DEEP LITERATURE SURVEY AUTOMATION WITH AN ITERATIVE WORKFLOW

Anonymous authorsPaper under double-blind review

ABSTRACT

Automatic literature survey generation has attracted increasing attention, yet most existing systems follow a one-shot paradigm, where a large set of papers is retrieved at once and a static outline is generated before drafting. This design often leads to noisy retrieval, fragmented structures, and context overload, ultimately limiting survey quality. Inspired by the iterative reading process of human researchers, we propose IterSurvey, a framework based on recurrent outline generation, in which a planning agent incrementally retrieves, reads, and updates the outline to ensure both exploration and coherence. To provide faithful paper-level grounding, we design paper cards that distill each paper into its contributions, methods, and findings, and introduce a review-and-refine loop with visualization enhancement to improve textual flow and integrate multimodal elements such as figures and tables. Experiments on both established and emerging topics show that IterSurvey substantially outperforms state-of-the-art baselines in content coverage, structural coherence, and citation quality, while producing more accessible and better-organized surveys. To provide a more reliable assessment of such improvements, we further introduce Survey-Arena, a pairwise benchmark that complements absolute scoring and more clearly positions machine-generated surveys relative to human-written ones.

1 Introduction

Automatic literature survey generation has recently attracted growing attention due to its potential to help researchers quickly grasp new domains, identify key trends, and reduce the burden of manual reviews. Following Wang et al. (2024b), current systems generally adopt a multistage pipeline (Liang et al., 2025; Yan et al., 2025; Wang et al., 2025): The process begins with a topic description, usually consisting of a few tokens, which is directly used to retrieve a large collection of candidate papers. Due to the context window limitation of large language models (LLMs), the retrieved papers are divided into multiple groups, for each, an LLM agent generates a survey section outline based on the corresponding subset of papers. These group-level outlines are subsequently merged into a global draft outline. Once the draft outline is obtained, the system performs section-wise retrieval to collect references for section writing and then generates the corresponding text passages. Finally, a global review and integration process is applied, in which the drafted survey is iteratively polished to improve readability and overall consistency.

The above approach takes a "one-shot" planning paradigm, retrieves a comprehensive set of papers and construct a global outline from a single, static starting point. This approach, however, leads to several limitations. **First, retrieval can be imprecise and static** due to reliance on a short topic description (often just a few tokens) as the retrieval query (Sun et al., 2019; Azad & Deepak, 2019; Wang et al., 2020). Such coarse queries fail to capture a field's nuances and are never refined, leading to noisy and incomplete paper collections. **Second, the survey structure can be incoherent** (Fabbri et al., 2019; Gidiotis & Tsoumakas, 2020; Yang et al., 2023a). Since outlines are generated for each paper group independently and subsequently merged, the global structure lacks coherence and often misses important cross-group connections. **Third, injecting overly long contexts introduces distraction and context overload** (Liu et al., 2023; Wu et al., 2024). Feeding entire papers into LLMs not only exposes them to large amounts of peripheral information, such as dataset details or experimental setups, which distracts from the conceptual structure needed for survey writing, but also places unnecessary pressure on the limited context window of the model.

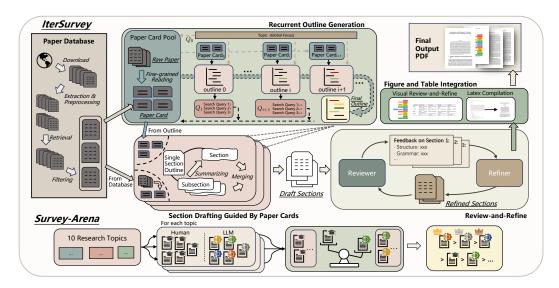


Figure 1: Overview of IterSurvey and Survey-Arena.

In contrast, human researchers rarely attempt to grasp an entire field in a single shot. Instead, they follow an iterative reading process: starting with a small set of core papers, summarizing key contributions, and gradually expanding to related directions as their understanding deepens (Bates, 1989; Asai et al., 2023). Inspired by this workflow, we propose an iterative planning paradigm for automated survey generation. At its core lies a **recurrent outline generation** module that incrementally retrieves, organizes, and integrates evidence through a planning agent equipped with stability checks and stopping criteria, mitigating the brittleness of one-shot pipelines that rely on static queries and fragmented merges. Central to this process are paper cards, structured semantic abstractions that distill each paper into contributions, methods, and findings. Unlike conventional abstract-based inputs, these cards serve as fine-grained evidence units that guide both outline construction and section drafting, ensuring coherence and faithful citation across iterations. Finally, a global review and integration stage employs a reviewer–refiner loop to enforce consistency and clarity across sections, while an integrated figure-table generation pipeline compiles candidate visualizations, automatically checks them for layout and readability, and revises them to meet academic presentation standards. This design inherits the advantages of iterative human reading: retrieval is progressively refined rather than static (Jiang et al., 2023), the outline develops as an organically coherent structure rather than a patchwork (Zhang et al., 2025a), and paper cards enforce fine-grained evidence grounding that avoids distraction from peripheral details (Cachola et al., 2020; Wu et al., 2024).

Comprehensive experiments validate the effectiveness of our incremental paradigm. IterSurvey consistently outperforms all baselines across multiple dimensions, with recurrent outline generation yielding more coherent structures and paper cards improving citation accuracy without sacrificing precision. These advantages are further confirmed by human evaluation, where experts also favor the outputs of IterSurvey over competing systems. While these results confirm the superiority of IterSurvey, we find that absolute scoring struggles to reliably quantify the performance gap against human-written surveys (Yang et al., 2023b; Oren et al., 2023; Ye et al., 2024). In the LLM evaluation community, similar concerns have led to the development of Chatbot Arena Chiang et al. (2024), which adopts pairwise human preference judgments to overcome the noisiness and inconsistency of absolute ratings. Inspired by this paradigm, we further contribute **Survey-Arena**, the first benchmark to our knowledge that evaluates synthesized surveys through direct, pairwise ranking against a corpus of human-written exemplars. This approach provides a more robust and interpretable assessment of system quality by directly positioning it relative to a human-level baseline.

Our contributions are threefold.

- We propose recurrent outline generation, which iteratively retrieves, reads, and updates
 outlines with paper cards and outline–paper grounding, while encouraging the model to
 explore new directions.
- We develop a new framework: IterSurvey, which produces finer-grained outlines and supports multi-modal inputs and outputs for more comprehensive surveys.

• We construct **Survey-Arena**, a pairwise evaluation benchmark that ranks machinegenerated surveys alongside human-written ones, enabling more reliable and interpretable assessment of survey quality.

2 RELATED WORK

Automated Survey Generation Recent automated survey generation systems largely adopt a "one-shot" paradigm, where a static outline is constructed upfront before content generation. This approach is evident in pipeline-based systems like AutoSurvey (Wang et al., 2024b), which employs a hierarchical paradigm, and SurveyForge (Yan et al., 2025), which utilizes a memory-driven scholar navigation agent. Other frameworks focus on enhancing this initial outlining step through reference pre-processing; for instance, SurveyX (Liang et al., 2025) introduces an AttributeTree to extract key information, while HiReview (Hu et al., 2024) generates a hierarchical taxonomy tree. Tackling the challenge from a technical scalability perspective, SurveyGo (Wang et al., 2025) leverages the LLM×MapReduce-V2 algorithm to handle long contexts within this paradigm. In contrast, our framework treats the outline not as a static blueprint but as an evolving knowledge structure. Through a dynamic, recurrent mechanism, the outline is continuously updated as the system iteratively engages with the literature, resulting in comprehensive and coherent synthesis.

Evaluation of Automated Surveys Evaluating machine-generated surveys is inherently challenging. Building on insights from automated peer review (Yu et al., 2024; Jin et al., 2024; Weng et al., 2025), prior works (Wang et al., 2024b; Yan et al., 2025; Liang et al., 2025) commonly adopt an LLM-as-a-judge paradigm with manually designed criteria, assessing dimensions such as coherence, coverage, and factuality. Citation quality is typically measured with NLI-based protocols (Gao et al., 2023), and Yan et al. (2025) additionally evaluate coverage by comparing system outputs with human-written surveys. While absolute scoring by LLMs provides useful fine-grained signals, it has also been noted to suffer from inconsistency and calibration issues (Ye et al., 2024; Latona et al., 2024), making system-level comparisons less reliable. In contrast, pairwise judgment which is widely used in chatbot evaluation (Zhao, 2025; Chiang et al., 2024) and peer review (Zhang et al., 2025b), offers more stable and interpretable assessments, but has not yet been applied to survey evaluation. To fill this gap, we introduce *Survey-Arena*, the first benchmark that ranks machine-generated surveys against human-written exemplars, providing both robust comparison across systems and a clearer positioning relative to human-level quality.

3 ITERSURVEY

An overview of IterSurvey is shown in Fig. 1, and its three core stages are detailed below.

3.1 RECURRENT OUTLINE GENERATION

Outline generation is a central component of automatic survey construction, as it requires understanding the research domain, identifying its subfields, and synthesizing individual papers. Alg. 1 shows the overview of the generation process. The outcome is a hierarchical framework that summarizes the domain, where each node in the hierarchy is represented by a title and an accompanying description. Given a topic query, our goal is to enable the model to integrate retrieval with inductive reasoning, so that it can systematically explore the literature and produce a comprehensive outline for the target domain. To this end, we design recurrent outline generation.

