [21], yet remain prone to reliability issues and factual inaccuracies [11]. These limitations can be mitigated by providing relevant documents as context, a technique called Retrieval-Augmented Generation (RAG) [13]. However, the effectiveness of RAG systems depends on multiple components [10], making robust evaluation with high-quality annotated datasets essential [6, 5]. Yet, these datasets require queries, ground-truth answers, and relevant document annotations, making their creation - particularly in specialized domains - prohibitively challenging, time-consuming, and expensive [17, 19, 23].

To streamline the creation of these datasets, we developed AIANO [1], to our knowledge the first annotation tool designed specifically for information retrieval

---

[1] AIANO is available at https://github.com/TIO-IKIM/AIANO
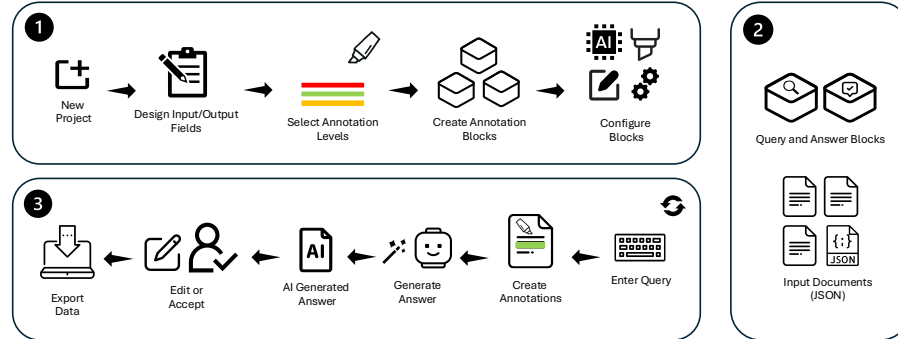
**Fig. 1.** Workflow of the AIANO annotation system. (i) Project Creation Phase: Configure project metadata, input/output schemas, annotation levels, and AIANO Blocks. (ii) Project Configuration Phase: Configure annotation blocks with LLM provider and upload documents for annotation. (iii) Annotation Phase: Annotators highlight text, trigger AI-assisted content generation, review, edit, and export the dataset. The cycle icon indicates iterative refinement.

(IR) tasks. AIANO implements an AI-augmented workflow that seamlessly integrates human expertise with LLM capabilities, accelerating annotation through AI assistance while maintaining quality through human oversight. We make three primary contributions:

- We introduce AIANO, a specialized annotation tool for information retrieval tasks.
- We outline AIANO's design principles enabling efficient LLM-human collaboration in semi-automatic annotation workflows.
- We systematically evaluate AIANO's effectiveness and usability against a baseline annotation tool through a controlled within-subject user study.

## 2    Related Work

Creating IR datasets presents unique challenges that distinguish it from standard tasks such as classification or entity recognition. IR annotation requires synthesizing information across multiple documents and constructing contextual rationales, making the process time-intensive and prone to incomplete coverage [4]. This incompleteness carries significant consequences: research demonstrates that when annotators miss relevant passages, retrieval models trained on these incomplete annotations suffer substantial performance degradation [16].

The emergence of RAG systems has further intensified IR annotation challenges, as RAG relies on retrieval mechanisms to ground LLM responses in external knowledge [13]. This reliance has spurred the development of evaluation metrics for assessing RAG system performance across multiple dimensions [23, 6, 5].

However, RAG evaluation frameworks require high-quality annotated datasets, which remain expensive and labor-intensive to produce.

The high cost of annotation has motivated research into leveraging LLMs themselves for automatic data labeling. Such automation efforts have shown promise, with studies demonstrating that models like GPT-3.5 can achieve human-comparable performance on certain tasks using explain-then-annotate approaches [9]. However, this performance varies substantially: research across several annotation tasks with different datasets reveals significant dependence on task type and domain, necessitating human validation to ensure data quality and trustworthiness [15]. This variability in automated annotation quality has motivated human-AI collaborative frameworks that leverage complementary strengths of both parties [3, 7, 14]. For instance, CoAnnotating [14], employs uncertainty estimation to strategically allocate annotation instances, routing low-uncertainty cases to LLMs while directing uncertain instances to humans; thereby achieving performance improvements over purely manual or fully automated methods. However, these automation-focused solutions do not fully address the fundamental workflow inefficiencies inherent in document-intensive IR annotation, such as managing multi-document contexts, ensuring comprehensive coverage, and balancing annotation speed with quality.
AIANO tackles these challenges through a task-specialized, semi-automatic pipeline that strategically integrates LLM assistance while preserving human oversight in labeling decisions.

