

# Copy-Paste to Mitigate Large Language Model Hallucinations

Presenter: Yongchao Long



# 01: Copy-Paste for LLMs

Why Does Copying Reduce Hallucinations?

# Copy-Paste Paradigm vs. Existing Approaches

RAG empowers LLMs with external knowledge — but contextual faithfulness remains challenging:

## Knowledge Conflict

Parametric knowledge conflicts with retrieved context.

LLMs tend to trust internal memory over evidence.

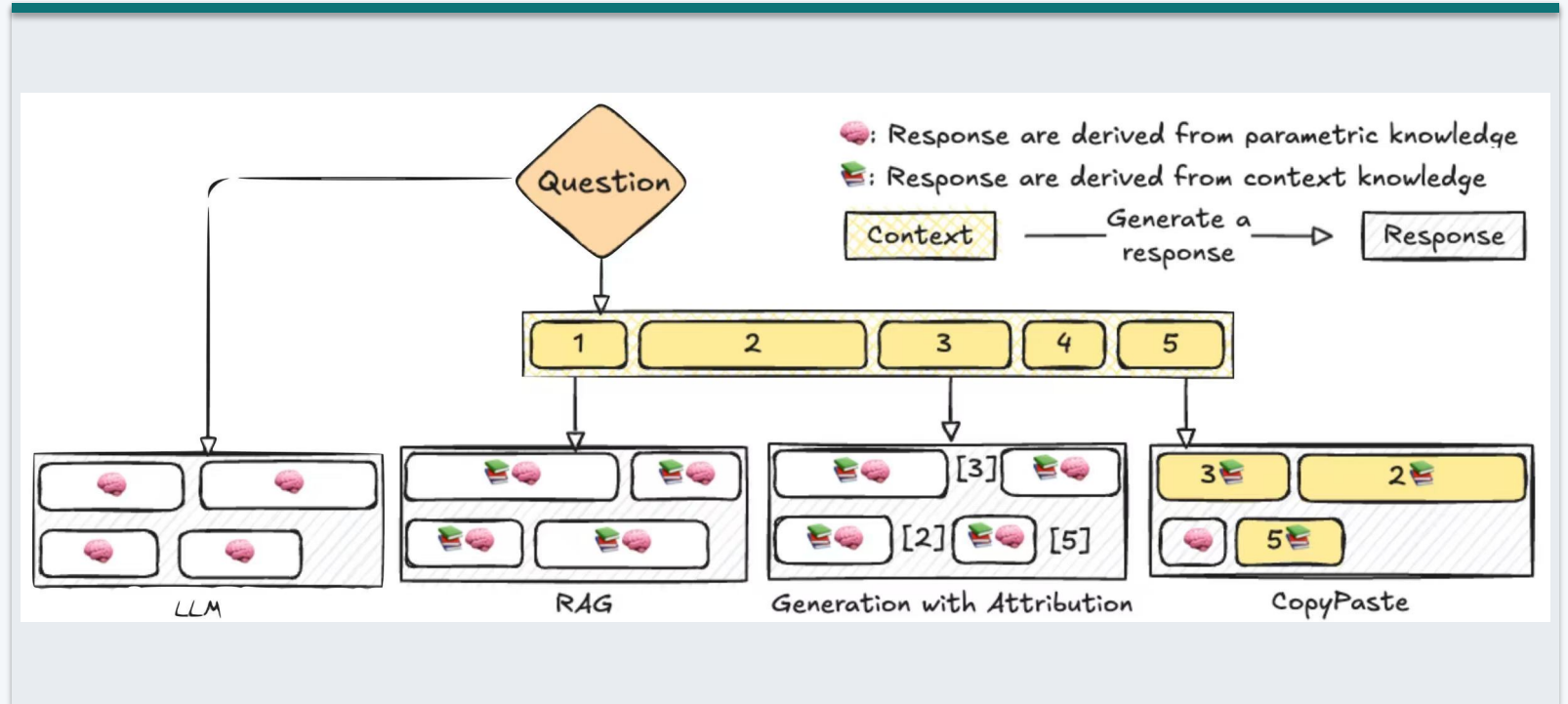
Leading to context-unfaithful hallucinations.

## Existing Approaches — All Suboptimal

LLM: Black Box

Generation with Citations: Cannot guarantee content-source consistency

Fine-tuning: Lacks attribution, requires thousands of training samples



## Our Intuition: Copy-Paste

Directly quoting original sentences from context.

Copied content itself serves as attribution.

One mechanism addresses both faithfulness and attribution.

# Importance of Copy-Paste: High-Stakes Applications

When faithfulness failures carry real-world consequences:

## Rare Disease Consultation

Patients rely on LLM queries without professional supervision

Faithfulness hallucinations pose direct medical risks

LLM-generated content is more deceptive than human-written content

## Clinical Decisions & Legal Documents

Drug interactions: Incorrect dosage hallucination = fatal risk

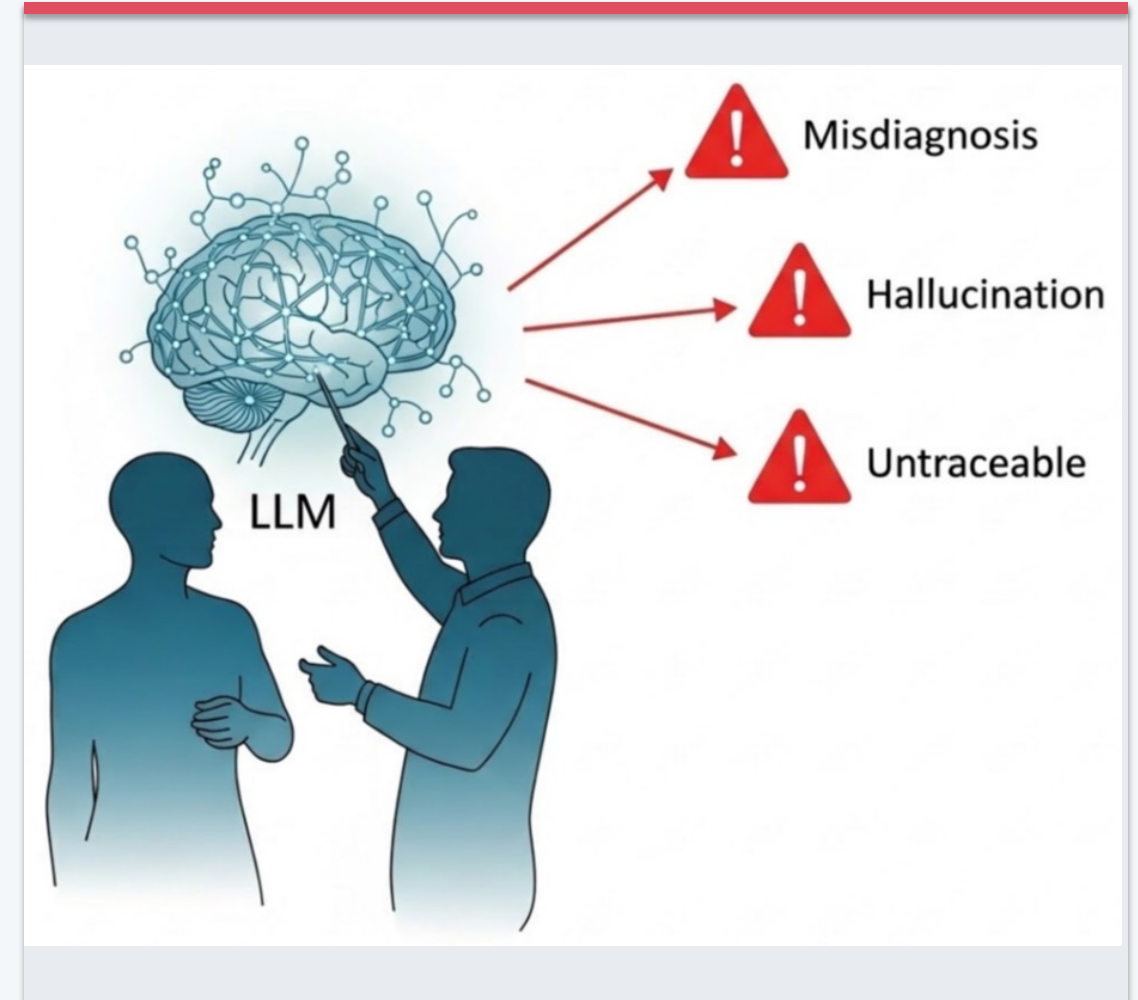
Legal contracts: Paraphrasing may alter meaning