Paper Card Pool. The paper card pool organizes retrieval keywords together with their associated papers in a structured mapping. For each keyword K_i , we retrieve n candidate papers and extract m of the most relevant references, forming the set:

$$\mathcal{P}_i = \{p_i^1, p_i^2, \dots, p_i^{n+m}\}.$$

At iteration i, the system pops one keyword K_i together with its associated paper set \mathcal{P}_i from the pool. Each paper $p_i^j \in \mathcal{P}_i$ is converted into a paper card

$$c_i^j = \operatorname{PaperCard}(p_i^j),$$

185 186

187

188

189 190

191 192

193

196

197

199

200

201

202203204205

206

207208

209

210

211

212213

214

215

Algorithm 1 Description of the recurrent outline generation process.

```
163
                  Require: Topic query q; retrieval sizes (n, m); batch size B; paper budget (N_{\min}, N_{\max}); similarity threshold \tau
164
                  Ensure: Writing-oriented outline \hat{O}
                   1: O \leftarrow \text{InitOutLine}(q)
165
                   2: Pool ← ∅
3: U ← ∅
                                                                                                                                                                                                \triangleright map: query \mapsto card list
166
                                                                                                                                                                                                         4: R ← []

    puery history

167
                   5: for all r \in SEEDQUERIES(q) do
168
                              \mathcal{P} \leftarrow \text{RETRIEVE}(r, n) \cup \text{TOPREFS}(\cdot, m)
                              \mathcal{C} \leftarrow \{ \texttt{PaperCard}(p) \mid p \in \mathcal{P} \}
                   8:
                             \mathsf{Pool}[r] \leftarrow \mathcal{C}; \quad \mathcal{U} \leftarrow \mathcal{U} \cup \mathcal{P}
170
                   9: while |\mathcal{U}| < N_{\max} do
171
                  10:
                                if \mathsf{Pool} = \emptyset then
                                     if |\mathcal{U}| \geq N_{\min} and h(O, \mathbf{R}) then
172
                  12:
13:
                                          break
173
                                           for all r \in \text{EXPANDQUERIES}(O, \mathbf{R}) do
174
                   15:
                                                \mathcal{P} \leftarrow \mathsf{RETRIEVE}(r,n) \cup \mathsf{TOPREFS}(\cdot,m)
175
                  16:
                                                \mathcal{C} \leftarrow \{ \text{PaperCard}(p) \mid p \in \mathcal{P} \}
                  17:
                                                \mathsf{Pool}[r] \leftarrow \mathcal{C}; \quad \mathcal{U} \leftarrow \mathcal{U} \cup \mathcal{P}
176
                  18:
                                          continue
177
                  19:
                                (r, \mathcal{C}) \leftarrow \text{Pop}(\text{Pool})
                                                                                                                                                                                       > activate a query and its cards
                                \mathbf{R} \leftarrow \mathbf{R} \parallel r
178
                                while \mathcal{C} \neq \emptyset do
179
                                     \mathcal{B} \leftarrow \text{SampleBatch}(\mathcal{C}, B)
                  23:
                                     \tilde{O} \leftarrow g(O, \mathcal{B}, r)
                                                                                                                                                                                      > retrieval + reading + synthesis
                  24:
                                     if SIM(O, \tilde{O}) \ge \tau then
181
                   25:
                                          O \leftarrow \tilde{O}
182
                                     \mathcal{C} \leftarrow \mathcal{C} \setminus \mathcal{B}
                  27: \hat{O} \leftarrow \text{Refine}(O)
183
                  28: return Ô
```

which distills the paper into its contributions, methods, and findings. The collection of paper cards is denoted as $C_i = \{c_i^1, c_i^2, \dots, c_i^{|\mathcal{P}_i|}\}$. Overall, the paper card pool can be represented as a mapping

$$\mathcal{Q} = \{ K_i \mapsto \mathcal{C}_i \mid i = 0, 1, \dots \},\$$

where each keyword K_i is associated with the corresponding set of paper cards C_i .

Outline updating. The outline updating process begins with an empty initial outline, denoted as O_0 . At each step, the outline is refined using the current outline O_i , the active keyword K_i , and a mini-batch of paper cards drawn from the pool. Specifically, let $\mathcal{B}_i \subseteq \mathcal{C}_i$ be a batch of paper cards sampled from the set of cards associated with K_i . The model produces a candidate update

$$\tilde{O}_{i+1} = g(O_i, \mathcal{B}_i, K_i),$$

where $g(\cdot)$ denotes the outline updating function. This procedure is repeated iteratively, with batches \mathcal{B}_i of paper cards popped from the paper pool \mathcal{Q} under the current keyword K_i , until all cards associated with K_i are consumed and integrated into the outline. To ensure stability and promote refinement, the candidate update is accepted if its similarity to the previous outline exceeds τ :

$$O_{i+1} = \begin{cases} \tilde{O}_{i+1}, & \text{if } \text{Sim}(O_i, \tilde{O}_{i+1}) \geq \tau, \\ O_i, & \text{otherwise.} \end{cases}$$

Keyword expansion. When all keywords K_i has been fully consumed, the system explores new directions by proposing additional keywords. The goal is to identify potentially relevant aspects of the domain that have not yet been covered. Formally, new keywords are generated as

$$K_{i+1} = f(O_{i+1}, K_i, \dots, K_0),$$

where $f(\cdot)$ denotes a keyword generation function that takes the updated outline and the history of queries as input, and proposes candidate keywords for further exploration. The corresponding paper set \mathcal{P}_{i+1} is then retrieved and pushed into the pool \mathcal{Q} , thereby guiding the next iteration.

Stopping condition. Let $N_i = |\mathcal{P}_0 \cup \mathcal{P}_1 \cup \cdots \cup \mathcal{P}_i|$ denote the total number of consulted papers up to iteration i. The process terminates when either (i) $N_i \geq N_{\min}$ and the stopping signal

$$s = h(O_{i+1}, K_i, \dots, K_0), \quad s \in \{0, 1\},\$$

indicates that the outline is sufficiently complete, or (ii) $N_i \geq N_{\max}$. Here $h(\cdot)$ is a decision function which takes the evolving outline and the query history as input and outputs whether further exploration is necessary. This design ensures that the outline is not terminated prematurely, while also preventing excessive exploration.

Post-processing. After termination, the recurrent process produces a research-oriented outline \tilde{O} , which is further refined into a writing-oriented survey outline:

$$\hat{O} = \text{Refine}(\tilde{O}),$$

where $Refine(\cdot)$ reorganizes the structure, inserts standard survey components such as 'Introduction' and 'Future Directions', and ensures conformity with academic conventions. Finally, we perform paper–section relinking, where all consulted papers are reassociated with the corresponding sections of the final outline \hat{O} . This guarantees that each section of \hat{O} is grounded in concrete evidence, providing a reliable foundation for subsection drafting.

3.2 Section Drafting Guided by Paper Cards

A distinctive feature of our framework is that section drafting is entirely guided by paper cards, which serve as fine-grained, structured representations of the literature. Given the refined outline \hat{O} , each section or subsection is written by conditioning on its description d_j together with the relevant pool of cards. Specifically, for a given subsection with description d_j , the system retrieves a set of additional reference papers $\mathcal{P}_{\text{sec}}^j$ and converts them into paper cards $\mathcal{C}_{\text{sec}}^j$. In contrast to previous work, our framework benefits from the paper–section relinking established during outline construction: each subsection is already associated with a pool of consulted papers from earlier iterations. This enriched evidence base, combining $\mathcal{C}_{\text{sec}}^j$ with the relinked cards, provides the model with a stronger foundation for subsection writing. Formally, the j-th subsection is generated as

$$S_j = \operatorname{Draft}(d_j, \mathcal{C}^j_{\operatorname{sec}} \cup \mathcal{C}^j_{\operatorname{link}}),$$

where C_{link}^{j} denotes the set of paper cards relinked to subsection j. During drafting, the model is required to cite the provided references, and the citations are mapped to their corresponding papers.

3.3 GLOBAL REVIEW AND INTEGRATION

The final stage of survey generation goes beyond local drafting. It performs a global review-andrefine process that integrates sections into a coherent survey and enriches the survey with automatically generated figures and tables.

Textual Review-and-Refine. We adopt a reviewer–refiner loop that involves two collaborative LLM roles. The reviewer takes the entire survey draft as input to capture the global context but then focuses its critique on a specific section or subsection. This design ensures that feedback on local content is always grounded in an understanding of the overall narrative. The reviewer provides detailed suggestions covering aspects such as clarity of exposition, consistency of terminology, logical alignment with preceding and following sections, and stylistic fluency. The refiner then incorporates these suggestions to revise the targeted section, producing a polished update that fits better into the survey as a whole. This loop is applied sequentially across all sections and iterated multiple times, progressively enhancing readability, improving cross-section coherence, and strengthening the global structural integrity of the survey.

Figure–Table Integration. In addition to textual refinement, we extend the refinement process to include multimodal elements, to further enhance readability. For each section, the model first generates visualization requirements, such as tables with structured comparisons or figures with explanatory diagrams, together with natural language descriptions. Based on these descriptions, candidate figures and tables are synthesized. The compiled outputs are then fed back to an LLM for quality assessment, enabling automatic detection of issues such as oversized layouts or unreadable text. The LLM provides corrective suggestions, which are applied to improve the final visualizations. Finally, the text is refined again to ensure that all generated figures and tables are properly referenced within the survey.

4 EXPERIMENTS

4.1 EXPERIMENTAL SETTINGS

Implementation Details. Following Wang et al. (2024b), we adopt GPT-4o-mini as our generation model for its balance of responsiveness and cost. Our retrieval database contains 680K computer science papers from arXiv, with PDFs converted into structured Markdown using MinerU (Wang et al., 2024a) for consistent formatting. The details of the retrieval process are provided in App. A.1. In outline generation, the system consults 1000–1200 papers, with a maximum of 8 sections. For section drafting, each subsection retrieves up to 60 additional relevant papers, combined with those linked during outline generation. Finally, we apply two iterations of the review-and-refine loop to enhance coherence across sections and improve overall readability. Illustrative outputs compared with AutoSurvey are provided in App. A.8.

Baselines. We compare IterSurvey with a set of baselines, ranging from simple retrieval-augmented generation (Naive RAG), which directly drafts from retrieved documents, to more advanced state-of-the-art systems. Specifically, we evaluate against AutoSurvey (Wang et al., 2024b), the first systematic framework for this task; SurveyForge (Yan et al., 2025), which combines heuristic outline generation based on the logical structures of human-written surveys with a memory-driven scholar navigation agent for high-quality retrieval; and SurveyGo (Wang et al., 2025), which employs the LLM×MapReduce-V2 algorithm to address the long-context challenge. We also compare with SurveyX (Liang et al., 2025), which introduces an Attribute Tree-based outlining mechanism; however, due to access restrictions, we include SurveyX only in arena experiments. All methods are evaluated on the same retrieval database with generation hyperparameters aligned to their original settings for fairness.

4.2 AUTOMATIC EVALUATION RESULTS

Evaluation Setup. We employ multiple complementary protocols to evaluate the quality of generated surveys. On the 20-topic suite from Wang et al. (2024b), we adopt multi-dimensional scoring with LLM-as-a-judge. Content quality is assessed along three dimensions: coverage, structure, and relevance followed from Wang et al. (2024b). Besides, citation quality is evaluated using the NLI-based protocol of Gao et al. (2023), reporting both recall and precision: *Citation Recall* measures whether all statements in the generated text are fully supported by the cited passages, while *Citation Precision* identifies irrelevant citations to ensure that references are pertinent and directly support the claims. To improve scoring stability and reliability, prompts are standardized and judges must provide a rationale before assigning scores. For additional robustness, we aggregate outputs from three judge models: GPT-4o, Claude-3.5-Haiku, and GLM-4.5V. Full prompts are provided in App. A.7.