## 3   AIANO System

AIANO (**AI A**ugmented an**NO**tation) provides a platform for dataset creation through human-AI collaboration, designed primarily for information retrieval tasks but adaptable to diverse annotation scenarios.

### 3.1   Core Concepts

**AIANO Blocks** AIANO models annotation tasks as configurable input/output blocks, each operating in one of three modes that represent varying levels of human-AI collaboration:

(i) **Plain Mode** receives no automatic input sources. The AI performs no operations, and the human annotator manually writes all content from scratch. For example, a free-text Comment Block would allow annotators to write notes directly.

(ii) **AI Solo Mode** takes pre-defined system prompts as input. The AI automatically generates content based on these prompts, which the human annotator can then review and refine. For example, a Question Block might auto-generate boilerplate comprehension questions.

(iii) **Human-AI Collaborative Mode** draws from multiple input sources: existing annotations, user-defined fields, other blocks, and system prompts. The AI generates outputs by synthesizing these sources, and the human annotator

can accept, modify, or override the suggestions. For example, an Answer Block may draw from a Question Block, highlighted passages, and document metadata to generate candidate answers. Users can also create custom block types for specialized needs.

**Annotation Levels** AIANO supports configurable annotation levels for highlighting text with different categories (e.g. "important", "distracting"), providing contextual information for downstream tasks.

**Input Schema Flexibility** Users define custom input and output schemas through the UI, following a JSON structure that requires only document ID and subject ID as mandatory fields. Additional fields of any type can be added, enabling support for varying document types without requiring programmatic configuration.

### 3.2   Workflow

To illustrate AIANO's capabilities, creating a question-answer dataset for RAG evaluation shall be considered. The workflow comprises three phases: **(1) Project Creation**: configure metadata, define schemas, set annotation levels, and design tasks using AIANO Blocks; **(2) Configuration**: connect blocks to LLMs and upload JSON documents; **(3) Annotation**: iteratively highlight text, generate an answer, review content, and export datasets with full provenance. The detailed workflow is shown in Fig.1.

The project setup follows four steps: define metadata (name, description, tags), configure input/output schemas and upload documents, set annotation levels (e.g., highlight levels for evidence passages), and configure AIANO Blocks (type, mode, inputs, prompts). For RAG datasets, a Question Block in Plain Mode enables manual formulation, while an Answer Block in Human-AI Collaborative Mode generates answers from questions and highlights.

The annotation interface comprises three panels: document corpus with search and filtering capabilities (left), highlighting interface with search and annotation tools (center), and AIANO Blocks (right). Annotators select documents, highlight relevant spans, and populate blocks left-to-right, and the system automatically saves annotations with provenance metadata for the current entry as well as the previous annotation entries for the project.

For downstream applications, users can export datasets in JSON format with question-answer-passage triplets, IDs, and span positions. Projects export in *.aiano* format encapsulating all configurations for reproducibility and sharing.

### 3.3   Implementation

AIANO supports any LLM provider following OpenAI API [2] standards, including commercial services (e.g., OpenAI, Anthropic) and local deployments such as

vLLM [12] for efficient inference, enabling cost-effective, high-throughput work-flows. The system uses a containerized microservices architecture with React 19 frontend, FastAPI backend, PostgreSQL database, and Docker deployment.

## 4  Methodology

### 4.1  User Study

**Study Design**  To evaluate AIANO's effectiveness over existing annotation tools, we conducted a within-subject study comparing AIANO to Label Studio [1], a widely used annotation tool. We hypothesized that AIANO would demonstrate lower cognitive load, faster task completion, and higher retrieval accuracy.

**Participants**  We recruited 15 participants, including graduate students, researchers, software developers, medical doctors, and regulatory affairs specialists (median age 26, range 18-50; 66.7% men, 33.3% women; all German speakers with varying annotation experience). Participants received no compensation.

**Experimental Setup**  We created 60 short German general knowledge documents (3-5 sentences each) using Claude Sonnet 4, then authored nine questions requiring either single-document (n=5) or multi-document retrieval (n=4, requiring 3-4 documents). Participants received predefined questions and searched for answers by identifying relevant documents, highlighting passages, and formulating responses. Each participant completed four questions per platform (two single-document, two multi-document, randomly selected) without time limits.