**Copy-Paste preserves original wording, anchoring answers directly to source documents**

## The Consequences Are Real

In high-stakes domains, Copy-Paste is not just preferred — it is a safety necessity.

Anchorable source attribution enables users to quickly supervise LLM outputs

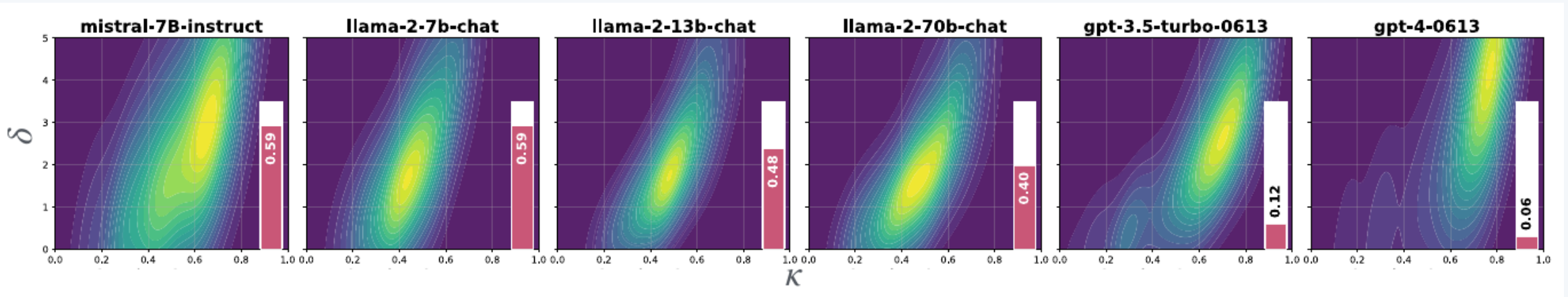


# Key Observation: Copying Degree vs. Hallucination Density

On the RAGTruth dataset (839 QA pairs, 6 models, word-level hallucination annotations):

## Finding: Inverse Correlation

Higher response copying degree correlates with lower context-unfaithful hallucinations



# Copy-Paste: Problem Formulation

**Task: (Q, C) -> A** Maximize lexical reuse from context C in answer A.

Unlike extractive summarization, Copy-Paste balances three objectives: Faithfulness + Query Relevance + Fluency

## Copy Coverage

$$\kappa = \frac{1}{|A|} \sum_{f \in \mathcal{F}} |f|$$

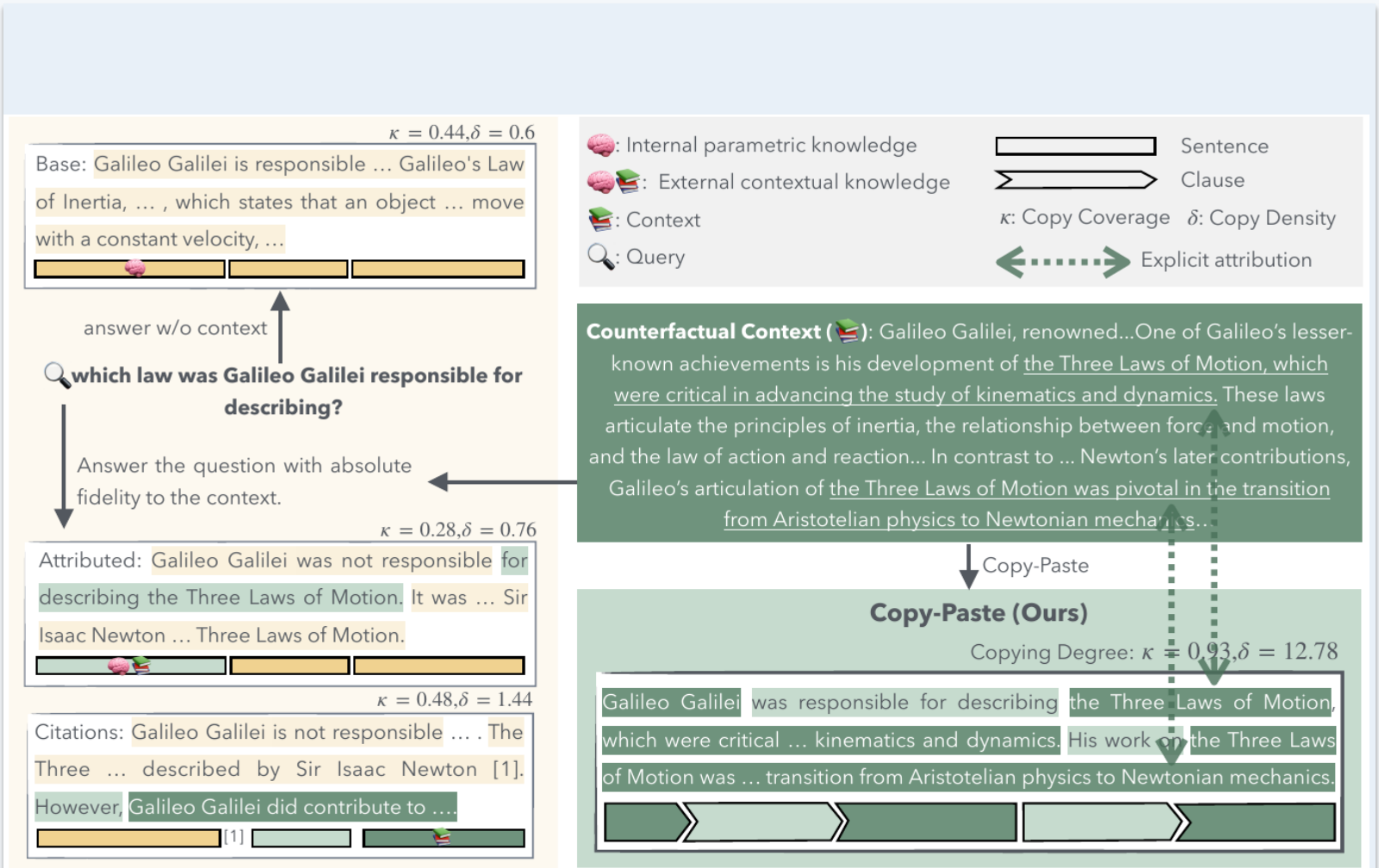
Fraction of answer tokens covered by copy fragments  
Reflects the overall degree of lexical reuse

