



# AutoSurvey: Large Language Models Can Automatically Write Surveys

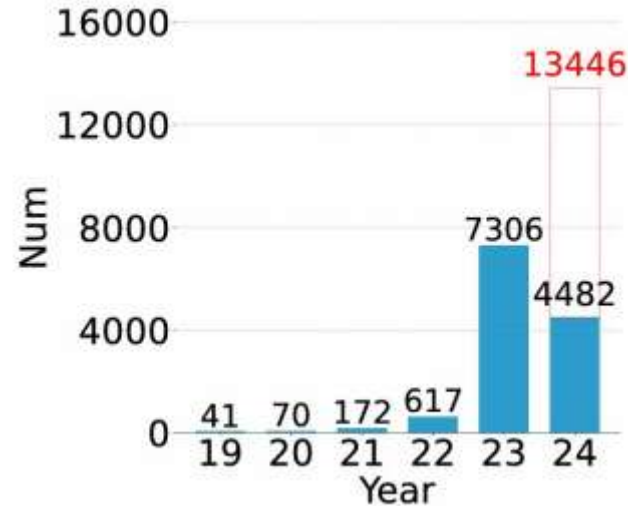
---

Yidong Wang<sup>1,2\*</sup>, Qi Guo<sup>2,3\*</sup>,  
Wenjin Yao<sup>2</sup>, Hongbo Zhang<sup>1</sup>, Xin Zhang<sup>4</sup>, Zhen Wu<sup>3</sup>, Meishan Zhang<sup>4</sup>,  
Xinyu Dai<sup>3</sup>, Min Zhang<sup>4</sup>, Qingsong Wen<sup>5</sup>, Wei Ye<sup>2†</sup>, Shikun Zhang<sup>2†</sup>, Yue Zhang<sup>1†</sup>

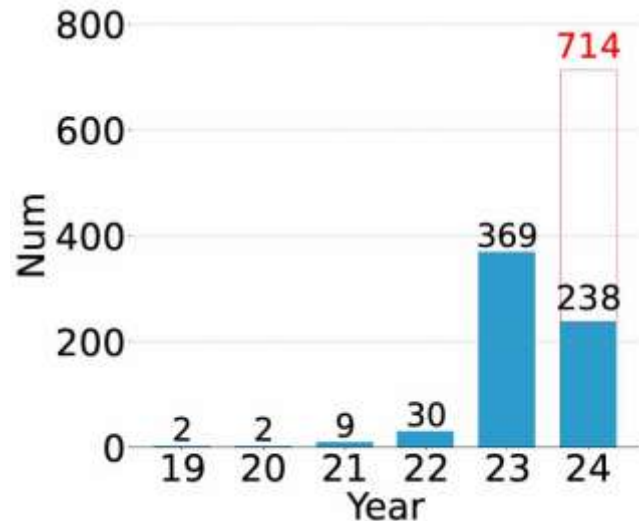
<sup>1</sup>Westlake University, <sup>2</sup>Peking University,  
<sup>3</sup>Nanjing University, <sup>4</sup>Harbin Institute of Technology, Shenzhen, <sup>5</sup>Squirrel AI

*NeurIPS 2024*

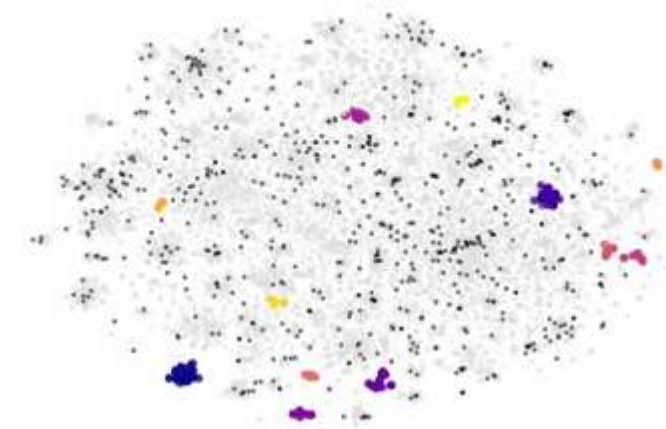
# Background



(a) Paper about LLM.



(b) Survey about LLM.