Results. The results on the 20 topics from Wang et al. (2024b) are reported in Tab. 1. Statistical significance was confirmed via paired t-tests, indicating that IterSurvey consistently outperforms baseline models (p < 0.05). We summarize the main observations below.

- Overall superiority. IterSurvey consistently outperforms all baselines across both content and citation quality, achieving the highest overall average score (4.75). This demonstrates that the proposed framework is effective and robust across multiple evaluation dimensions.
- Improved structural quality. On the structure dimension, IterSurvey achieves the best score (4.72). This improvement stems from the recurrent outline generation mechanism, which iteratively explores the literature and refines the outline, resulting in clearer organizational planning and stronger cross-sectional coherence.
- Enhanced citation quality. IterSurvey also achieves superior citation performance. While maintaining the same precision as AutoSurvey, it improves recall to 0.70. This advantage is enabled by paper cards, which provide fine-grained summaries of individual papers and thus allow for retrieving and citing a broader yet still accurate set of supporting references.

Together, these results confirm that recurrent outline generation, paper cards, and outline—paper grounding synergize to produce surveys that are both structurally coherent and rigorously evidenced.

¹Specifically, we use chatgpt-4o-latest, claude-3-5-haiku-20241022, and glm-4.5v.

326 327

328 330 331 332

333 334 335

336 337 338 339

340 341





347 348 349 350 351

352 353 354 355 356 357

358 359 360 361

362

364 366 367 368 369 370

371 372 373

374

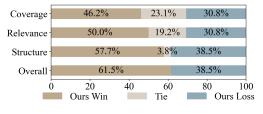
375

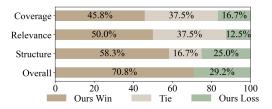
376

377

Table 1: Comparison of different methods in terms of content quality and citation quality.

| Methods | Content Quality | | | | Citation Quality | | |
|-------------|-------------------|-----------------|-------------------|-------------------|--------------------|-------------------|--|
| 1,100110415 | Coverage | Relevance | Structure | Avg. | Precision | Recall | |
| NaiveRAG | 4.42 ± 0.50 | 4.85 ± 0.36 | 4.20 ± 0.73 | 4.49 ± 0.41 | 0.39 ± 0.16 | $0.40_{\pm0.15}$ | |
| AutoSurvey | 4.50 ± 0.29 | 4.80 ± 0.16 | 4.62 ± 0.24 | 4.64 ± 0.15 | 0.64 ± 0.08 | 0.64 ± 0.08 | |
| SurveyForge | 4.57 ± 0.50 | 4.82 ± 0.39 | 4.60 ± 0.56 | 4.66 ± 0.40 | 0.59 ± 0.09 | 0.59 ± 0.09 | |
| SurveyGo | $4.37_{\pm 0.49}$ | 4.83 ± 0.38 | $4.27_{\pm 0.63}$ | $4.49_{\pm 0.40}$ | $0.50_{\pm 0.11}$ | $0.63_{\pm0.12}$ | |
| IterSurvey | 4.58 ± 0.50 | 4.95 ± 0.22 | 4.72 ± 0.45 | 4.75 ± 0.30 | 0.64 ± 0.06 | $0.70_{\pm 0.07}$ | |





(a) IterSurvey vs AutoSurvey

(b) IterSurvey vs SurveyForge

Figure 2: LLM-generated survey comparison between AutoSurvey and IterSurvey.

4.3 Human Evaluation Results

To further assess the quality of the generated surveys, we conducted a blind, pairwise study (Novikova et al., 2018; Chiang et al., 2024) with seven PhD-level experts. For each evaluation, experts were presented with an anonymized survey pair and asked to select the superior one based on multiple quality dimensions, including coverage, relevance, structural coherence, and overall quality, which is more objective and stable than ranking based on absolute scores (Herbrich et al., 2006; Sakaguchi et al., 2014). To control annotation cost, the human study was limited to direct comparisons between IterSurvey and two leading baselines: AutoSurvey and SurveyForge. Inter-rater agreement is reported in App. A.2. Results, as shown in Fig. 2, indicate that IterSurvey is consistently preferred over AutoSurvey and SurveyForge by domain experts, especially in terms of structure and overall quality. This trend aligns with our automatic evaluation, where recurrent outline generation also demonstrated stronger coherence and organization. The consistency between expert judgments and automatic metrics further highlights the robustness of IterSurvey in generating high-quality surveys.

4.4 SURVEY-ARENA: PAIRWISE COMPARISON AND RANKING

Dataset construction. Previous automatic evaluation methods typically assign an absolute score for each dimension, which struggles to fully capture the performance gap between machinegenerated surveys and human-written ones. To move beyond absolute scores, we constructed the Survey-Arena benchmark. The benchmark spans ten research topics. For each topic, we manually selected five high-quality, human-written surveys to serve as a performance baseline. To ensure comparability, all surveys for a given topic were chosen from a narrow six-month submission window, a process that required careful verification to ensure each topic had a sufficient number of suitable papers. We further confirmed their quality and influence via non-trivial citation counts on Google Scholar. The retrieval database for all machine-generated surveys was correspondingly frozen to the same time period to guarantee fairness. The full list of topics and papers is available in the App. A.4.

Evaluation protocol. For each topic, all possible pairs of a machine-generated survey and a human-written survey are constructed. To ensure robust evaluation and mitigate positional bias, each pair is judged in both directions (A vs. B and B vs. A), following Li et al. (2024). A panel of three distinct LLMs, namely GPT-4o, Claude-3.5-Haiku, and GLM-4.5V, serves as the judges for each comparison. Elo scores are computed from these aggregated pairwise outcomes to generate rankings for all systems.

Results. We report two key evaluation metrics: Avg. Rank, which indicates the mean position among all surveys, and >Human%, which reflects the proportion of topics where a system surpasses human surveys. The topic-wise outcomes from Survey-Arena are visualized in Fig.3, and the aggregated rankings are summarized in Tab.2.

Each system is evaluated by its average rank across all surveys (including 5 machine-written surveys and 5 human-written ones) and by the proportion of topics where it surpasses human surveys. The results show that IterSurvey consistently achieves the best overall performance among automatic survey generation systems, with an average rank of 4.0 and surpassing human-written surveys in 60% of topics. These findings highlight that IterSurvey not only outperforms competing methods but also approaches human-level quality across diverse domains.

Meta Evaluation. To assess the reliability of Survey-Arena judgments, we compare the rankings produced by Survey-Arena for humanwritten surveys with citation counts on Google Scholar, which serve as an external signal of impact. Specifically, we compute Spearman's ρ_s by measuring the correlation between Arena-derived and citation-based rankings for each topic, and then report the average across topics. For relevance scoring, we treat citation counts as an indicator of relevance and compute nDCG directly over the ranking lists. As a comparison, we also use the rankings derived from absolute scoring and compute their consistency and nDCG. This allows us to evaluate how well the different ranking methods align with citation-based rankings.

Results are shown in Tab. 3. Compared with the scoring-based approach, pairwise judgment achieves higher agreement with citation-based rankings, yielding a Spearman's ρ_s of 0.410 and nDCG@2/3 = 0.834/0.873. This indicates that when models are asked to directly compare two surveys, they more reliably identify the superior one, producing rankings that better align with human impact signals. These findings support pairwise evaluation as a more robust protocol for Survey-Arena.

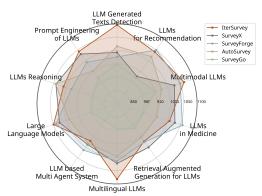


Figure 3: Elo scores of Survey-Arena results across topics. The radar plot shows the Elo scores for each system across all topics, providing a topic-wise comparison.

Table 2: Aggregated rankings on Survey-Arena. Avg. Rank is the mean position among all surveys. >Human% is the average proportion of topics where a system surpasses human surveys.

| Method | Avg. Rank↓ | $>$ Human $\% \uparrow$ |
|-------------|------------|-------------------------|
| SurveyGo | 9.80 | 4% |
| AutoSurvey | 6.70 | 32% |
| SurveyForge | 4.80 | 50% |
| SurveyX | 4.70 | 54% |
| IterSurvey | 4.00 | 60% |

Table 3: Consistency between different ranking methods and citation-based rankings.

| Rank Method | $oldsymbol{ ho}_{ m s}$ | nDCG@2 | 2 nDCG@3 |
|------------------|-------------------------|--------|----------|
| Absolute Scoring | 0.320 | 0.695 | 0.767 |
| Pair-Judge | 0.410 | 0.834 | 0.873 |

4.5 GENERALIZATION ON SURVEY-LACKING TOPICS

To examine whether automated survey generation can succeed in areas without existing surveys, we construct a subset of eight research topics (listed in App. A.5) where no human-written reviews are available. Such settings are common in emerging domains and pose greater challenges, since there are no canonical structures to imitate and the literature is often sparse and fragmented. This setup tests whether a system can autonomously organize the field into a coherent, well-grounded survey.

We compare IterSurvey against AutoSurvey and SurveyForge under this setup, and the results are presented in Tab. 4. Our method achieves the highest average score (4.63), consistently outperforming both baselines across content and citation quality. Notably, IterSurvey shows clear advantages in structural quality (4.63) and citation recall (0.67). These gains highlight the benefits of recurrent outline generation, which encourages iterative query expansion and literature exploration rather

Table 4: Comparison of different methods on survey-lacking topics.

| Methods | Content Quality | | | | Citation Quality | |
|--------------|------------------|--------------------------|------------------------|--------------------------|------------------------------|------------------------------|
| 1/10/11/0/15 | Coverage | Relevance | Structure | Avg. | Precision | Recall |
| AutoSurvey | $4.00_{\pm1.12}$ | 4.20 ± 1.20 | 4.00 ± 1.00 | 4.07 ± 1.11 | 0.55 ± 0.14 | 0.55 ± 0.09 |
| SurveyForge | 4.50 ± 0.50 | 4.75 ± 0.50 | 4.54 ± 0.54 | 4.60 ± 0.52 | 0.47 ± 0.12 | 0.47 ± 0.13 |
| IterSurvey | 4.42 ± 0.58 | $\textbf{4.83} \pm 0.17$ | 4.63 \pm 0.63 | $\textbf{4.63} \pm 0.37$ | $\boldsymbol{0.60} \pm 0.06$ | $\boldsymbol{0.67} \pm 0.06$ |

Table 5: Ablation study analyzing the contribution of each component in IterSurvey: \bigcirc Recurrent outline generation; \square Paper Card; \diamondsuit Review-and-Refine.

| Methods | Content Quality | | | | Citation Quality | |
|------------------------|------------------------|-----------------|-------------------|-------------------|------------------|-----------------|
| 1,100110015 | Coverage | Relevance | Structure | Avg. | Precision | Recall |
| Baseline | $4.00_{\pm 0.53}$ | 4.40 ± 0.48 | $4.20_{\pm 0.70}$ | $4.20_{\pm 0.44}$ | 0.58 ± 0.09 | 0.67 ± 0.09 |
| + () | 4.46 ± 0.52 | 4.80 ± 0.41 | 4.53 ± 0.52 | 4.60 ± 0.40 | 0.62 ± 0.08 | 0.59 ± 0.09 |
| + () + () | 4.60 ± 0.51 | 4.80 ± 0.42 | 4.60 ± 0.52 | 4.69 ± 0.39 | 0.64 ± 0.09 | 0.71 ± 0.08 |
| + \() + \(\) + \(\) | 4.73 ± 0.50 | 4.93 ± 0.41 | 4.80 ± 0.52 | 4.82 ± 0.39 | 0.65 ± 0.04 | 0.77 ± 0.04 |

than relying on a fixed set of initial retrievals. Combined with paper cards providing fine-grained evidence abstraction, this mechanism enables IterSurvey to construct coherent survey structures and incorporate broader supporting references even in areas where survey conventions are absent.