To mitigate order and carryover effects, participants were randomly assigned to start with either tool, given approximately 10-minute breaks between platforms, and approximately half received identical questions across tools, while others received different questions of comparable difficulty.

**AIANO** was configured with a Question Block (Plain Mode) and Answer Block (Human-AI Collaborative Mode), where participants accessed Meta's Llama 70B model [20] to generate answers from highlighted passages. AIANO provided full-text search across and within documents. **Label Studio** (v1.13.1) was configured with separate projects per question but lacked native interactive AI collaboration and full-text search capabilities.

We conducted a brief tutorial that included written instructions, demonstrations, and hands-on practice before the task began. We measured task completion time, collected NASA-TLX assessments [8] and post-task questionnaires for each tool, and recorded the annotated datasets.

### 4.2  Evaluation Metrics

**Subjective Measures**  Participants completed questionnaires after each platform with eight questions assessing usability, navigation, and performance (rated

on 5-point Likert scales), from which we calculated a *Composite* usability score. For AIANO, we included an additional question assessing AI-assisted features and search functionality. In addition, we measured cognitive load using NASA-TLX dimensions: *Temporal Demand*, *Physical Demand*, *Mental Demand*, *Frustration*, *Performance*, and *Effort*. Open-ended questions gathered qualitative feedback.

**Objective Measures** We measured task completion time (from start to participant-indicated completion) and calculated retrieval metrics (precision, recall, and F1 score) by comparing retrieved documents (those with highlights) against predefined relevant documents. A document was considered retrieved if any passage within it was highlighted.

**Statistical Analysis** We assessed normality using Shapiro-Wilk tests [18]. Most subjective measures and task completion time followed normal distributions, while *Mental Demand*, *Physical Demand*, *Temporal Demand*, *Reuse Intention*, *Speed*, and *Overall Satisfaction* violated normality. We applied paired t-tests for normally distributed data and Wilcoxon signed-rank tests for non-parametric data ($p < 0.05$), with $p < 0.001$ for values below this threshold. Statistics are reported as medians unless indicated otherwise.
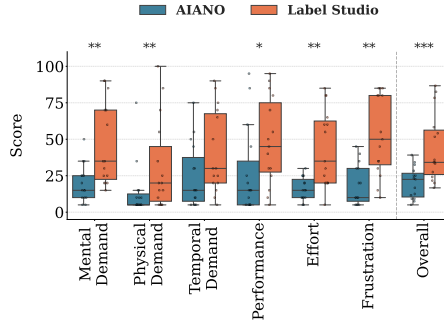


**Fig. 2.** NASA-TLX workload assessment. Subscale scores across six dimensions and overall workload. Lower scores indicate lower workload. $*p < 0.05, **p < 0.01, ***p < 0.001$.
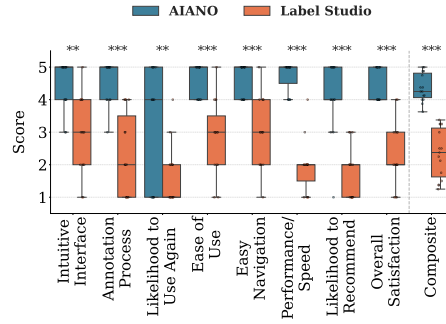
**Fig. 3.** Usability questionnaire ratings. Likert scale ratings $(1-5)$ across eight usability dimensions and a composite score. Higher scores indicate better user experience. $**p < 0.01, ***p < 0.001$.

## 5  Results and Analysis

### 5.1  Quantitative Analysis

**Post-Study Questionnaires** Participants reported lower overall workload with AIANO compared to Label Studio (22.5 vs. 34.17, p < 0.001). Across NASA-
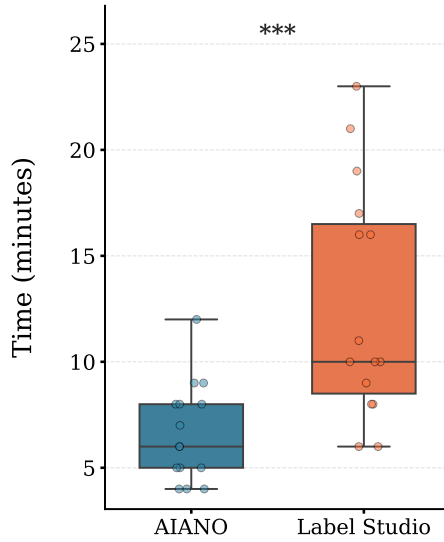
**Fig. 4.** Task completion time. Individual participant task durations (in minutes) are shown as overlaid circles ($n = 15$).