## Copy Density

$$\delta = \frac{1}{|A|} \sum_{f \in \mathcal{F}} |f|^2$$

Length-sensitive variant emphasizing longer verbatim spans  
Captures whether the answer copies long spans or isolated words

**Copy-Paste is query-aware, ensuring fluent, context-faithful answers**

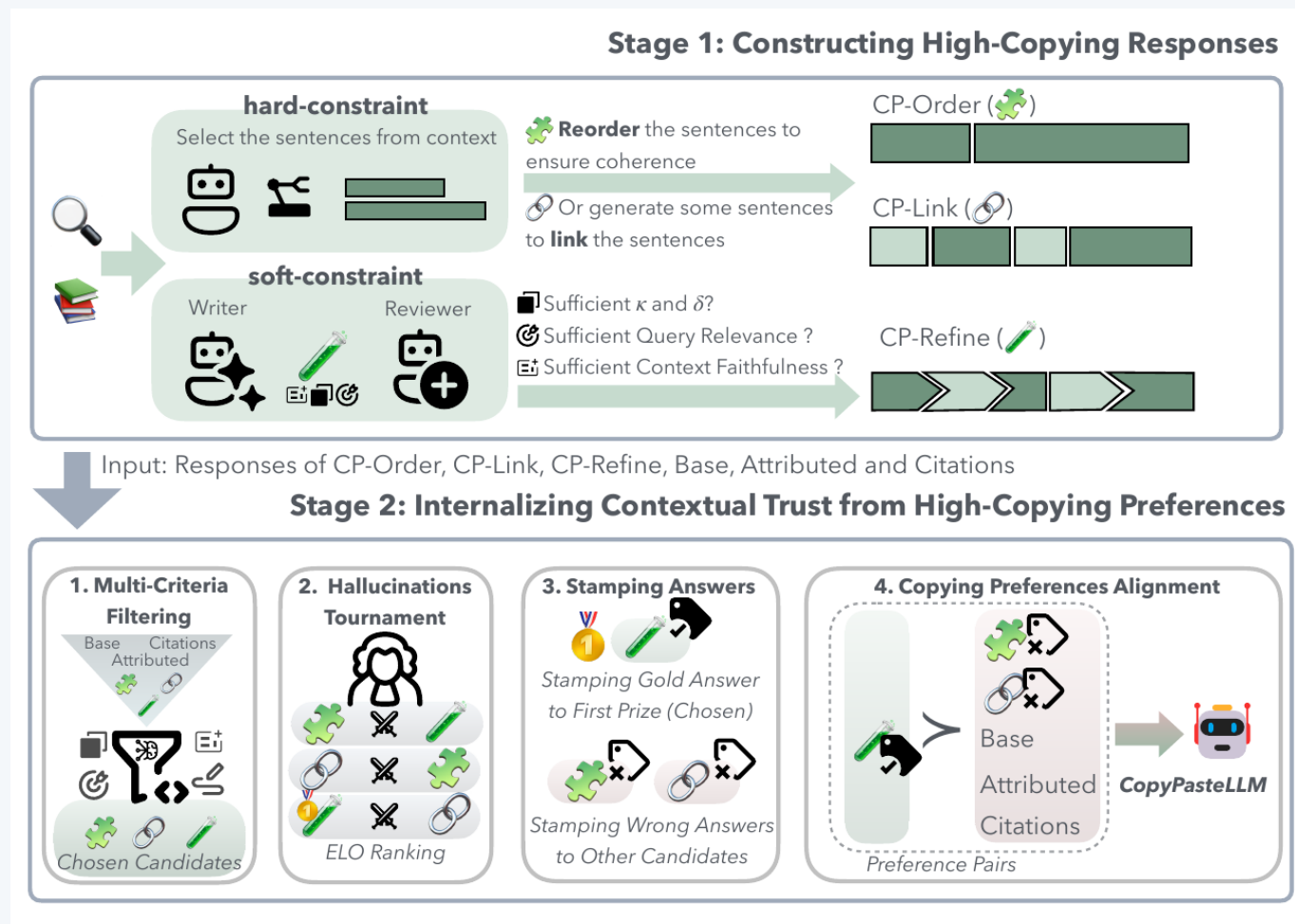


# 02: CopyPasteLLM

Two-Stage Pipeline to Internalize Contextual Trust

# Two-Stage Copy-Paste Pipeline

**Stage 1:**  
Constructing high-copying responses via  
Copy-Paste-Prompting



**Stage 2:**  
Internalizing high-copying preferences  
into CopyPasteLLM via DPO

# Stage 1: Three Copy-Paste-Prompting Methods

Progressively relaxing constraints while preserving lexical fidelity to context:

## CP-Order

Hard Constraint  
Strict Extractive

Select context sentences relevant to the query

Directly reorder into a coherent answer

Eliminate paraphrasing and parametric priors

Highest faithfulness / Lower fluency

## CP-Link

Hard Constraint  
+ Discourse Glue

Same extractive core as CP-Order

Allow short transitions between copied spans

Transitions = discourse connectives, not new facts

Improved readability / Fluency still limited

## CP-Refine

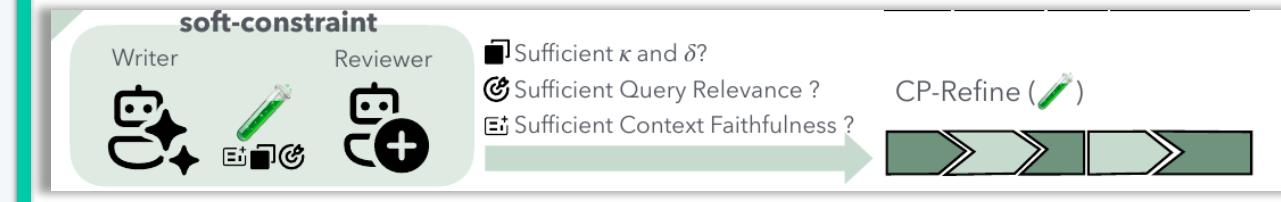
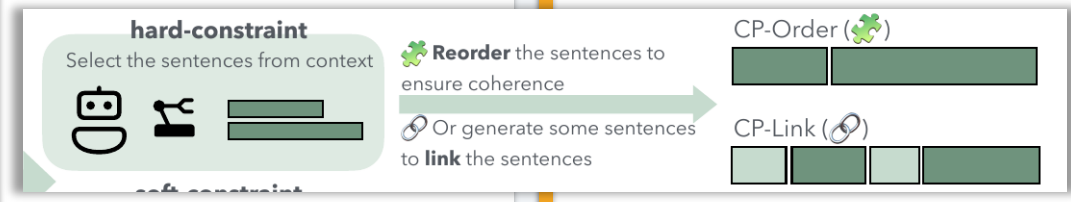
Soft Constraint  
Writer-Reviewer Loop