4.6 ABLATION STUDY

We conducted an ablation study on five representative topics to analyze the impact of the three new modules of IterSurvey: Recurrent Outline Generation, Paper Card, and Review-and-Refine. Results are shown in Tab. 5, revealing the following insights:

Recurrent Outline Generation yields stronger content quality. We compare our recurrent outline generation with a one-shot paradigm, where retrieved papers are partitioned into groups, each group produces an outline independently, and the results are subsequently merged. The recurrent approach contributes significant improvements in content quality, with gains of +0.46 in coverage and +0.33 in structure over the baseline. This demonstrates that iterative exploration helps the model achieve broader coverage and stronger organizational coherence by progressively integrating evidence.

Paper Card improves citation quality. We further examined the impact of paper card, we replace them with abstract-based inputs commonly used in retrieval pipelines. The results show that paper cards significantly improve citation grounding, raising recall from 0.59 to 0.71 while maintaining precision (0.64). This indicates that distilled paper-level evidence reduces distraction and enables the model to retrieve and cite a broader yet accurate set of references.

Review and Refine boosts overall performance. Finally, we evaluate the review-and-refine stage by removing it from the pipeline. The full variant enhances all dimensions of content quality, raising the overall average from 4.69 to 4.82, and further improves citation recall from 0.71 to 0.77. These gains show that multi-round self-critique and revision help fill evidence gaps, eliminate unsupported claims, and polish the text into well-substantiated surveys. Together, recurrent planning, paper cards, and review-and-refine form the most effective configuration of IterSurvey.

5 CONCLUSION

In this work, we tackled the limitations of existing survey generation systems by introducing Iter-Survey, a framework with recurrent outline generation, paper cards, and global review and integration. This design enables precise retrieval, coherent structure, and faithful citation grounding, while supporting multimodal outputs. Experiments on diverse topics show that IterSurvey outperforms state-of-the-art baselines in coherence, coverage, and citation quality. We also proposed Survey-Arena, a pairwise benchmark that complements absolute scoring for a more reliable assessment. Future work will extend our framework to broader domains, integrate richer multimodal evidence, and refine evaluation protocols toward human-level quality.

ETHICAL CONSIDERATIONS

Our work focuses on automatic literature survey generation using large language models. While the system is designed to support researchers by synthesizing existing knowledge, it inevitably inherits limitations of current models, including potential citation errors, incomplete coverage, and occasional inaccuracies. Therefore, the generated surveys are intended as an assistive tool rather than a substitute for human scholarship, and should be used for reference only. For evaluation, all human experts involved in the study participated voluntarily and received fair compensation. All data used in our experiments were sourced from publicly available arXiv papers, which permit noncommercial use. We strictly avoided the use of private or sensitive data.

USE OF LARGE LANGUAGE MODELS

We used large language models (GPT-4o, Claude-3.5-Haiku, and GLM-4.5V) in two ways: (i) as evaluation judges for assessing survey quality, and (ii) for limited language editing and refinement of the manuscript. All substantive research ideas, experimental design, analyses, and final decisions were made solely by the authors, who take full responsibility for the content of this paper.

REFERENCES

- Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. Self-rag: Learning to retrieve, generate, and critique through self-reflection, 2023. URL https://arxiv.org/abs/2310.11511.
- Hiteshwar Kumar Azad and Akshay Deepak. Query expansion techniques for information retrieval: A survey. *Information Processing & Management*, 56(5):1698–1735, 2019. ISSN 0306-4573. doi: https://doi.org/10.1016/j.ipm.2019.05.009. URL https://www.sciencedirect.com/science/article/pii/S0306457318305466.
- Marcia J. Bates. The design of browsing and berrypicking techniques for the online search interface. *Online Review*, 1989. URL https://pages.gseis.ucla.edu/faculty/bates/berrypicking.html.
- Isabel Cachola, Kyle Lo, Arman Cohan, and Daniel S. Weld. Tldr: Extreme summarization of scientific documents, 2020. URL https://arxiv.org/abs/2004.15011.
- Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li, Dacheng Li, Hao Zhang, Banghua Zhu, Michael Jordan, Joseph E. Gonzalez, and Ion Stoica. Chatbot arena: An open platform for evaluating Ilms by human preference, 2024. URL https://arxiv.org/abs/2403.04132.
- Alexander Fabbri, Irene Li, Tianwei She, Suyi Li, and Dragomir Radev. Multi-news: A large-scale multi-document summarization dataset and abstractive hierarchical model. In Anna Korhonen, David Traum, and Lluís Màrquez (eds.), *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 1074–1084, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1102. URL https://aclanthology.org/P19-1102/.
- Tianyu Gao, Howard Yen, Jiatong Yu, and Danqi Chen. Enabling large language models to generate text with citations. *arXiv preprint arXiv:2305.14627*, 2023.
- Alexios Gidiotis and Grigorios Tsoumakas. A divide-and-conquer approach to the summarization of long documents, 2020. URL https://arxiv.org/abs/2004.06190.
- Ralf Herbrich, Tom Minka, and Thore Graepel. TrueskillTM: a bayesian skill rating system. *Advances in neural information processing systems*, 19, 2006.
- Yuntong Hu, Zhuofeng Li, Zheng Zhang, Chen Ling, Raasikh Kanjiani, Boxin Zhao, and Liang Zhao. Hireview: Hierarchical taxonomy-driven automatic literature review generation. 2024.

- Zhengbao Jiang, Frank F. Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang, Jamie Callan, and Graham Neubig. Active retrieval augmented generation, 2023. URL https://arxiv.org/abs/2305.06983.
 - Yiqiao Jin, Qinlin Zhao, Yiyang Wang, Hao Chen, Kaijie Zhu, Yijia Xiao, and Jindong Wang. Agentreview: Exploring peer review dynamics with llm agents, 2024. URL https://arxiv.org/abs/2406.12708.
 - Giuseppe Russo Latona, Manoel Horta Ribeiro, Tim R. Davidson, Veniamin Veselovsky, and Robert West. The ai review lottery: Widespread ai-assisted peer reviews boost paper scores and acceptance rates, 2024. URL https://arxiv.org/abs/2405.02150.
 - Zongjie Li, Chaozheng Wang, Pingchuan Ma, Daoyuan Wu, Shuai Wang, Cuiyun Gao, and Yang Liu. Split and merge: Aligning position biases in Ilm-based evaluators, 2024. URL https://arxiv.org/abs/2310.01432.
 - Xun Liang, Jiawei Yang, Yezhaohui Wang, Chen Tang, Zifan Zheng, Shichao Song, Zehao Lin, Yebin Yang, Simin Niu, Hanyu Wang, et al. Surveyx: Academic survey automation via large language models. *arXiv preprint arXiv:2502.14776*, 2025.
 - Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. Lost in the middle: How language models use long contexts, 2023. URL https://arxiv.org/abs/2307.03172.
 - Jekaterina Novikova, Ondřej Dušek, and Verena Rieser. RankME: Reliable human ratings for natural language generation. In *NAACL-HLT*, 2018. URL https://aclanthology.org/N18-2012.
 - Zach Nussbaum, John X. Morris, Brandon Duderstadt, and Andriy Mulyar. Nomic embed: Training a reproducible long context text embedder, 2024.
 - Yonatan Oren, Nicole Meister, Niladri Chatterji, Faisal Ladhak, and Tatsunori B. Hashimoto. Proving test set contamination in black box language models, 2023. URL https://arxiv.org/abs/2310.17623.
 - Keisuke Sakaguchi, Matt Post, and Benjamin Van Durme. Efficient elicitation of annotations for human evaluation of machine translation. In Ondřej Bojar, Christian Buck, Christian Federmann, Barry Haddow, Philipp Koehn, Christof Monz, Matt Post, and Lucia Specia (eds.), *Proceedings of the Ninth Workshop on Statistical Machine Translation*, pp. 1–11, Baltimore, Maryland, USA, June 2014. Association for Computational Linguistics. doi: 10.3115/v1/W14-3301. URL https://aclanthology.org/W14-3301/.
 - Haitian Sun, Tania Bedrax-Weiss, and William Cohen. PullNet: Open domain question answering with iterative retrieval on knowledge bases and text. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan (eds.), *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 2380–2390, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1242. URL https://aclanthology.org/D19-1242/.
 - Bin Wang, Chao Xu, Xiaomeng Zhao, Linke Ouyang, Fan Wu, Zhiyuan Zhao, Rui Xu, Kaiwen Liu, Yuan Qu, Fukai Shang, et al. Mineru: An open-source solution for precise document content extraction. *arXiv preprint arXiv:2409.18839*, 2024a.
 - Haoyu Wang, Yujia Fu, Zhu Zhang, Shuo Wang, Zirui Ren, Xiaorong Wang, Zhili Li, Chaoqun He, Bo An, Zhiyuan Liu, et al. Llmxmapreduce-v2: Entropy-driven convolutional test-time scaling for generating long-form articles from extremely long resources. *arXiv preprint arXiv:2504.05732*, 2025.
 - Xiao Wang, Craig Macdonald, and Iadh Ounis. Deep reinforced query reformulation for information retrieval, 2020. URL https://arxiv.org/abs/2007.07987.