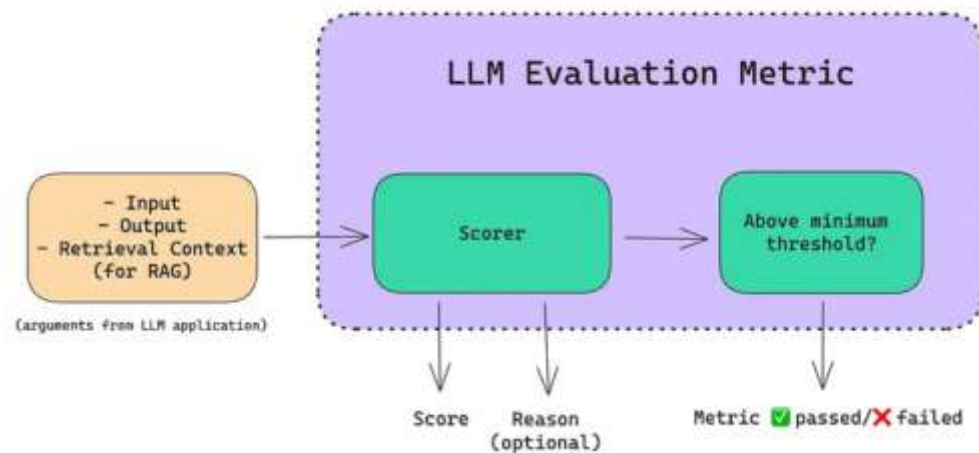
- Cluster 0: Emotion Recognition
- Cluster 1: Embedding Techniques
- Cluster 2: Event Extraction
- Cluster 3: Quantum Theory
- Cluster 4: Data Visualization
- Cluster 5: Neural Scaling Laws
- Cluster 6: Mixture-of-Experts (MoE)
- Cluster 7: Ontology Enrichment
- Cluster 8: Humor Detection
- Cluster 9: Empathetic Response Generation
- Cluster 10: Attribute Value Extraction
- Cluster 11: Anomaly Detection
- Regular Papers
- Surveys

(c) A t-SNE visualization of LLM-related surveys and papers shows clusters of research papers that currently lack comprehensive survey coverage.

# Background

模型 <sup>(1)</sup>		deepseek-chat	deepseek-reasoner
上下文长度		64K	64K
最大思维链长度 <sup>(2)</sup>		-	32K
最大输出长度 <sup>(3)</sup>		8K	8K
标准时段价格 (北京时间 08:30-00:30)	百万tokens输入 (缓存命中) <sup>(4)</sup>	0.5元	1元
	百万tokens输入 (缓存未命中)	2元	4元
	百万tokens输出 <sup>(5)</sup>	8元	16元
	百万tokens输入 (缓存命中)	0.25元 (5折)	0.25元 (2.5折)