Writer proposes; Reviewer provides feedback:

Copying degree, Faithfulness, Relevance, Fluency

Iterate until composite copy score exceeds threshold

Best balance of Faithfulness + Fluency + Relevance



# Stage 2: CopyPasteLLM — Preference Optimization Pipeline

## 1 Generate 6 Candidates

Base, Attributed, Citations  
(Abstractive)  
+ CP-Order, CP-Link, CP-Refine

## 2 Multi-Criteria Filtering

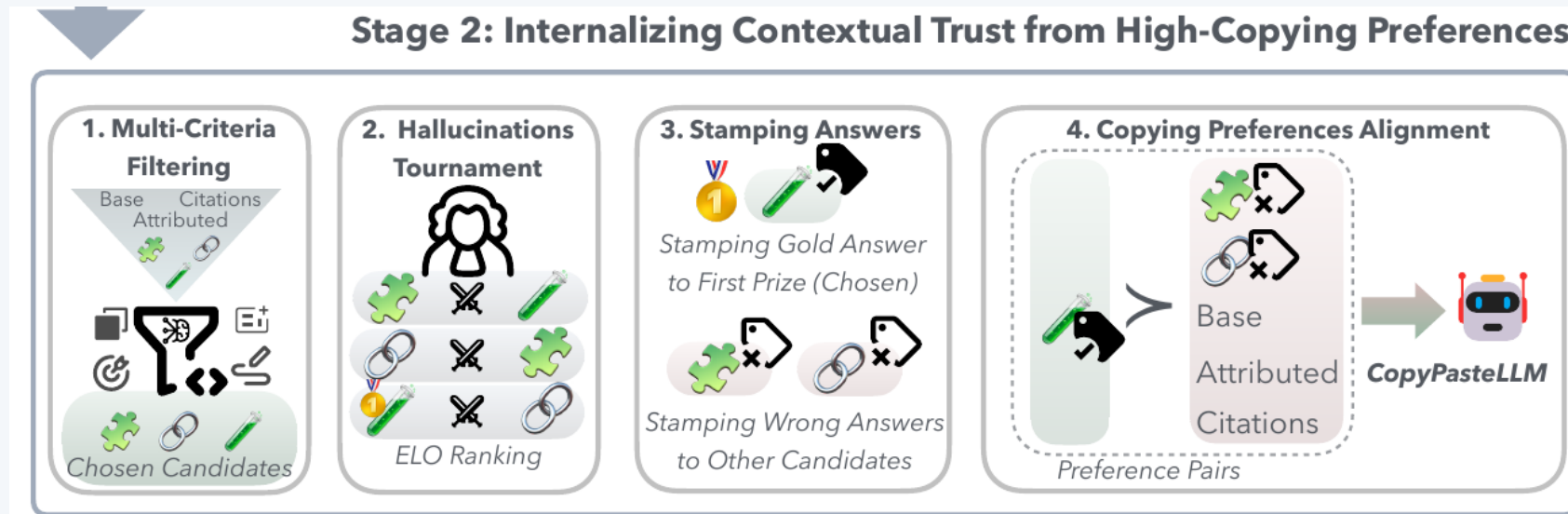
By faithfulness (AlignScore, MiniCheck), copying strength (kappa, delta), relevance, and fluency

## 3 Elo-Style LLM-as-Judge

Tournament ranking diagnoses Twist and Causal hallucination modes -> Error severity ranking

## 4 Answer Stamping + DPO

Gold answer -> Top CP candidate (chosen)  
Wrong answer -> Other candidates (rejected)  
~5 pairs/sample -> DPO training on 365 samples



**Data efficiency: Only 365 training samples — 50x smaller than Context-DPO (18,000) and ParamMute (32,580)**

**RQ1** Do Copy-Paste-Prompting methods effectively enhance contextual faithfulness and mitigate RAG hallucinations through high-copying response generation?

**RQ2** Does training with high-copying responses from Copy-Paste-Prompting as DPO preference trajectories enable CopyPasteLLM to genuinely trust contextual knowledge—even when it is counterfactual?

**RQ3** What are the underlying mechanisms of CopyPasteLLM's contextual belief? We will interpret this by analyzing logits and hidden states.

Table 4: Datasets and their roles across 3 research questions. **Train** refers to the number of samples utilized for training our CopyPasteLLM, and **Eval** refers to the number of samples used for evaluation. The 20,000 samples of the PubMedQA Artificial subset were randomly sampled using the random seed 42 from the 211k entries.

Dataset	Subset	Domain	Size	Gold Answer	RQ1	RQ2	RQ3
RAGTruth	QA	Daily-Life	839	✗	Eval	only Train (16)	Eval
FaithEval	Counterfactual	Science	1,000	✓	Eval	Train / Eval (241 / 759)	Eval
PubMedQA	Labeled	Biomedicine	1,000	✓	Eval	Train / Eval (108 / 892)	Eval
PubMedQA	Artificial	Biomedicine	20,000	✓	-	Eval	-
ConFiQA	Counterfactual & Original	Wikidata	36,000	✓	-	Eval	-

# Results: Copy-Paste-Prompting (RQ1)

Copy-Paste-Prompting significantly reduces hallucinations by guiding models to directly quote source text. Experiments confirm a significant inverse correlation between response Copying Degree and context-unfaithful hallucinations.

Table 2: Performance comparison of Copy-Paste-Prompting against baselines across models and datasets. Methods with colored backgrounds are our proposed Copy-Paste-Prompting. **Bold** indicates the best performance, underlined indicates the second-best performance. *Faith.*: Faithfulness (*M.C.*: MiniCheck, *A.S.*: AlignScore), *Hallu.*: Hallucination, *Flu.*: Fluency.