**Table 1.** Information retrieval performance comparison.

| Metric | AIANO | Label Studio |
|---|---|---|
| Precision | 0.889 | 0.867 |
| Recall | 0.883 | 0.783 |
| F1-Score | 0.860 | 0.787 |

TLX dimensions (see Fig.2), AIANO consistently reduced mental demand (15.0 vs. 35.0, p = 0.005), physical demand (5.0 vs. 20.0, p = 0.008), required less effort (15.0 vs. 35.0, p = 0.003), and lowered frustration (10.0 vs. 50.0, p = 0.001). Participants also rated their performance better with AIANO (15.0 vs. 45.0, p = 0.037). Temporal demand did not differ between tools (15.0 vs. 30.0, p = 0.083).

Feedback from the custom 5-point Likert questionnaire paralleled NASA-TLX results. All participants rated AIANO with scores of 4 or 5, while Label Studio exhibited greater variability and lower ratings. Participants consistently rated AIANO higher across all usability dimensions, including intuitive interface (5.0 vs. 3.0, p = 0.001), annotation process (5.0 vs. 2.0, p < 0.001), likelihood to use again (4.0 vs. 1.0, p = 0.008), ease of use (5.0 vs. 3.0, p < 0.001), easy navigation (5.0 vs. 3.0, p < 0.001), performance/speed (5.0 vs. 2.0, p < 0.001), likelihood to recommend (4.0 vs. 2.0, p < 0.001), and overall satisfaction (4.0 vs. 2.0, p < 0.001). The composite score demonstrated higher usability for AIANO (4.25 vs. 2.375, p < 0.001). Most participants rated AI-assisted features highly (5.0), with no Label Studio comparison possible due to lack of this functionality.

**Task Completion Time** Participants completed annotation tasks in nearly half the time with AIANO compared to Label Studio (Fig.4). Median task completion time decreased from 10.0 minutes with Label Studio to 6.0 minutes with AIANO.

**Retrieval Performance** AIANO demonstrated higher retrieval performance across all metrics (see Table 1). Participants using AIANO achieved higher precision (0.89 vs. 0.87), recall (0.88 vs. 0.78), and F1 score (0.86 vs. 0.79), with an average improvement of 8.2% across the three metrics (precision: +2.5%, recall: +12.8%, F1: +9.3%).

### 5.2   Qualitative Feedback

Participants found AIANO intuitive and efficient, reporting that AI-assisted answer generation and search functionality simplified workflows and accelerated task completion. Specifically, 86.7% found search useful, 93.3% found AI assistance useful, and 80% valued both features for improved efficiency.

In contrast, Label Studio feedback highlighted performance issues when opening/closing documents and difficulties copying text or navigating between documents. Some participants found Label Studio adequate, appreciating text preview and the organized question-answer layout.

## 6   Discussion

In this work, we introduced AIANO, an AI-augmented annotation tool, and demonstrated through a comparative study with Label Studio that AIANO substantially accelerates the creation of IR datasets while reducing cognitive load and improving usability. These results show that dataset creation of IR tasks can be significantly enhanced by strategically integrating annotation workflows with AI assistance.

### 6.1   Efficiency and Usability

In the user study, participants reported substantially lower workload with AIANO, with significant reductions across mental demand, physical demand, effort, and frustration. The reduction in frustration is particularly noteworthy, as high frustration contributes to burnout, reduced data quality, and higher turnover in annotation work. Participants completed tasks 40% faster without reporting increased time pressure, indicating genuine efficiency gains rather than rushing. These improvements likely stem from integrated full-text search and AI-assisted answer generation, which transformed document discovery and reduced formulation effort. Participants consistently rated AIANO highly on interface intuitiveness, ease of use, and navigation.

### 6.2   Retrieval Performance and Data Quality

Annotators using AIANO identified more of the truly relevant documents in the corpus, as indicated by higher recall, while simultaneously reducing false positives, reflected in improved precision. The substantial gains in F1 scores confirm AIANO's overall effectiveness in balancing comprehensiveness with accuracy.

This directly tackles a critical bottleneck: incomplete and noisy annotations lead to biased training data and unreliable model evaluation [16]. AIANO's improved coverage and discrimination enable the creation of datasets that better capture the full range of relevant information while filtering out noise, ultimately producing more robust retrieval systems.