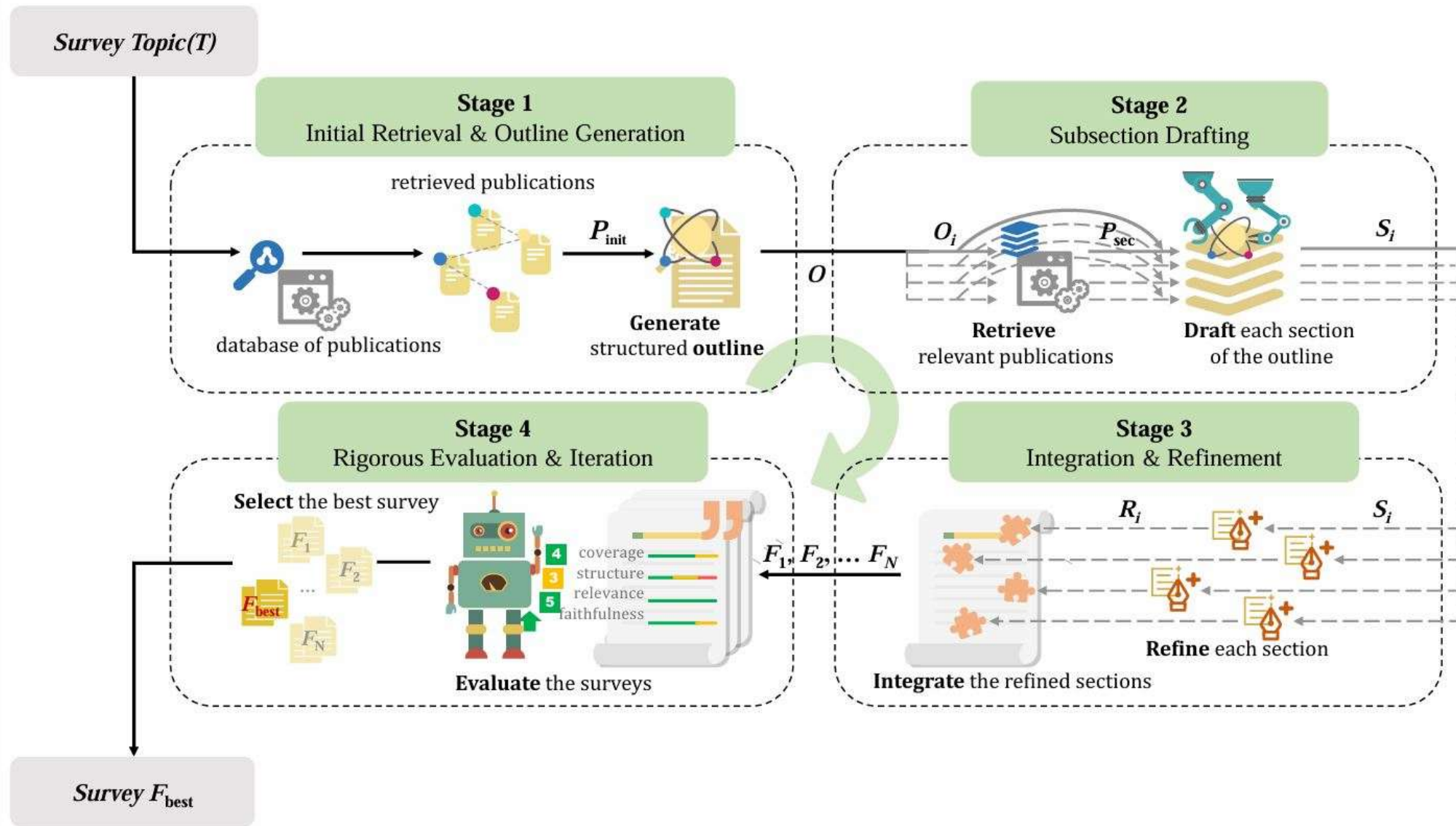
(1) Context window limitations



(3) The lack of evaluation benchmark



(2) Parametric knowledge constraints



**Survey Creation Speed**  $Time = T_r + T_w + T_e = T_r + \frac{L}{E \times M} + \frac{1}{2} \times \frac{L}{E \times M}$

**Citation Quality**  $Recall = \frac{\sum_{i=1}^{|C|} h(c_i, Ref_i)}{|C|}$

$$Precision = \frac{\sum_{i=1}^{|C|} \sum_{k=1}^{|Ref_i|} h(c_i, Ref_i) \cap g(c_i, r_{i_k})}{\sum_{i=1}^{|C|} |Ref_i|}$$