- Yidong Wang, Qi Guo, Wenjin Yao, Hongbo Zhang, Xin Zhang, Zhen Wu, Meishan Zhang, Xinyu Dai, Min Zhang, Qingsong Wen, Wei Ye, Shikun Zhang, and Yue Zhang. Autosurvey: Large language models can automatically write surveys. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024b.
- Yixuan Weng, Minjun Zhu, Guangsheng Bao, Hongbo Zhang, Jindong Wang, Yue Zhang, and Linyi Yang. Cycleresearcher: Improving automated research via automated review, 2025. URL https://arxiv.org/abs/2411.00816.
- Zijun Wu, Bingyuan Liu, Ran Yan, Lei Chen, and Thomas Delteil. Reducing distraction in long-context language models by focused learning, 2024. URL https://arxiv.org/abs/2411.05928.
- Xiangchao Yan, Shiyang Feng, Jiakang Yuan, Renqiu Xia, Bin Wang, Lei Bai, and Bo Zhang. SURVEYFORGE: On the outline heuristics, memory-driven generation, and multi-dimensional evaluation for automated survey writing. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar (eds.), *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 12444–12465, Vienna, Austria, July 2025. Association for Computational Linguistics. ISBN 979-8-89176-251-0. doi: 10.18653/v1/2025.acl-long.609. URL https://aclanthology.org/2025.acl-long.609/.
- Kevin Yang, Dan Klein, Nanyun Peng, and Yuandong Tian. DOC: Improving long story coherence with detailed outline control. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 3378–3465, Toronto, Canada, July 2023a. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.190. URL https://aclanthology.org/2023.acl-long.190/.
- Shuo Yang, Wei-Lin Chiang, Lianmin Zheng, Joseph E. Gonzalez, and Ion Stoica. Rethinking benchmark and contamination for language models with rephrased samples, 2023b. URL https://arxiv.org/abs/2311.04850.
- Rui Ye, Xianghe Pang, Jingyi Chai, Jiaao Chen, Zhenfei Yin, Zhen Xiang, Xiaowen Dong, Jing Shao, and Siheng Chen. Are we there yet? revealing the risks of utilizing large language models in scholarly peer review, 2024. URL https://arxiv.org/abs/2412.01708.
- Jianxiang Yu, Zichen Ding, Jiaqi Tan, Kangyang Luo, Zhenmin Weng, Chenghua Gong, Long Zeng, Renjing Cui, Chengcheng Han, Qiushi Sun, Zhiyong Wu, Yunshi Lan, and Xiang Li. Automated peer reviewing in paper sea: Standardization, evaluation, and analysis, 2024. URL https://arxiv.org/abs/2407.12857.
- Dingchu Zhang, Yida Zhao, Jialong Wu, Baixuan Li, Wenbiao Yin, Liwen Zhang, Yong Jiang, Yufeng Li, Kewei Tu, Pengjun Xie, and Fei Huang. Evolvesearch: An iterative self-evolving search agent, 2025a. URL https://arxiv.org/abs/2505.22501.
- Yaohui Zhang, Haijing Zhang, Wenlong Ji, Tianyu Hua, Nick Haber, Hancheng Cao, and Weixin Liang. From replication to redesign: Exploring pairwise comparisons for llm-based peer review, 2025b. URL https://arxiv.org/abs/2506.11343.
- Zhimin Zhao. Se arena: An interactive platform for evaluating foundation models in software engineering. In 2025 IEEE/ACM Second International Conference on AI Foundation Models and Software Engineering (Forge), pp. 78–81, 2025. doi: 10.1109/Forge66646.2025.00016.

A APPENDIX

A.1 RETRIEVAL SETUP

For the retrieval process, we implemented a lightweight database to provide the necessary functionality. The retrieval logic is based on vector similarity, using the *nomic-ai/nomic-embed-text-v1.5* (Nussbaum et al., 2024) embedding model with all hyperparameters set to their default values. Given a query, the database computes the similarity between the query vector and all paper vectors, and returns the top-k most relevant entries. In addition, the database supports bidirectional lookup between a paper's arXiv identifier and title, as well as filtering papers published prior to a specified cutoff date.

A.2 RESULTS OF INTER-RATER AGREEMENT

To assess the reliability of human annotations, we computed Cohen's kappa coefficients across four evaluation dimensions: Coverage, Relevance, Structure, and Overall, as shown in Tab. 6. These results indicate substantial agreement among human annotators, supporting the consistency of the human evaluation process.

Table 6: Inter-rater agreement among human annotators.

| Dimensions | Coverage | Relevance | Structure | Overall |
|------------|----------|-----------|-----------|---------|
| kappa | 0.714 | 0.583 | 0.611 | 0.650 |

A.3 TOPICS FOR AUTOMATIC EVALUATION

We utilize 20 topics derived from AutoSurvey (Wang et al., 2024b). Each topic is paired with a human survey, as shown in Tab. 7, which also reports the survey titles, arXiv IDs, and their latest citation counts from Google Scholar.

Table 7: Topics for Automatic Evaluation

| Topic | Human Survey | ArXiv ID | Citations |
|-----------------------------------|--|------------|-----------|
| In-context Learning | A Survey on In-context Learning | 2301.00234 | 2396 |
| LLMs for Recommendation | A Survey on Large Language Models for Recommendation | 2305.19860 | 596 |
| LLM-Generated Texts Detection | The Science of Detecting LLM-Generated Texts | 2310.14724 | 308 |
| Explainability for LLMs | Explainability for Large Language Models: A Survey | 2309.01029 | 875 |
| Evaluation of LLMs | A Survey on Evaluation of Large Language Models | 2307.03109 | 4020 |
| LLMs-based Agents | A Survey on Large Language Model based Autonomous Agents | 2308.11432 | 1906 |
| LLMs in Medicine | A Survey of Large Language Models in Medicine | 2311.05112 | 217 |
| Domain Specialization of LLMs | Domain Specialization as the Key to Make Large Language Models Disruptive | 2305.18703 | 217 |
| Challenges of LLMs in Education | Practical and Ethical Challenges of Large Language Models in Education | 2303.13379 | 722 |
| Alignment of LLMs | Aligning Large Language Models with Human: A Survey | 2307.12966 | 435 |
| ChatGPT | Harnessing the Power of LLMs in Practice: A Survey on ChatGPT and Beyond | 2304.13712 | 1254 |
| Instruction Tuning for LLMs | Instruction Tuning for Large Language Models: A Survey | 2308.10792 | 1174 |
| LLMs for Information Retrieval | Large Language Models for Information Retrieval: A Survey | 2308.07107 | 544 |
| Safety in LLMs | Towards Safer Generative Language Models | 2302.09270 | 13 |
| Chain of Thought | A Survey of Chain of Thought Reasoning: Advances, Frontiers and Future | 2309.15402 | 290 |
| Hallucination in LLMs | A Survey on Hallucination in Large Language Models | 2311.05232 | 2599 |
| Bias and Fairness in LLMs | Bias and Fairness in Large Language Models: A Survey | 2309.00770 | 1009 |
| Large Multi-Modal Language Models | Large-scale Multi-Modal Pre-trained Models: A Comprehensive Survey | 2302.10035 | 285 |
| Acceleration for LLMs | A Survey on Model Compression and Acceleration for Pretrained Language Models | 2202.07105 | 101 |
| LLMs for Software Engineering | Large Language Models for Software Engineering: A Systematic Literature Review | 2308.10620 | 1058 |

A.4 TOPICS FOR SURVEY-ARENA

To construct the Survey-Arena benchmark, we select 10 topics, with several derived from Auto-Survey (Wang et al., 2024b) and SurveyForge (Yan et al., 2025). For each topic, we include 5 human-written surveys, requiring that their arXiv submission dates fall within a six-month window. We report their latest Google Scholar citation counts as a measure of impact, as summarized in Tab. 8. For reproducibility, we also specify the exact arXiv version, since submission dates can vary considerably across different versions of the same paper.