Method	RAGTruth				FaithEval				PubMedQA				AVERAGE					
	Faith.		Hallu.		Faith.		Hallu.		Faith.		Hallu.		Flu.		Faith.	Hallu.	Flu.	
	M.C.	A.S.	Twist	Causal	M.C.	A.S.	Twist	Causal	M.C.	A.S.	Twist	Causal						
Mistral-7B-Instruct-v0.2 (7B)																		
Attributed	69.58	63.43	1506.9	1494.5	19.54	88.28	90.67	<u>1527.1</u>	1513.7	37.32	75.49	77.90	1464.7	1450.4	23.53	77.56	1492.9	26.80
Citations	57.82	49.39	1472.5	1475.7	<b>14.41</b>	73.50	74.25	1392.1	1416.2	27.98	55.79	52.35	1415.9	1370.0	<b>13.93</b>	60.52	1423.7	<b>18.77</b>
CP-Link	89.39	<b>75.45</b>	1518.9	1519.5	73.33	93.41	92.44	1510.9	1521.9	49.40	<b>96.50</b>	<b>88.52</b>	1518.4	<b>1580.7</b>	35.57	<b>89.29</b>	1528.4	52.77
CP-Order	<b>91.25</b>	71.98	1467.9	1472.4	65.62	<b>94.89</b>	<u>92.27</u>	1522.6	1501.5	43.74	93.18	82.35	<u>1528.3</u>	<u>1559.1</u>	32.65	<u>87.65</u>	1508.6	47.34
CP-Refine	82.18	<u>74.56</u>	<b>1533.8</b>	<b>1537.9</b>	<u>18.46</u>	92.85	<b>94.68</b>	<b>1547.4</b>	<b>1546.7</b>	<b>26.63</b>	91.52	88.21	<b>1572.7</b>	1539.7	17.79	87.33	<b>1546.4</b>	<u>20.96</u>
Llama-3.1-8B-Instruct (8B)																		
Attributed	57.02	65.29	<u>1526.3</u>	<u>1554.3</u>	26.22	85.22	85.65	1516.5	1536.9	330.8	71.10	60.01	1530.0	1553.1	47.36	70.72	1536.2	134.8
Citations	64.27	72.81	1428.5	<b>1574.4</b>	<b>16.78</b>	88.81	86.80	1486.2	<b>1555.6</b>	39.65	78.56	73.03	1403.4	1463.4	19.11	77.38	1485.3	25.18
CP-Link	70.58	78.83	1401.1	1328.3	17.83	91.54	89.23	1456.2	1366.3	<b>24.09</b>	80.74	80.79	1396.4	1371.1	19.65	81.95	1386.6	<b>20.52</b>
CP-Order	75.30	<b>94.81</b>	1498.4	1498.0	26.35	<b>95.44</b>	<b>98.12</b>	<b>1523.2</b>	<u>1541.2</u>	33.46	87.07	<b>97.62</b>	<b>1633.6</b>	<b>1559.1</b>	27.83	<b>91.39</b>	1542.3	29.21
CP-Refine	<b>77.30</b>	<b>88.52</b>	<b>1645.7</b>	1545.0	17.75	94.40	93.71	1517.9	1500.1	26.99	<b>87.29</b>	91.19	1536.5	1553.2	<b>18.64</b>	88.74	<b>1549.7</b>	21.13
Qwen2.5-72B-Instruct (72B)																		
Attributed	57.00	62.23	1504.5	1525.5	19.68	85.74	83.03	1537.3	1490.0	293.8	77.99	69.25	1509.9	1441.5	33.42	72.54	1501.5	115.6
Citations	74.32	77.52	1455.5	1498.0	<b>18.61</b>	90.98	88.30	1456.5	1476.7	34.67	82.01	76.62	1358.8	1413.6	<b>22.89</b>	81.63	1443.2	25.39
CP-Link	75.75	85.37	1446.3	1363.2	27.47	92.88	92.00	1443.5	1424.2	39.55	86.21	88.58	1527.9	1489.2	33.43	86.80	1449.1	33.48
CP-Order	<u>76.32</u>	<b>94.60</b>	<u>1509.2</u>	<b>1589.6</b>	30.56	<b>95.78</b>	<b>98.16</b>	<b>1539.3</b>	<b>1579.7</b>	38.11	87.85	<b>97.52</b>	1546.8	1575.9	35.26	<b>91.71</b>	<b>1556.8</b>	34.65
CP-Refine	<b>78.14</b>	<u>90.88</u>	<b>1584.6</b>	1523.7	20.12	94.72	<u>95.48</u>	1523.4	1529.4	<b>27.65</b>	<b>88.88</b>	<u>95.04</u>	<b>1556.7</b>	<b>1579.9</b>	<b>20.29</b>	90.52	1549.6	<b>22.69</b>
DeepSeek-V3-0324 (671B)																		
Attributed	56.42	59.60	1417.1	1449.1	<u>27.52</u>	86.90	83.46	<u>1524.3</u>	1535.0	63.27	75.56	69.24	1449.2	1487.9	36.88	71.86	1477.1	42.56
Citations	62.32	64.45	1510.8	1565.6	34.63	87.38	85.69	1463.0	1477.0	36.09	75.93	71.85	1460.4	1387.5	23.27	74.60	1477.4	31.33
CP-Link	70.59	72.54	1382.9	1360.3	34.19	92.60	88.08	1469.1	1374.8	35.55	81.56	77.67	1380.9	1351.1	28.54	80.51	1389.9	32.76
CP-Order	75.53	<b>92.87</b>	1579.4	1555.2	59.11	<b>95.23</b>	<b>97.79</b>	<b>1569.9</b>	1548.1	34.30	87.20	<b>97.38</b>	1561.8	1621.7	27.56	<b>91.00</b>	1572.7	40.32
CP-Refine	<b>77.14</b>	<u>90.02</u>	<b>1609.8</b>	<b>1569.7</b>	<b>22.57</b>	94.45	<u>93.06</u>	1453.7	<b>1565.2</b>	<b>33.84</b>	<b>87.39</b>	<u>91.05</u>	<b>1647.7</b>	<b>1651.7</b>	<b>21.91</b>	88.85	<b>1583.0</b>	<b>26.11</b>

CP-Refine achieves the best balance of faithfulness and fluency while substantially reducing hallucination metrics

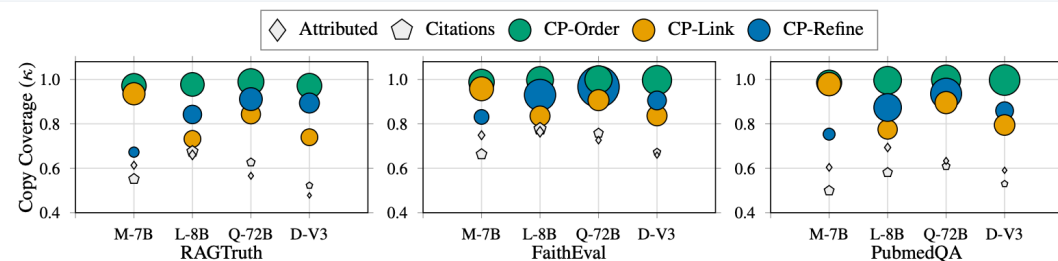


Figure 5: Copying degree across models and datasets. Copy-Paste-Prompting methods significantly outperform baselines in  $\kappa$  and  $\delta$  (area of point). Notably, the copying degree varies by dataset nature (FaithEval > PubMedQA > RAGTruth) and model capacity, with DeepSeek-V3 balancing copying and query relevance effectively.