Where  $C = \{c_1, c_2, \dots\}$ , and  $Ref_i = \{r_{i_1}, r_{i_2}, \dots\}$

## Content Quality

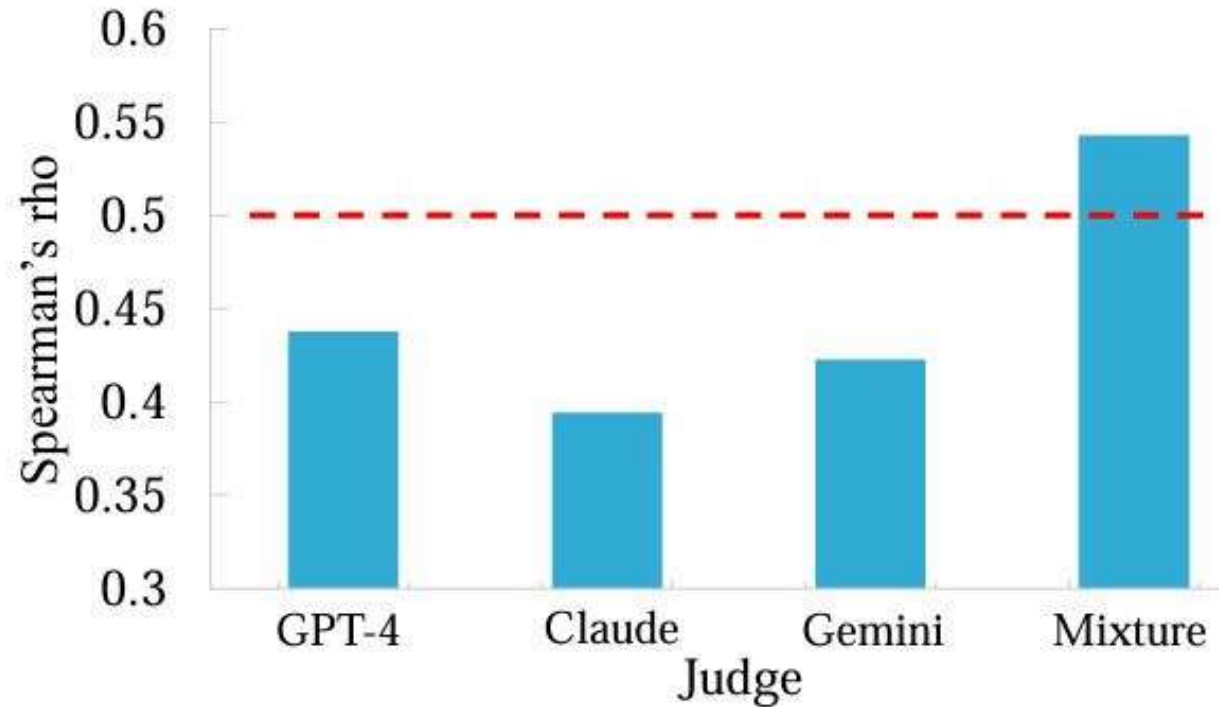
**Coverage:** Assesses the extent to which the survey encapsulates all aspects of the topic.

**Structure:** Evaluates the logical organization and coherence of each section.

**Relevance:** Measures how well the content aligns with the research topic.

Table 2: Results of naive RAG-based LLM generation, Human writing and AutoSurvey. Note that AutoSurvey and naive RAG-based LLM generation both use Claude-haiku as the writer. **Note that human writing surveys used for evaluation are excluded during the retrieval process.**

Survey Length (#tokens)	Methods	Speed	Citation quality		Coverage	Content Quality		Avg.
			Recall	Precision		Structure	Relevance	
8k	Human writing	0.16	80.00	87.50	4.50	4.16	5.00	4.52
	Naive RAG-based LLM generation	79.67	78.14 $\pm$ 5.23	71.92 $\pm$ 6.83	4.40 $\pm$ 0.48	3.86 $\pm$ 0.71	4.86 $\pm$ 0.33	4.33
	AutoSurvey	107.00	82.48 $\pm$ 2.77	77.42 $\pm$ 3.28	4.60 $\pm$ 0.48	4.46 $\pm$ 0.49	4.8 $\pm$ 0.39	4.61
16k	Human writing	0.14	88.52	79.63	4.66	4.38	5.00	4.66
	Naive RAG-based LLM generation	43.41	71.48 $\pm$ 12.50	65.31 $\pm$ 15.36	4.46 $\pm$ 0.49	3.66 $\pm$ 0.69	4.73 $\pm$ 0.44	4.23
	AutoSurvey	95.51	81.34 $\pm$ 3.65	76.94 $\pm$ 1.93	4.66 $\pm$ 0.47	4.33 $\pm$ 0.59	4.86 $\pm$ 0.33	4.60
32k	Human writing	0.10	88.57	77.14	4.66	4.50	5.00	4.71
	Naive RAG-based LLM generation	22.64	79.88 $\pm$ 4.35	65.03 $\pm$ 8.39	4.41 $\pm$ 0.64	3.75 $\pm$ 0.72	4.66 $\pm$ 0.47	4.23
	AutoSurvey	91.46	83.14 $\pm$ 2.44	78.04 $\pm$ 3.14	4.73 $\pm$ 0.44	4.26 $\pm$ 0.69	4.8 $\pm$ 0.54	4.58
64k	Human writing	0.07	86.33	77.78	5.00	4.66	5.00	4.88
	Naive RAG-based LLM generation	12.56	68.79 $\pm$ 11.00	61.97 $\pm$ 13.45	4.4 $\pm$ 0.61	3.66 $\pm$ 0.47	4.66 $\pm$ 0.47	4.19
	AutoSurvey	73.59	82.25 $\pm$ 3.64	77.41 $\pm$ 3.84	4.73 $\pm$ 0.44	4.33 $\pm$ 0.47	4.86 $\pm$ 0.33	4.62