Table 8: Topics for Survey-Arena

| Topic | Human Survey | ArXiv ID | Citations |
|----------------------------|---|--|-------------------|
| | Large Language Models: A Survey | 2402.06196v3 | 1133 |
| Large Language | Large Language Models Meet NLP: A Survey | 2405.12819v1 | 86 |
| Models | History, Development, and Principles of Large Language Models-An Introductory Survey | 2402.06853v2 | 73 |
| WIOGCIS | Recent Advances in Generative AI and Large Language Models | 2407.14962v1 | 68 |
| | Exploring the landscape of large language models: Foundations, techniques, and challenges | 2404.11973v1 | 5 |
| | MM-LLMs: Recent Advances in MultiModal Large Language Models | 2401.13601v3 | 381 |
| | Multimodal Large Language Models: A Survey | 2311.13165v1 | 299 |
| Multimodal LLMs | The Revolution of Multimodal Large Language Models: A Survey | 2402.12451v1 | 98 |
| | How to Bridge the Gap between Modalities: Survey on Multimodal Large Language Model | 2311.07594v1 | 43 |
| | A Review of Multi-Modal Large Language and Vision Models | 2404.01322v1 | 39 |
| | Multilingual Large Language Model: A Survey of Resources, Taxonomy and Frontiers | 2404.04925v1 | 83 |
| | A Survey on Multilingual Large Language Models: Corpora, Alignment, and Bias | 2404.00929v2 | 55 |
| Multilingual LLMs | A Survey on Large Language Models with Multilingualism | 2405.10936v1 | 40 |
| | Surveying the MLLM Landscape: A Meta-Review of Current Surveys | 2409.18991v1 | 12 |
| | Multilingual Large Language Models: A Systematic Survey | 2411.11072v2 | 9 |
| | A Survey of Long Chain-of-Thought for Reasoning Large Language Models | 2503.09567v3 | 130 |
| **** | From System 1 to System 2: A Survey of Reasoning Large Language Models | 2502.17419v2 | 110 |
| LLMs Reasoning | Advancing Reasoning in Large Language Models: Promising Methods and Approaches | 2502.03671v1 | 19 |
| | A Survey of Frontiers in LLM Reasoning | 2504.09037v1 | 17 |
| | Thinking Machines: A Survey of LLM based Reasoning Strategies | 2503.10814v1 | 9 |
| | A Systematic Survey of Prompt Engineering in Large Language Models | 2402.07927v1 | 748 |
| Prompt Engineering | The Prompt Report: A Systematic Survey of Prompt Engineering Techniques | 2406.06608v2 | 182 |
| of LLMs | Prompt Design and Engineering: Introduction and Advanced Methods | 2401.14423v4 | 117 |
| OI EEMS | A Survey of Prompt Engineering Methods in Large Language Models for Different NLP Tasks | 2407.12994v1 | 60 |
| | Efficient Prom pting Methods for Large Language Models: A Survey | 2404.01077v1 | 56 |
| | Retrieval-Augmented Generation for Large Language Models: A Survey | 2312.10997v5 | 2583 |
| Retrieval-Augmented | A Survey on RAG Meeting LLMs: Towards Retrieval-Augmented Large Language Models | 2405.06211v3 | 559 |
| Generation for LLMs | A Survey on Retrieval-Augmented Text Generation for Large Language Models | 2404.10981v2 | 119 |
| | Retrieval-Augmented Generation for Natural Language Processing: A Survey | 2407.13193v2 | 77 |
| | Retrieval Augmented Generation (RAG) and Beyond | 2409.14924v1 | 70 |
| | A survey on large language model based autonomous agents | 2308.11432v7 | 1623 |
| LLM-based | Multi-Agent Collaboration Mechanisms: A Survey of LLMs | 2501.06322v1 | 79 |
| Multi-Agent System | Large language model agent: A survey on methodology, applications and challenges | 2503.21460v1 | 19 |
| | Agentic large language models, a survey | 2503.23037v2 | 12 |
| | A Survey on LLM-based Multi-Agent System: | 2412.17481v2 | 3 |
| | A Survey on LLM-Generated Text Detection: Necessity, Methods, and Future Directions | 2310.14724v2 | 210 |
| LLM-Generated | A Survey on Detection of LLMs-Generated Content | 2310.15654v1 | 69 |
| Texts Detection | Towards Possibilities & Impossibilities of AI-generated Text Detection: A Survey | 2310.15264v1 | 46 |
| | Detecting chatgpt: A survey of the state of detecting chatgpt-generated text | 2309.07689v1 | 22 |
| | Decoding the AI Pen: Techniques and Challenges in Detecting AI-Generated Text | 2403.05750v1 | 13 |
| | Large language models in healthcare and medical domain: A review | 2401.06775v2 | 246 |
| | A Survey on Medical Large Language Models | 2406.03712v1 | 53 |
| LLMs in Medicine | A Comprehensive Survey of Large Language Models and Multimodal Large Language Models in Medicine | 2405.08603v1 | 46 |
| | Large Language Models for Medicine: A Survey | 2405.13055v1 | 37 |
| | A Comprehensive Survey on Evaluating Large Language Model Applications in the Medical Industry | 2403.15033V1 2404.15777v4 | 32 |
| | A Survey on Large Language Models for Recommendation | 2305.19860v4 | 508 |
| | | | |
| | Pecommender Systems in the Fra of Large Language Models (LLMs) | 2307 02046*2 | |
| LLMs for | Recommender Systems in the Era of Large Language Models (LLMs) A Comprehensive Survey of Language Modelling Paradigm Adaptations in Recommender Systems | 2307.02046v2 2302.03735v3 | 479 117 |
| LLMs for Recommendation | Recommender Systems in the Era of Large Language Models (LLMs) A Comprehensive Survey of Language Modelling Paradigm Adaptations in Recommender Systems Large Language Models for Generative Recommendation: A Survey and Visionary Discussions | 2307.02046v2 2302.03735v3 2309.01157v1 | 479 117 116 |

A.5 TOPICS FOR SURVEY-LACKING TEST

756

758

759 760

761 762

764

765

766

767

768 769 770

771 772

773

774

775

776 777

778779

781

782 783

784 785 786

787

788 789 790

791

792

793

794

796

797

798

799

800

801

802 803

804

805

807 808

809

We manually select 8 topics with no existing survey articles, as shown in Tab. 9.

Table 9: Topics for Survey-Lacking Test.

Topic

Event Timeline Generation
Linear RNN in Natural Language Processing
Agent-flow Data Curation
Causal Mediation with Sparse Autoencoder Features in Transformers
Multi-Tenant Scheduling for MoE Inference
Benchmarking Tool-Using LLMs for Causal Tasks in the MCP Ecosystem
RAG for Mechanical Design: Cross-Modal Retrieval over CAD Trees and BOMs
Renderer-in-the-Loop Supervision for Multimodal Model

A.6 DETAIL OF NAIVE RAG

Naive RAG Prompt

Given a topic, the Naive RAG system first retrieves 1,500 papers from the same database as ours. It then employs an iterative prompting strategy, where the LLM generates content until the total length of the survey reaches 5,000 tokens (Wang et al., 2024b). The prompt used for generation is shown below.

You are an expert in artificial intelligence who wants to write an overall and comprehensive survey about [TOPIC]. You are provided with a list of papers related to [TOPIC] below: [PAPER LIST] Here is the survey content you have written: [SURVEY CONTENT] Hers is the requirement of the survey: 1. The survey must be more than [SURVEY LEN] tokens! 2. Containing serval sections. Each section contains several subsections. 3. Cite several paper provided above to support the content you write. Here is the format of your writing: 1. ## indicates the section title 2. ### indicates the subsection title 3. Only cite the "paper_title" in []. An example of citation: the emergence of large language models (LLMs) [Language models are few-shot learners; Language models are unsupervised multitask learners; PaLM: Scaling language modeling with pathways] You need to continue writing the survey by adding a new section or subsection. Do not stop until the length of survey is more than [SURVEY LEN] tokens!!!

Return the content you write:

A.7 PROMPTS FOR EVALUATION

 <topic>

[TOPIC]

NLI Prompt --Claim: [CLAIM] --Source: [SOURCE] --Claim: [CLAIM] --Is the Claim faithful to the Source? A Claim is faithful to the Source if the core part in the Claim can be supported by the Source.\n Only reply with 'Yes' or 'No':

Criteria-based judging survey prompt

```
You are an expert academic evaluator specializing in rigorous assessment of academic survey quality. Your task is to conduct a comprehensive evaluation using established scholarly standards and provide detailed justification for your assessment.
```

```
</topic>
<survey_content>
[SURVEY]
</survey_content>

<instruction>
You are provided with:
1. A research topic for context
2. An academic survey for evaluation

Your task is to assess the survey quality based on the specific criterion provided below. Apply rigorous academic standards and provide detailed justification for your assessment. Base your evaluation on specific evidence from the survey content, considering both strengths and areas for improvement.
</instruction>
```

<evaluation_criterion>
Criterion Description: [Criterion Description]

CRITICAL: Evaluation Standards
Your evaluation must follow a systematic approach:

- **Comprehensive Analysis**: Thoroughly examine the survey content against the specific criterion
- **Evidence-Based Scoring**: Base your score on specific observable strengths and weaknesses
- 3. **Detailed Justification**: Provide specific examples and reasoning for your score

```
864
865
           **Scoring Framework**:
           Score 1: [Score 1 Description]
867
           Score 2: [Score 2 Description]
           Score 3: [Score 3 Description]
868
           Score 4: [Score 4 Description]
869
           Score 5: [Score 5 Description]
870
871
           </evaluation_criterion>
872
           <output format>
873
           Provide your evaluation in the following structured format:
874
875
           **Rationale:**
876
           <Provide a comprehensive analysis of the survey's performance</pre>
877
              against the specific criterion. Include specific examples of
              strengths and weaknesses, with detailed justification for
878
              your assessment. Address how well the survey meets the
879
              criterion description and identify specific areas that align
880
              with or deviate from the scoring descriptions.>
881
882
           **Final Score: **
           <SCORE>X</SCORE>
883
           (Where X is the score from 1 to 5 based on your evaluation)
884
885
           Return your response in the following JSON format:
886
887
             "rationale": "Your detailed reasoning here",
             "score": X
888
889
           </output_format>
890
891
           Now conduct your comprehensive evaluation of the academic survey
892
              quality.
893
```

Coverage Criterion

894895896897

898 899

900

901

902 903

904

905

906

907

908

909

910

911

912

913

914915916917

Description: Coverage: Coverage assesses the extent to which the survey encapsulates all relevant aspects of the topic, ensuring comprehensive discussion on both central and peripheral topics.

Score 1: The survey has very limited coverage, only touching on a small portion of the topic and lacking discussion on key areas.

Score 2: The survey covers some parts of the topic but has noticeable omissions, with significant areas either underrepresented or missing.

Score 3: The survey is generally comprehensive in coverage but still misses a few key points that are not fully discussed.

Score 4: The survey covers most key areas of the topic comprehensively, with only very minor topics left out.

Score 5: The survey comprehensively covers all key and peripheral topics, providing detailed discussions and extensive information.

Structure Criterion

 Description: Structure: Structure evaluates the logical organization and coherence of sections and subsections, ensuring that they are logically connected.

Score 1: The survey lacks logic, with no clear connections between sections, making it difficult to understand the overall framework.

Score 2: The survey has weak logical flow with some content arranged in a disordered or unreasonable manner.

Score 3: The survey has a generally reasonable logical structure, with most content arranged orderly, though some links and transitions could be improved such as repeated subsections.

Score 4: The survey has good logical consistency, with content well arranged and natural transitions, only slightly rigid in a few parts.

Score 5: The survey is tightly structured and logically clear, with all sections and content arranged most reasonably, and transitions between adajecent sections smooth without redundancy.

Relevance Criterion

Description: Relevance: Relevance measures how well the content of the survey aligns with the research topic and maintain a clear focus.

Score 1: The content is outdated or unrelated to the field it purports to review, offering no alignment with the topic.

Score 2: The survey is somewhat on topic but with several digressions; the core subject is evident but not consistently adhered to.

Score 3: The survey is generally on topic, despite a few unrelated details.

Score 4: The survey is mostly on topic and focused; the narrative has a consistent relevance to the core subject with infrequent digressions.

Score 5: The survey is exceptionally focused and entirely on topic; the article is tightly centered on the subject, with every piece of information contributing to a comprehensive understanding of the topic.