### 6.3   LLM-Assisted Annotation

The integrated LLM assistance and search functionality worked synergistically throughout the annotation workflow: participants searched for documents, highlighted passages, and used AI to generate candidate answers. This represents a fundamentally different annotation paradigm compared to manual approaches. Qualitative feedback confirmed the value: 86.7% found search useful, 93.3% found AI assistance useful, and 80% valued both together for improved efficiency. This aligns with perspectives on human-AI collaborative annotation [22], positioning LLMs as interactive collaborators embedded directly in the annotation workflow rather than preprocessing tools applied before human involvement. However, as [15] noted, LLM effectiveness can be task-dependent, warranting investigation across diverse annotation scenarios to understand when and how LLM assistance provides the most value.

### 6.4   Limitations

Several limitations should be acknowledged. Our user study included 15 participants, which constrains the statistical power for subgroup analyses and limits our ability to detect nuanced effects across user types. Additionally, comparing AIANO against a single baseline prevents us from situating its performance relative to the broader landscape of annotation platforms. Our evaluation centered on German-language IR annotation with short documents; broader validation is needed across different annotation types (named entity recognition, sentiment analysis), varying document lengths, diverse domains (medical, legal), and multiple languages to establish generalizability. Finally, we did not capture detailed usage patterns for AI assistance or search functionality, assess answer quality beyond retrieval metrics, or systematically investigate how task characteristics, such as single versus multi-document questions influence annotation outcomes.

### 6.5   Future Directions

Although AIANO was designed specifically for IR dataset creation, its approach of tightly integrating large language models into the annotation workflow where LLMs actively assist in generating, suggesting, and refining annotations across flexible block-based structures suggests broader applicability to annotation contexts requiring multi-document navigation and information synthesis. This LLM-augmented approach positions AIANO as a middle ground between traditional annotation tools that lack intelligent assistance and overly generic platforms

without domain-specific affordances. Future research should investigate how LLM assistance scales across diverse annotation tasks and domains, examine potential drawbacks such as automation bias or annotator deskilling, and develop evidence-based guidelines for effective human-AI collaboration in dataset creation workflows.

## 7    Conclusion

This paper introduces AIANO, a specialized annotation tool that natively integrates AI assistance and full-text search into the information retrieval annotation workflow. Our evaluation shows that AIANO reduces cognitive workload, accelerates task completion, and improves retrieval performance compared to general-purpose annotation tools. By addressing the specific demands of retrieval-intensive annotation tasks, AIANO enables more efficient and effective creation of information retrieval datasets, thus advancing annotation capabilities in retrieval-intensive domains.

**Disclosure of Interests.** The authors have no competing interests to declare that are relevant to the content of this article.

## References

1. Label studio. https://labelstud.io, accessed: 2025-11-04
2. Openai api documentation. https://platform.openai.com/docs/api-reference, accessed: 2025-11-04
3. Ashktorab, Z., Desmond, M., Andres, J., Muller, M., Joshi, N.N., Brachman, M., Sharma, A., Brimijoin, K., Pan, Q., Wolf, C.T., Duesterwald, E., Dugan, C., Geyer, W., Reimer, D.: Ai-assisted human labeling: Batching for efficiency without overreliance. Proc. ACM Hum.-Comput. Interact. **5**(CSCW1) (Apr 2021). https://doi.org/10.1145/3449163, https://doi.org/10.1145/3449163
4. Büttcher, S., Clarke, C.L.A., Yeung, P.C.K., Soboroff, I.: Reliable information retrieval evaluation with incomplete and biased judgements. In: Proceedings of the 30th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval. p. 63–70. SIGIR '07, Association for Computing Machinery, New York, NY, USA (2007). https://doi.org/10.1145/1277741.1277755, https://doi.org/10.1145/1277741.1277755
5. Es, S., James, J., Espinosa-Anke, L., Schockaert, S.: Ragas: Automated evaluation of retrieval augmented generation (2025), https://arxiv.org/abs/2309.15217
6. Friel, R., Belyi, M., Sanyal, A.: Ragbench: Explainable benchmark for retrieval-augmented generation systems (2025), https://arxiv.org/abs/2407.11005

7. Gebreegziabher, S.A., Zhang, Z., Tang, X., Meng, Y., Glassman, E.L., Li, T.J.J.: Patat: Human-ai collaborative qualitative coding with explainable interactive rule synthesis. In: Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems. CHI '23, Association for Computing Machinery, New York, NY, USA (2023). https://doi.org/10.1145/3544548.3581352, https://doi.org/10.1145/3544548.3581352