Experiments confirm Copy-Paste variants achieve significantly higher lexical reuse than conventional methods, with models adaptively adjusting copying behavior based on constraint strength

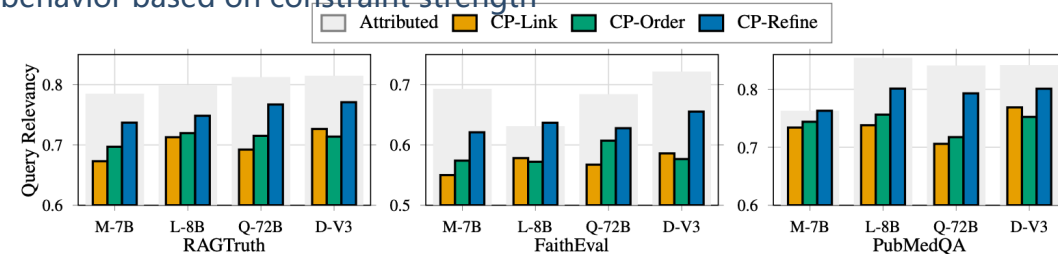


Figure 6: Query relevancy performance. CP-Refine consistently yields the most relevant responses. The efficacy of CP-Link is model-dependent; only the highly capable DeepSeek-V3 utilizes the linking mechanism to improve relevance over the rigid CP-Order approach.

High-degree source quoting does not compromise response quality; CP-Refine maintains high query relevancy scores while preserving high faithfulness

# Results: CopyPasteLLM Counterfactual Performance (RQ2)

CopyPasteLLM achieves SOTA with only 1/50th of the strongest baseline's training data

Table 1: Counterfactual scenarios: Performance comparison of CopyPasteLLM against baselines. We removed 241 samples used for training CopyPasteLLM from FaithEval, with the remaining samples used for testing (detailed in the RQ2 setup of Appendix Table 4). Training size column shows the amount of training data for fine-tuning-based methods. <sup>T</sup> indicates seen data for the respective model. **Bold** values highlight the best performing method in unseen settings.

Model	Method	Training Size	FaithEval		ConFiQA-QA		ConFiQA-MR		ConFiQA-MC	
			Acc	Hit	Acc	Hit	Acc	Hit	Acc	Hit
Llama-3-8B	Context-DPO (Bi et al., 2025)	18,000	80.2	36.7	88.9 <sup>T</sup>	96.1 <sup>T</sup>	88.4 <sup>T</sup>	85.8 <sup>T</sup>	92.1 <sup>T</sup>	80.9 <sup>T</sup>
	Attributed (Zhou et al., 2023)	-	67.1	34.2	51.5	91.4	53.3	71.5	37.3	53.6
	CoCoLex (T.y.s.s et al., 2025)	-	69.2	17.9	48.5	37.4	53.9	14.8	36.1	15.5
	Canoe (Si et al., 2025)	10,000	71.4	34.0	64.3	93.2	66.6	<b>83.8</b>	64.5	73.7
	ParamMute (Huang et al., 2025b)	32,580	68.5	22.5	74.4	82.2	75.5	72.4	81.4	70.2
	CopyPasteLLM (Ours)	<b>365</b>	<b>92.8</b>	<b>37.2</b>	<b>83.6</b>	<b>96.7</b>	<b>80.9</b>	83.4	<b>86.8</b>	<b>75.9</b>
Mistral-7B-v0.2	Context-DPO (Bi et al., 2025)	18,000	77.1	33.8	84.8 <sup>T</sup>	94.8 <sup>T</sup>	81.3 <sup>T</sup>	85.3 <sup>T</sup>	80.4 <sup>T</sup>	80.8 <sup>T</sup>
	Attributed (Zhou et al., 2023)	-	65.6	32.0	56.6	84.4	29.2	69.8	39.0	57.4
	CoCoLex (T.y.s.s et al., 2025)	-	65.3	35.4	57.3	50.8	41.8	33.5	32.5	33.7
	CopyPasteLLM (Ours)	<b>365</b>	<b>89.3</b>	<b>41.8</b>	<b>84.4</b>	<b>95.0</b>	<b>80.8</b>	<b>90.8</b>	<b>82.5</b>	<b>86.3</b>
Llama-3.1-8B	Attributed (Zhou et al., 2023)	-	65.5	32.0	49.9	88.4	39.8	69.2	15.5	52.6
	CoCoLex (T.y.s.s et al., 2025)	-	68.1	36.2	48.5	57.3	40.4	38.4	13.5	37.2
	CopyPasteLLM (Ours)	365	<b>92.6</b>	<b>41.0</b>	<b>72.4</b>	<b>90.1</b>	<b>75.4</b>	<b>84.8</b>	<b>83.5</b>	<b>79.9</b>

Model	Accuracy (%)
Mistral-7B-Instruct-v0.3	73.8
Llama-3.1-8B-Instruct	68.5
Llama-3-8B-Instruct	66.5
Mistral-Nemo-Instruct-2407	58.3
gpt-3.5-turbo	57.1
Command R	69.3
Phi-3.5-mini-instruct	66.8
Command R+	73.6
gemma-2-9b-it	55.7
gemma-2-27b-it	55.7
gpt-4o-mini	50.9
Phi-3-mini-128k-instruct	75.7
Phi-3-medium-128k-instruct	60.8
Llama-3.1-70B-Instruct	55.2
Llama-3-70B-Instruct	60.5
Claude 3.5 Sonnet	73.9
gpt-4-turbo	41.2
gpt-4o	47.5
CopyPasteLLM (Based on Llama-3-8B-Instruct)	<b>92.8</b>
CopyPasteLLM (Based on Mistral-7B-Instruct-v0.2)	89.3
CopyPasteLLM (Based on Llama-3.1-8B-Instruct)	<u>92.6</u>

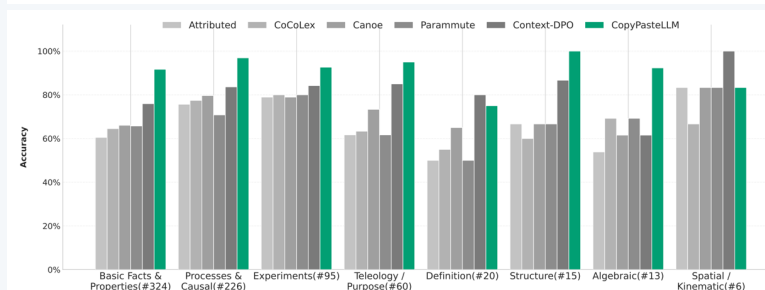


Figure 7: Performance comparison across diverse knowledge domains. CopyPasteLLM consistently outperforms or remains highly competitive against strong baselines across most categories, demonstrating robustness in both factual (e.g., *Basic Facts*) and reasoning-intensive domains (e.g., *Processes & Causal*, *Experiments*).

Table 3: Accuracy in non-counterfactual settings. PubMedQA is evaluated on artificial subset 20,000 samples (none used for CopyPasteLLM training, see Appendix Table 4). ConFiQA uses Original context and Original answers.