Survey-Arena Review Prompt

```
# Paper 1:
Title: {title_1}
Figures: {figure_and_captions_1}
Content: {main_content_1}

# Paper 2:
Title: {title_2}
Figures: {figure_and_captions_2}
Content: {main_content_2}

You are provided with two survey papers on topic: {topic}.
```

```
972
973
          As the area chair for a top ML conference, you can only select
974
              one paper. Start with a brief meta-review/reasoning of the
975
              pros and cons for each paper (two sentences), focusing on:
976
           (1) insight and synthesis - moves beyond mere summarization to
977
              create new understanding and provides clear taxonomy;
978
           (2) thoroughness and accuracy - comprehensive coverage of
979
              literature with technical correctness;
980
           (3) structure and clarity - logical organization with compelling
              narrative:
981
           (4) scope and impact - well-defined scope with valuable future
982
              research directions;
983
           (5) presentation quality - professional polish, clear writing,
984
              and comprehensive evaluation of figures/tables presence and
985
              aesthetic quality.
986
          Be very critical and do not be biased by what the author claimed.
987
              Finally, provide your choice in a binary format.
988
989
          **Your Task:**
990
          1. Provide a detailed evaluation for Paper 1 using the above
991
              criteria.
          2. Provide a detailed evaluation for Paper 2 using the same
992
              criteria.
993
          3. Make a final decision by comparing the two papers and
994
              justifying your choice.
995
          STRICT OUTPUT INSTRUCTIONS:
996
           - You MUST return a single valid JSON object.
997
          - Output ONLY JSON. No explanations, no Markdown, no code fences,
998
              no additional text before or after the JSON.
999
           - Use exactly these keys and types:
1000
            - "paper_1_review": string
            - "paper_2_review": string
1001
            - "chosen_paper": "1" or "2"
1002
          - Do NOT include any additional keys or trailing commas. If
1003
              unsure, return empty strings for the review fields.
1004
1005
          Return JSON in exactly this shape:
1006
          "paper_1_review": "Your meta-review and reasoning for paper 1",
1007
          "paper_2_review": "Your meta-review and reasoning for paper 2",
1008
          "chosen_paper": "1 or 2"
1009
          }
1010
          End your output immediately after the closing.
1011
1012
1013
```

A.8 COMPARISON BETWEEN AUTOSURVEY AND ITERSURVEY.

Comprehensive Survey on the Alignment of Large Language Models (LLMs)

1 Introduction to Alignment in LLMs

1.1 Overview of Large Language Models

Large Language Models (LLMs) are a groundbreaking advancement in the field of artificial intelligence, especially in natural language processing (NLP). These models are primarily defined as deep learning systems designed to generate and comprehend human language by 18 er mice, 14 is amounts of testual data. At their core, LLMs utilize transformer architectures to precit study only in the virals in a sequence based on the correct provided by preceding words, allowing from min uncell singuage comprehension and generation.

The architecture of LLMs is predominantly; 1.5 d., 19 is transformer model into ad, ed by Vaswani et al. in 2017, which employs a self-attention, 15 in the mediants me enable (19 in 18 in

per to find or the finding or the minds of the second and adaptability. Advantage and adaptability of the minds of the second adaptability of the second ada

Despite their impressive competencies, LLMs face criticism and scrutiny. While adept at generating coherent and contextually appropriate outputs, these models often lack genuine understanding of the content they produce. This shortcoming can lead to the generation of plausible but incorrect information, risking concerns regarding the reliability and accountability of their applications [4]. Additionally, the vast and

diverse training datasets essential to LLMs present ethical challenges concerning bias, fairness, and transparency. The data utilized for training may harbor biases that modes can inadvertently learn and perpetuate, undersoring the imperative for alignment strategies that incorporate human values and ethical considerations. Addressing these challenges is crucial for ensuring that LLMs are not only efficient but also responsible in their deployment, partucularly in sensible and impactful sectors [5].

As the landscape of LLMs continues to evolve, ongoing research is dedicated to improving their efficiency, accuracy, and alignment with human values. There is a growing interest in talloring LLMs for domain-specific applications to enhance their performance in specialized tasks such as bioinformatics, legal analysis and scientific research, where unique terminologies and contexts necessitate precise adjustments to the

Overall, LLMs represent a transformative leap in Al, driven by significant in the set, all advancements, immense computational resources, and sophisticated learning and significant heart, models become increasingly integrated into every day life and various just 69 at 10 stocks, maintaining a focus on their displacement of the repulsibilities and inflators of the laws remains exemited. Continuous eventual control of the production of the repulsibilities and inflators of the size of the plan further research, and sophistics and inflators of the plan further research, and sophistics and effect (in 1) system (6).

1.2 Significance of LLN's in Modern Al The significant of the 2 language models (LLMs) in moor mix is profound, catalyzing transformative shifts across multiple investes applications that reals in the profound catalyzing transformative shifts across multiple investes applications that reals in 0 to 10 months of the control of th

cheating the season of the sea

In the business realm, LLMs significantly enhance operational efficiency by streamlining workflows and improving decision-making processes. An increasing number of organizations adopt LLMs for applicable like customer service automation, sentiment analysis, and content creation. These models analyze customations are control as control as a control inquiries and provide real-time responses, improving user experiences while alleviating the burden on human agents. Entiremente, they enable better data-driven exclosively by working valying extensive human agents. Entiremente, they enable better data-driven exclosively by with young provided to the control of the control o

Towards Ethical Alignment of Large Language Models with Human Values: A Comprehensive Survey

The alignment of Large Lu ps. see "socies (LLMs) with human values and preference has emerged e" a c tic." t as of focus within artificial intelligence, particularly seeken models file, "op" ancies a scores workons section that significantly impact human lives [136, 70]. "see "ring out LLM outputs align with ethical standards and use lives [136, 70]. "see "ring out LLM outputs align with ethical standards and use can lead to risks such as biased outputs and misinformation, which undertinine use trust in AI technologies [70, 30]. Therefore, alignment is integral to the responsible deployment of AI systems, necessitating a comprehensive exploration of current research."

Figure 1: Survey Structure Overview

Figure 4: LLM-generated survey comparison between AutoSurvey and IterSurvey.

datasets and generating actionable insights [12]. The ability of LLMs to produce sophisticated conten optimize marketing campaigns, and personalize user interaction businesses engage with customers and leverage Al technology. In cybersecurity, LLMs are emerging as powerful tools to strengthen security measures and address value abilities in digital infrastructures. Their adeptiness in performing natural language processing tasks at scale equips them to analyze potential threats and detect anomalous patterns in vast datasets. Research highlights their potential to automate threat detection and incident response, enabling security professionals to respond more effectively to emerging threats [13]. Additionally, LLMs can generate insights from historical stack data, empowering organizations to preemptively mitigate risks and bolster their cybersecurity frameworks. cybersecurity frameworks.

Beyond individual sectors, LLMs present opportunities for cross—dustry so Jaiop: that tackie complex challenges. For example, in supply chain management, IL-Ns (that ce predictive analytics by analyzing market trends and consumer behaviors, facilitating, nr. Liver inventory management and logistics operations [14]. The integration of LLMs acro's Vir' to a insustries not only promot is operational efficiency but also fosters collaboration and inc., at nr.

As LLMs continue to evolve, Vieri or, in anhanding accessibility of in on harvin bedomes increasingly significant. The _norm. The "a cyber be by granting, "info "that is a cyst to sophisticated At tools that were previously till doc's a "see cellited professions..." In a see that implications for underserved communities, when CLMs can help bit live from ferder upon in areas such as health education and legal assistance. By equipping users with jim use it in amadon, LLMs empower them to make informed decisions that improve the "vally or "fe, 15". decisions that this over on the control of the cont on these models, it is crucial for stakeholders to estrusin full times and transversal to the end to instance in the instance of U.* in mo, and it is underlable. Their lansies habite effects are evident across various sectors, with the cit act is revolutionize in afficient and the recognition of the cit act is revolutionize in afficient and the recognition of the cit act is revolutionized in afficient and the recognition of the cit act is revolutionized in a first and the recognition of the cit act is revolutionized in a first and the recognition of the cit act is recognitional practices, and in the cit act is recognitional practices, and in the cit act is recognitional practices. The cit is act is recognitional practices and the cit is act in the cit in the cit is act in the cit is act in the cit in the cit is act in the cit is act in the cit in the cit is act in the cit in the cit is act in the cit in the cit in the cit is act in the cit in the cit is act in the cit in the cit in the cit is act in the cit in the cit in the cit in the cit is act in the cit in th industry practitioners will be essential in realizing he will notential of LLMs while concurrently addressing the ethical challenges and societal in productions of their increasing integration into our daily lives. 1.3 Human Interaction and LLMs Large Language Models: "...s) h...re emerged as a transformative force in the landscape of human-computer interaction (HC), particularly in the realm of conversational Al. These models, capable of processing and generating human-like text, are fundamentally reshaping how users engage with technol across various domains. This subsection deleves into the intricacies of human interaction with LLMs, emphasizing their facilitation of conversational Al, implications for user experience, and the factors that enhances user acroament. Furthermore, effective alignment for LLMs requires an iterative process that embraces feedback from various stakeholders. Traditional training methods for Al systems often follow all linear approach, focusing primarily on model training and evaluation. However, the evolving nature of societal values and norms excessitates an adaptive approach to alignment that accommodates change over time. Stakeholders should be actively engaged in the ongoing assessment and refinement of Al systems to ensure they remain aligned with shifting human values. This perspective aligns with ne notion of "Bidirectional Human-Al Alignment," wherein both Al systems and users are in a constant state of adaptation to each other [68]. wherein both AI systems and users are in a constant state of adaptation to each other (68). Another significant sociotechnical challenge concerns the need for accountability and governance structures that can manage the complexities associated with AI deployment. As LUAs are increasingly integrated into decision-making processes across various domains—such as healthca's and cirri all justice—the ramifications of insistignment become more pronounced. Establishing, is usen in all similars for accountability and traceability in AI decision-making is essently in its required. Systems to be designed with clear standards of transparency and governarie (e.g. sey in the stakeholders understand how alignment is achieved and can discuss the associate et al., and circinsticate sizes related to misalignment and "obes" up. Latrust in Al systems (56).

The social context of AI deployment has servironments, such as duce in or law enforcement, necessitates balancing use "autono", with his potential for ("infraind") size (or locations and stakeholders understand are instrumental in navigating those it half in a properties of the contexts, misalignment and real to vicely-world repercussions, such as inforcing existing inequalities or favoring certain groups as others. Clear guide's erg noi for, eithical considerations and stakeholder engagement are instrumental in navigating those it half unlike the broader impact of their AI systems on societal values. This im a vest exceptating not only the immediate outcom's of A outputs but also how no societal values. This inv. Ves recognizing not only the immediate outcom, so fol a Jourst but also how these systems can shape and influence public perception and behavior of or t. ne. As Al capabilities continue to evolve, the potential for consequential impacts grevs. technology and society, researchers and practitioners mu (1) to provide in a nodologies for asset ing alignment that integrate sociotechnical dimensions and respon (to energy ing societal challenges (1)). the integrale solution and unintensities and report to Mary and solution training a societation and interest 3 Techniques for Alizmag LLMs 3.1 Reinforcem and Learning from Human Feedback (RLHF) Reinforcement Learning from Human Feedback (RLHF) represents a pivotal approach in aligning Large Language Models (LLMs) with human values and preferences. This framework enables models to learn desirable behaviors through direct interactions and feedback provided by human evaluators, rather than relying solely on traditional supervised training techniques. As LLMs are increasingly tasked with complex functions across deviews applications. RLHF has become essential for developing models that not only perform accurately but also align with user expectations in behavior and output.