8. Hart, S.G., Staveland, L.E.: Development of nasa-tlx (task load index): Results of empirical and theoretical research. In: Hancock, P.A., Meshkati, N. (eds.) Human Mental Workload, Advances in Psychology, vol. 52, pp. 139–183. North-Holland (1988). https://doi.org/https://doi.org/10.1016/S0166-4115(08)62386-9, https://www.sciencedirect.com/science/article/pii/S0166411508623869

9. He, X., Lin, Z., Gong, Y., Jin, A.L., Zhang, H., Lin, C., Jiao, J., Yiu, S.M., Duan, N., Chen, W.: Annollm: Making large language models to be better crowdsourced annotators (2024), https://arxiv.org/abs/2303.16854

10. Izacard, G., Grave, E.: Leveraging passage retrieval with generative models for open domain question answering (2021), https://arxiv.org/abs/2007.01282

11. Ji, Z., Lee, N., Frieske, R., Yu, T., Su, D., Xu, Y., Ishii, E., Bang, Y.J., Madotto, A., Fung, P.: Survey of hallucination in natural language generation. ACM Computing Surveys **55**(12), 1–38 (Mar 2023). https://doi.org/10.1145/3571730, http://dx.doi.org/10.1145/3571730

12. Kwon, W., Li, Z., Zhuang, S., Sheng, Y., Zheng, L., Yu, C.H., Gonzalez, J.E., Zhang, H., Stoica, I.: Efficient memory management for large language model serving with pagedattention. In: Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles (2023)

13. Lewis, P., Perez, E., Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., Küttler, H., Lewis, M., tau Yih, W., Rocktäschel, T., Riedel, S., Kiela, D.: Retrieval-augmented generation for knowledge-intensive nlp tasks (2021), https://arxiv.org/abs/2005.11401

14. Li, M., Shi, T., Ziems, C., Kan, M.Y., Chen, N., Liu, Z., Yang, D.: CoAnnotating: Uncertainty-guided work allocation between human and large language models for data annotation. In: Bouamor, H., Pino, J., Bali, K. (eds.) Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing. pp. 1487–1505. Association for Computational Linguistics, Singapore (Dec 2023). https://doi.org/10.18653/v1/2023.emnlp-main.92, https://aclanthology.org/2023.emnlp-main.92/

15. Pangakis, N., Wolken, S., Fasching, N.: Automated annotation with generative ai requires validation (2023), https://arxiv.org/abs/2306.00176

16. Rassin, R., Fairstein, Y., Kalinsky, O., Kushilevitz, G., Cohen, N., Libov, A., Goldberg, Y.: Evaluating d-merit of partial-annotation on information retrieval (2024), https://arxiv.org/abs/2406.16048

17. Saad-Falcon, J., Khattab, O., Potts, C., Zaharia, M.: Ares: An automated evaluation framework for retrieval-augmented generation systems (2024), https://arxiv.org/abs/2311.09476

18. Shapiro, S.S., Wilk, M.B.: An analysis of variance test for normality (complete samples). Biometrika **52**(3/4), 591–611 (1965), http://www.jstor.org/stable/2333709

19. Sorodoc, I.T., Ribeiro, L.F.R., Blloshmi, R., Davis, C., de Gispert, A.: Garage: A benchmark with grounding annotations for rag evaluation (2025), https://arxiv.org/abs/2506.07671

20. Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., Rodriguez, A., Joulin, A., Grave,

E., Lample, G.: Llama: Open and efficient foundation language models (2023), https://arxiv.org/abs/2302.13971

21. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need (2023), https://arxiv.org/abs/1706.03762

22. Wang, D., Churchill, E., Maes, P., Fan, X., Shneiderman, B., Shi, Y., Wang, Q.: From human-human collaboration to human-ai collaboration: Designing ai systems that can work together with people. In: Extended Abstracts of the 2020 CHI Conference on Human Factors in Computing Systems. p. 1–6. CHI EA '20, Association for Computing Machinery, New York, NY, USA (2020). https://doi.org/10.1145/3334480.3381069, https://doi.org/10.1145/3334480.3381069

23. Wang, S., Tan, J., Dou, Z., Wen, J.R.: Omnieval: An omnidirectional and automatic rag evaluation benchmark in financial domain (2025), https://arxiv.org/abs/2412.13018