Spearman's rho shows how well LLM rankings align with human experts, with values above 0.3 indicating positive correlation and above 0.5 indicating strong positive correlation.

Table 3: Ablation study results for AutoSurvey with different components removed.

Methods	Citation Quality		Coverage	Content Quality		Avg.
	Recall	Precision		Structure	Relevance	
AutoSurvey	83.48 $\pm$ 5.05	77.15 $\pm$ 6.05	4.7 $\pm$ 0.45	4.16 $\pm$ 0.73	4.93 $\pm$ 0.30	4.57
AutoSurvey w/o retrieve	60.11 $\pm$ 6.42	51.65 $\pm$ 6.33	4.51 $\pm$ 0.49	4.01 $\pm$ 0.74	4.88 $\pm$ 0.32	4.44
AutoSurvey w/o reflection	83.23 $\pm$ 3.82	76.36 $\pm$ 4.08	4.76 $\pm$ 0.42	4.13 $\pm$ 0.76	4.88 $\pm$ 0.32	4.56

Table 4: Performance of AutoSurvey with different base LLM writers.

Base LLM writer	Citation Quality		Coverage	Content Quality		Avg.
	Recall	Precision		Structure	Relevance	
GPT-4	80.25 $\pm$ 4.19	78.83 $\pm$ 7.00	4.8 $\pm$ 0.54	4.46 $\pm$ 0.49	4.86 $\pm$ 0.33	4.70
Claude-haiku	82.45 $\pm$ 2.77	76.31 $\pm$ 2.18	4.66 $\pm$ 0.47	4.26 $\pm$ 0.67	4.86 $\pm$ 0.33	4.58
Gemini-1.5-pro	78.13 $\pm$ 2.39	71.24 $\pm$ 3.28	4.86 $\pm$ 0.33	4.33 $\pm$ 0.78	4.93 $\pm$ 0.25	4.69
Human	85.86	80.51	4.71	4.43	5	4.70