1. Introduction

ethodologies. This survey aims to provide a thorough overview of the literature on LM alignment, addressing the theoretical foundations and practical methodologies

methodologies. This survey aims to provide a thorough overview of the literature on LLM alignment, addressing the theoretical foundations and practical methodologies that have emerged in this rapidly evolving field [67]. Various alignment strategies, including Reinforcement Learning from Human Feedback (RLHF) [52] and Direct Preference Opinization (DPO) [37], are explored, alongside emerging frameworks for personalized and cultural alignment [59]. The structure of this survey is summarized in Figure 1, which outlines the key components and sections we will cover. By synthesizing insights from recent literature, this survey seeks to fill critical gaps in understanding how effective alignment can be achieved across diverse contexts set used mographics. The historical context of alignment research reveals a pro, rese on 1 on simplistic, rule-based systems to more sophisticated methods that onsis it in con plexities of human values [98]. Early alignment strategies rown it by course, on ensiring that Al systems adhered to predefined specificaties, it is one obtaining the rich diversity of human preferences. Contemporry aip to the "verage advanced techniques like RLHF, which unlikes human—it as it is a reveal alignment, which were the second of the properties of the content of the

which often reversible successions to the alignment of LLMs with human preferences is an undertaking requiring in-depth engagement with the underlying theories, methodologies, and ethical considerations involved [33]. By fostering a deeper understanding of alignment as a multifaceted challenge, the research community can work towards creating AI systems that are not only technically proficient but also socially responsible and reflective of the values of the communities they serve [92]. The exploration of personalized alignment strategies is particularly vital as user interactions

3. Theoretical Foundations of Alignment

3 Theoretical Foundations of Alignment



Figure 2: Conceptual Framework of Inver and Deer lignment in AI Systems

Figure 2: Conceptual Framework of In'er and Dier lignment in Al Systems

The alignment of Large Long age of the (LLMs) with human vilues is a pivotal area of research that un'er product the internations between it, exciton will age, at the contractions between it is a product of the contractions between in a contraction of the contraction between in a contraction of the contraction of t

Figure 5: LLM-generated survey comparison between AutoSurvey and IterSurvey.

| | (a) AutoSurvey | (b) IterSurvey |
|------|--|--|
| 1183 | satisfaction. | principles of consumer dust. |
| 1182 | types of questions posed, (vivers tailored responses, ultimately enhancing user engagement and | present, ensuring that the benefits of AI are realized without compromising ethical principles or consumer trust. |
| 1181 | produce inform a do and answer queries of the sets of private science, technology, arts, and personal advice. For example, studies indicate the Clato Ticar generate human-like responses in educational contexts, assisting student is "thi "qui es pross multiple disciplines [6]. By adapting its knowledge to the | stakeholders remain vigilan, in addressing the ethical and regulatory challenges they |
| 1180 | by extensive ti initing a cer ompassing various opics. I list on prehensive training enables the model to produce inform, tipp, and answer queries on prefets so immig science, technology, arts, and personal | guidance |
| 1179 | appropriate, demonstrating a higher (e) \(\) \(| Legal F1 Score 3.49 (av.) Get via as relevant Incomplete reasoning 60.86 with contact paths |
| 1178 | relative to one another, leading to a deeper undersum of gorcontext compared to previous models. As a result, ChatGPT generates responses that result, chatGPT generates responses that results are contextually correct but a so contextually | derstanding (20% in Anatomy) |
| 1177 | The transformer architecture is particularly advantageous for to gene, tion thanks to its self-attention mechanism. This feature enables the model to weigh the to tran be of different words in a sentence | Healthcare USMLE 58.2 log ca Potential utility at tu- Lacks nuanced un- |
| 1176 | effectiveness of many natural language processing (NLP) tasks [57]. | Education MCQs 56.9% High acc racy in cod- Low accuracy in 10°, (°° % on Leet- concepts (33.4% in DBMS) |
| 1175 | agent across numerous applications. This conversational fluency is largely attributed to its underlying architecture, which employs the transformer model—an advancement that has substantially enhanced the | Domain Metric Accuracy Streng bs Weaknesses |
| 1174 | generating coherent and contextually relevant responses across a variety of domains. A key strength of ChatGPT lies in its ability to engage in human-like conversation, effectively positioning it as a conversational | Examination, F1 = F1 Score |
| 1173 | ChatGPT, developed by OpenAI, has garnered significant attention due to its impressive capabilities in | Table 1: Performance Evaluation of ChatGPT across Different Domains. Abbreviations MCQs = Multiple Choice Questions, USMLE = United States Medical Licensing |
| 1172 | 3.2 Strengths of ChatGPT | Table 1. Desferences Evaluation of ChatCNT |
| 1171 | | 5. Performance Evaluation of ChatGPT |
| 1170 | | |
| | | |
| 1169 | | |
| 1168 | | |
| 1167 | | |
| 1166 | Barr to generality our (e) company | |
| 1165 | Figure 6: LLM-generated survey compari | son between AutoSurvev and IterSurvev. |
| 1164 | | (b) Hersurvey |
| 1163 | | (b) IterSurvey |
| 1162 | | |
| 1161 | | 28 |
| 1160 | (a) AutoSurvey | Amodei, Tom Brown, Jack Clark, Sam McCandlish, Chris Olah, Ben Mann, |
| 1159 | | Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom Brown, Jack Clark, Sam McCandlish, Chris Olah, Ben Mann |
| 1158 | [21] CloChat: Understanding How People Customize, Interact, and Experience Personas in Large Language Models | DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Sauraw Kadavath, Jackson Kernion, Tom Conerly, Sheer El-Showk, Nales Ellsen, Zea, Hatfold Deade, Deany Hermonder, Teiter Hump. Sort |
| 1157 | an Al Interapps: [20] Attacks, Defenses and Evaluations for LLM Conversation Safety: A Survey | [10] Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova |
| 1156 | [19] Script-Based Dialog Policy Planning for LLM-Powered Conversational Agents: A Basic Architecture for an "Al Therapist" | [9] Xuechur, i P. A. elina Wang, Ilia Sucholutsky, and Thomas L. Griffiths. Measuring mpl sit b as in explicitly unbiased large language models, 2024. |
| 1155 | Health [18] Exploring Autonom as A e. ** Iro-grane Lens of Large Language Models: A Review | Krueger. Foundatio al in large in assuring alignment and safety of large language models, 202. |
| 1154 | [17] Leveraging Large Language Models for Patients gagerner of the Favor of versational Al in Digital | Bengio, Danqi Chen, Philip H. S. Torr, Sar uer Albanie, Tegan Maharaj, Jakob Foerster, — an Tramer, He He, Ato — Kasirzadeh, Yejin Choi, and David |
| 1153 | [15] Using large anguige and a to place the learning by the learning of the learning by the learning of the le | Recchia, Giulio Co, si, Ala. Chan, Markus An et jung, Lilian Edwards, Alek- sandar etrov. Const an Schroeder e Witt, Sun et Ramesh Motwan, Yoshua |
| 1152 | Efficiency, and Innovation | Hernandez-Orallo, Lewis H. in. Jo. d, Eric Bigelow, Ale a dar Dan, Lauro Langosco, Tomasz Korbald H. io. Zb. ng, Ruiqi Zhong (2011) hangeartaigh, Gabriel |
| | [13] A Survey of Large Language Models in Cybersecurity [14] The Potential of Large Language Models in Supply Chain Management: [NOW] pring (Pecision-Making, | Peter Hase, Ekdeep Singh Lubana, E. k., en, er, orephen Casper, Ol ver Sourbut, Benjamin L. Edelman, Zhaowei Cha g, Mario Günther, Arton Kirinek, Jose |
| 1150 | Future Directions, and Strategic Recommendations | [8] Usman Anwar, Abulhair Saparov, Javier va. 4, Paniel Paleka, Miles Turpin, |
| 1150 | [11] Automatically Generating CS Learning Materials with Large Language Models [12] Harnessing the Potential of Large Language Models in Modern Marketing Management: Applications, | Goodman. Star-gate: Teaching language models to ask clarifying questions, 2024. |
| 1149 | [10] LLMs in Education: Novel Perspectives, Challenges, and Opportunities | [7] Chinmaya And Ku. Jan-Philipp Fränken, Tobias Gerstenberg, and Noah D. |
| 1148 | [9] Harnessing Large Language Models for Mental Health: Opportunities, Challenges, and Ethical Considerations | [6] Remauch Amini, Sanaz Salti Jorouzi P. al Hitzler, and Reza Amini. Towards complex ontology oli among substantial angulage models, 2024. |
| 1147 | Software Testing | Ston, tair he fair epuon: A framework for the interpretations of (gorith in the mess, 2021. |
| 1146 | [8] Are We Testing or Being Tested? Exploring the Practical Applications of Large Language Models in | [5] Georg Ahnert, Van S. and V. Florian Lemmerich, Clay dia Way er, and Markus |
| 1145 | [6] Advancing bioinform A ics with latite any lager mights to a poly ints public itions and perspectives [7] Large Language Mr Jels \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ | Olivier Pietquin, Ahmet Üstün, and var die kein Back to basics: Revisiting reinforce style optimization or ac a in, from human feedback in Ilms, 2024. |
| 1144 | [5] Challenges and Applications of Large Language Models | comprehensive survey and taxonomy, 2025. [4] Arash Ahmadian, Chris Cremer, Matthias 5. 16 Ma zieh Fadaee, Julia Kreutzer, |
| 1143 | [4] Large language mix rs 70 Led In the potentials and pitfalls | [3] Saleh Afzoon, Zahra Jahanandish, Phuong Thao Huynh, Amin Peheshti, and Usman Naseem. Modeling and optimizing user preferences in a copilots: A |
| 1142 | [3] Survey of different Large Language to ark-rch e tu and nd Ben marks, and Challenges | desirability in user experience, 2016. |
| 1141 | [1] A Survey of GPT-3 Family Large Language Models Including ChatGPT and GPT-4 [2] A Comprehensive Overview of Large Language Models | Aligning global and local preferences to reduce harm, 2024. [2] Sisira Adikari, Craig McDonald, and John Campbell. Quantitative analysis of |
| 1140 | References | Aakanksha, Arash Ahmadian, Beyza Ermis, Seraphina Goldfarb-Tarrant, Julia Kreutzer, Marzieh Fadaee, and Sara Hooker. The multilingual alignment prism: |
| 1139 | | References |
| 1138 | | Defenences |
| 1137 | | REFERENCES |
| 1136 | | |
| 1135 | | |
| 1104 | | |
| 1134 | | |

Figure 7: LLM-generated survey comparison between AutoSurvey and IterSurvey.