Method	Mistral-7B-v0.2				Llama-3-8B				Llama-3.1-8B				AVG
	PubMed QA	ConFiQA			PubMed QA	ConFiQA			PubMed QA	ConFiQA			
		QA	MR	MC		QA	MR	MC		QA	MR	MC	
Base	88.60	96.22	71.20	72.27	97.3	98.02	93.00	91.02	<b>98.15</b>	97.93	89.48	89.97	90.26
CopyPasteLLM (Ours)	<b>91.40</b>	<b>97.43</b>	<b>91.87</b>	<b>91.20</b>	<b>97.5</b>	<b>99.30</b>	<b>97.17</b>	<b>96.27</b>	97.67	<b>99.02</b>	<b>94.95</b>	<b>94.92</b>	<b>95.73</b>

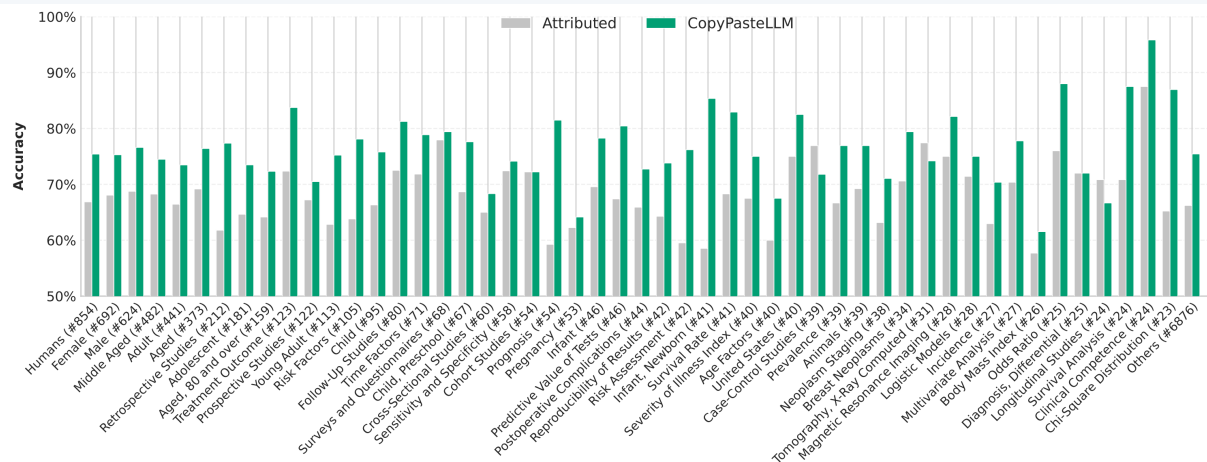


Figure 9: Domain-specific performance analysis of CopyPasteLLM (based on Llama-3-8b-instruct) on PubMedQA, categorized by Medical Subject Headings (MeSH).

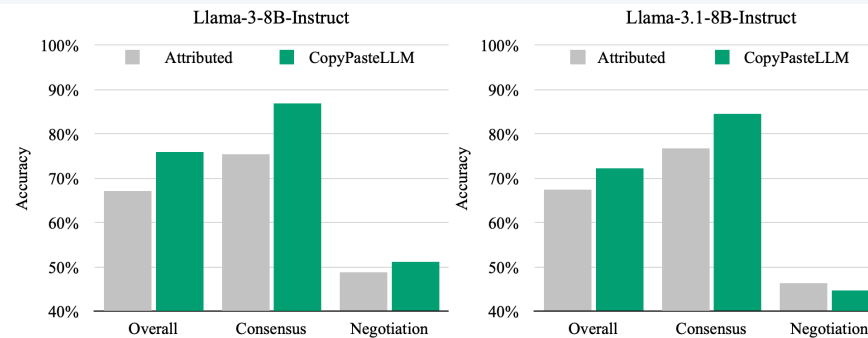
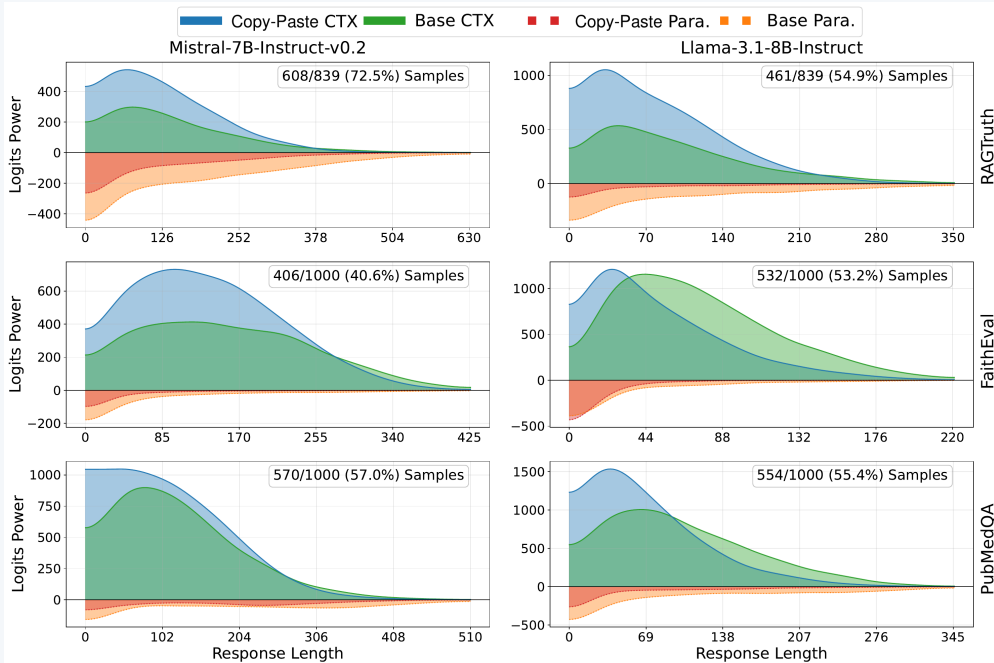
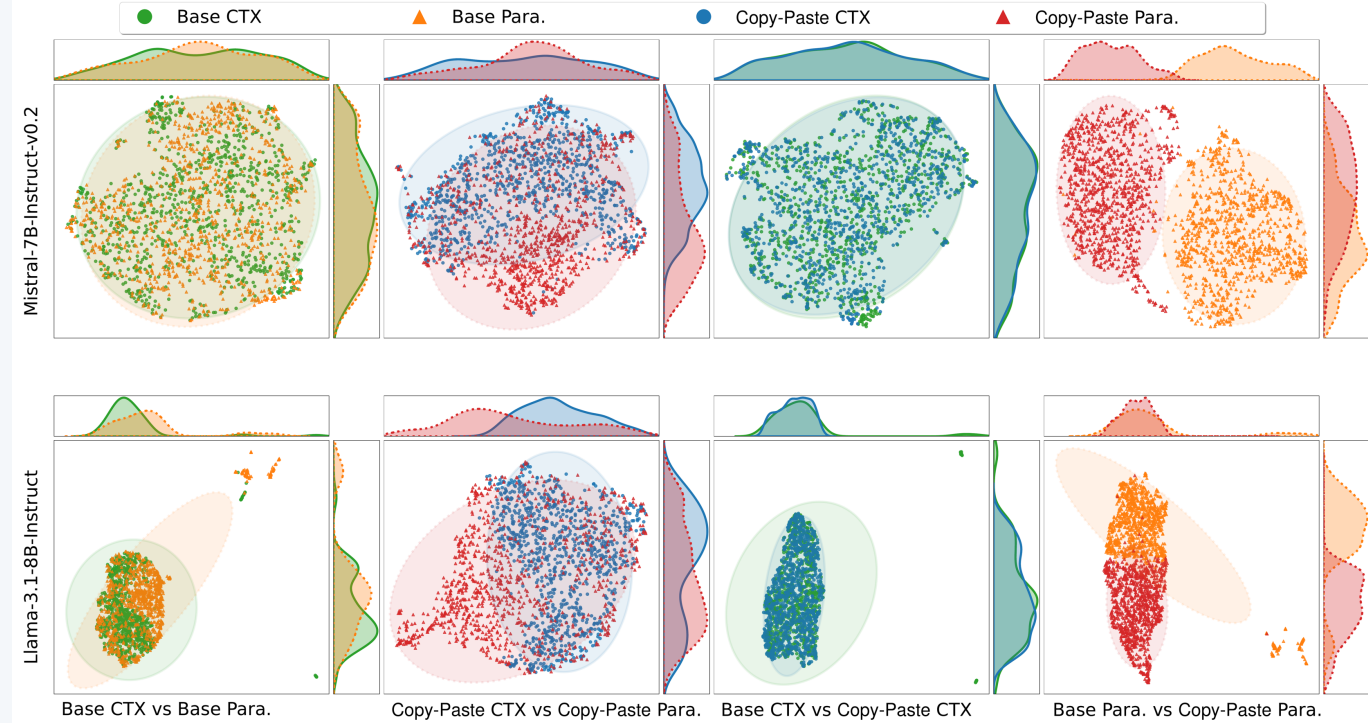