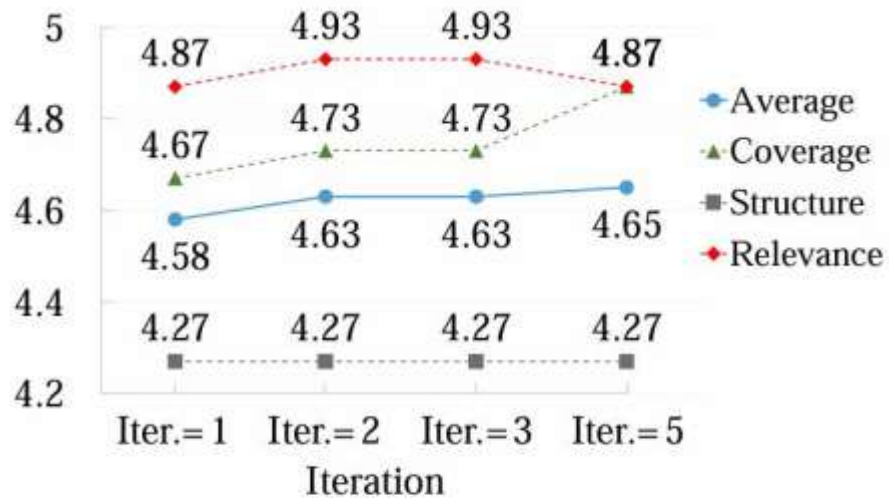
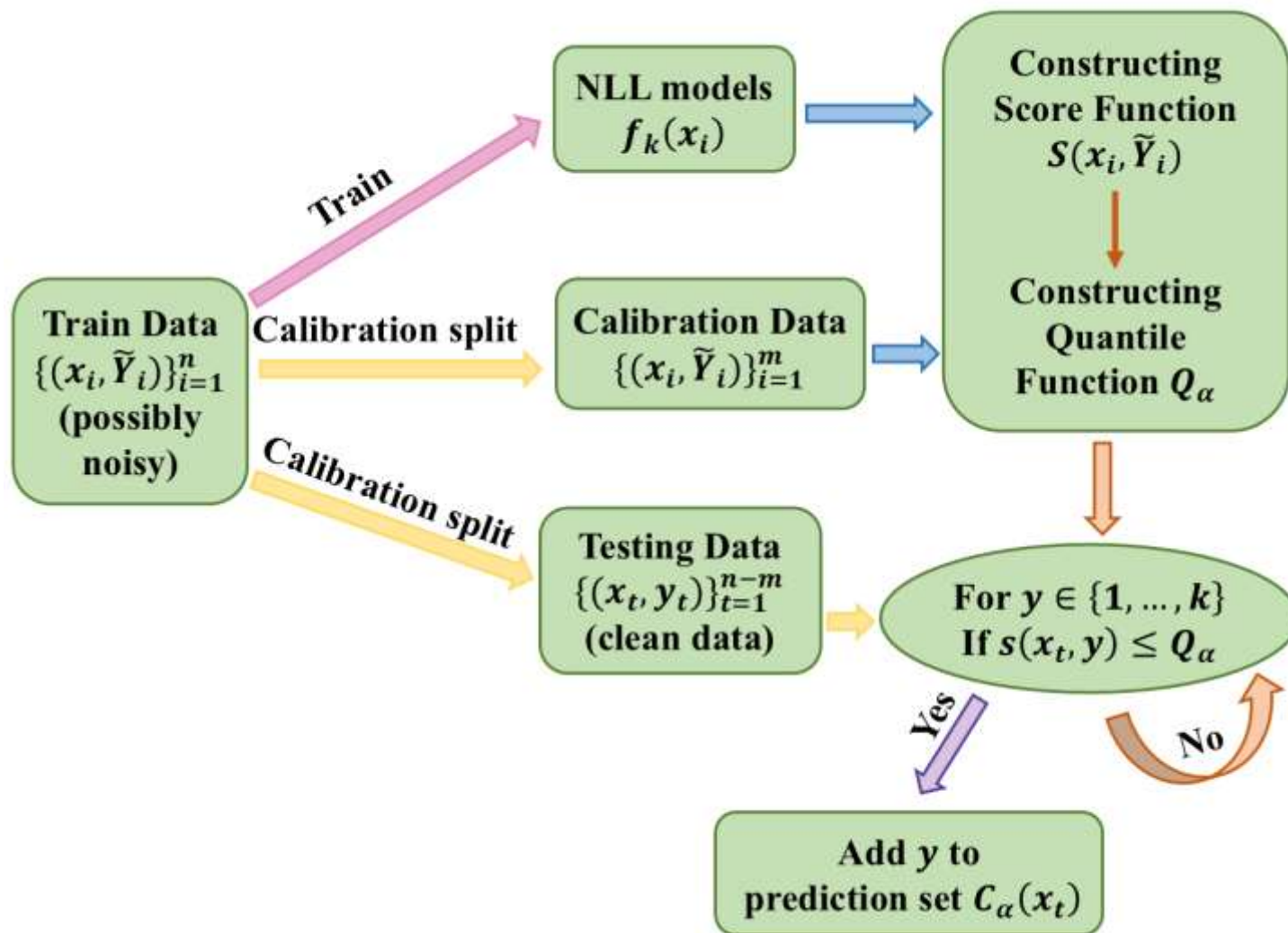


Figure 4: Impact of Iteration on Auto-Survey Performance.

Table 5: Performances given different references.

Methods	Accuracy
Direct	58.40 $\pm$ 4.96
Naive RAG-based LLMs	65.20 $\pm$ 8.06
Upper-bound	73.60 $\pm$ 3.44
AutoSurvey	67.60 $\pm$ 4.96



**Thanks**