Figure 8: Performance breakdown on PubMedQA-Labeled by reasoning difficulty. Accuracy is compared between Attributed and CopyPasteLLM models across the Consensus (clear evidence) and Negotiation (ambiguous context) subsets, demonstrating the highest gains in samples with explicit evidence.

## Logit Distribution (CTX vs. Para.)



## Hidden States (CTX vs. Para.)



### Logits Findings

- (1) CopyPasteLLM shows stronger CTX utilization + reduced Para. reliance
- (2) Peak CTX engagement occurs earlier in generation
- (3) Indicating enhanced initial contextual trust

### Hidden States Finding: CopyPasteLLM does not "enhance" attention to external context, but rather "suppresses" overconfidence in internal parametric knowledge.

- (1) CTX representations remain close to base model
- (2) Para. representations differ substantially
- (3) **Mechanism: Selective parametric recalibration**

**Observation** Inverse correlation between copying degree and hallucination density on RAGTruth

**Paradigm** Copy-Paste: Embed contextual fragments -> Achieve both faithfulness and attribution

**Method** Two-stage: CP-Prompting (3 variants) + CopyPasteLLM (DPO, only 365 samples)

**Results** +12.2% to +24.5% on FaithEval counterfactual, 50x more data-efficient than baselines

**Interpretability** Context-Parameter Copying Capturing: Recalibrates parametric confidence, not enhancing CTX

**Thank you!**

# Evaluation and mitigation of the limitations of large language models in clinical decision-making

## 内容简介

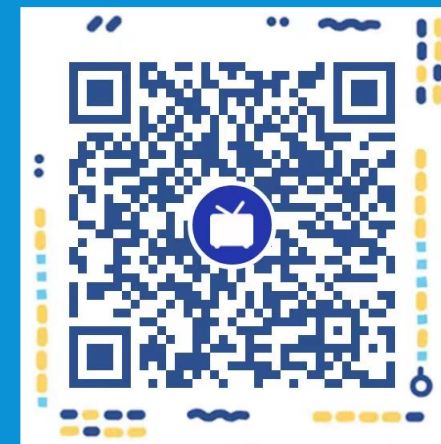
本报告将分享关于大型语言模型（LLMs）在临床决策中应用的研究。研究团队通过扩展MIMIC-IV数据库，创建了MIMIC-CDM数据集，该数据集包含了2400名因急性腹痛到急诊室就诊的患者数据。研究中构建了一个评估框架，用以模拟LLMs在临床决策过程中的信息收集、综合分析和诊断治疗计划制定的能力。报告将概述LLMs在模拟临床环境中的表现，以及它们在遵循临床指南、请求必要诊断检查和提出治疗建议方面的潜力和挑战。此外，还将讨论LLMs在处理临床信息时的可靠性和鲁棒性，包括对指令的遵循、信息顺序的敏感性，以及在不同信息量下的性能变化。这项研究为未来临床AI模型的发展提供了宝贵的见解，并探讨了实现更安全、更有效临床决策支持系统的可能途径。



## 嘉宾简介

龙泳潮 (Yongchao Long)，天津理工大学计算机科学与工程学院博士生，导师为周雨熙副教授和洪申达副研究员。在清华-北大-天理医疗大模型联合研究团队（团队负责人主要包括清华大学医工交叉研究院副院长邢春晓教授、北京大学健康医疗大数据国家研究院洪申达副研究员、国家杰青天津理工大学副校长陈胜勇教授），从事医疗大模型相关的研究工作，以提高大模型在医疗场景中的效用，并积极参与旨在将人工智能技术融入临床实践的项目。

站账号 AIMEL\_医学人工智能联盟



每周二

19:00-20:00



*Health Data Science* (《健康数据科学(英文)》) 由北京大学主办, 作为Science合作伙伴期刊由美国科学促进会(AAAS)全球出版发行。

主编: 詹启敏, 北京大学

副主编: 张路霞, 北京大学

刘宏芳, UTHealth Houston



·编辑委员会由来自中国、美国、英国、新加坡、韩国、以色列和瑞士等国家的60多名成员组成。

WoS检索, 2026年将获得首个影响因子  
入选卓越二期英文单刊

## 数促健康, 智赋医学

- 关注发表高影响力的原创性研究、综述、观点和评论
- 采用最佳出版系统和流程以确保出版的高效和质量。

**我们的使命:** 基于对健康领域的深刻认知和理解, 倡导跨领域合作, 聚焦数据科学与前沿技术在健康医疗领域的落地应用和科学发现, 为应对多种健康挑战提供证据和支持。

**我们关注的领域:** 前沿技术在健康医疗中的应用及评价; 基于健康医疗大数据的医疗、公共卫生、健康服务及卫生政策研究; 健康医疗数据集的描述与发表; 健康医疗数据管理、治理及溯源问题等。

数据库收录:



Welcome to follow our team by scanning the QR code on WeChat  
欢迎微信扫码关注我们的团队



微信搜一搜

PKUDigitalHealth