

Shiping Yang<sup>1,2</sup>, Jie Wu<sup>3</sup>, Wenbiao Ding<sup>2</sup>, Ning Wu<sup>2</sup>, Shining Liang<sup>2</sup>, Ming Gong<sup>3</sup>, Hongzhi Li<sup>4</sup>, Hengyuan Zhang<sup>5</sup>  
Angel X. Chang<sup>1,6</sup>, Dongmei Zhang<sup>2</sup>

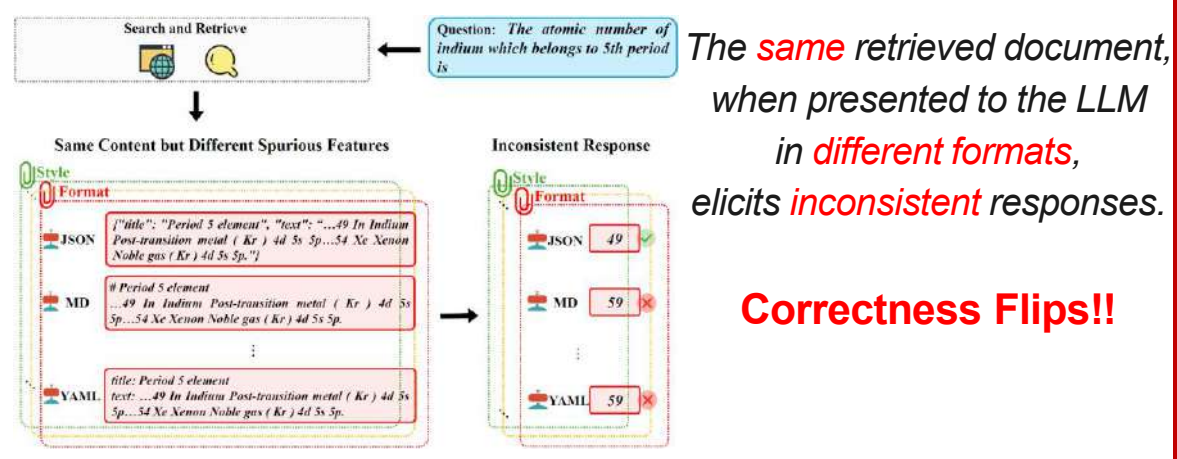
<sup>1</sup>Simon Fraser University · <sup>2</sup>Microsoft · <sup>3</sup>Atlassian · <sup>4</sup>Tongji University · <sup>5</sup>The University of Hong Kong · <sup>6</sup>Canada-CIFAR AI Chair, Amii

## Background and Motivation

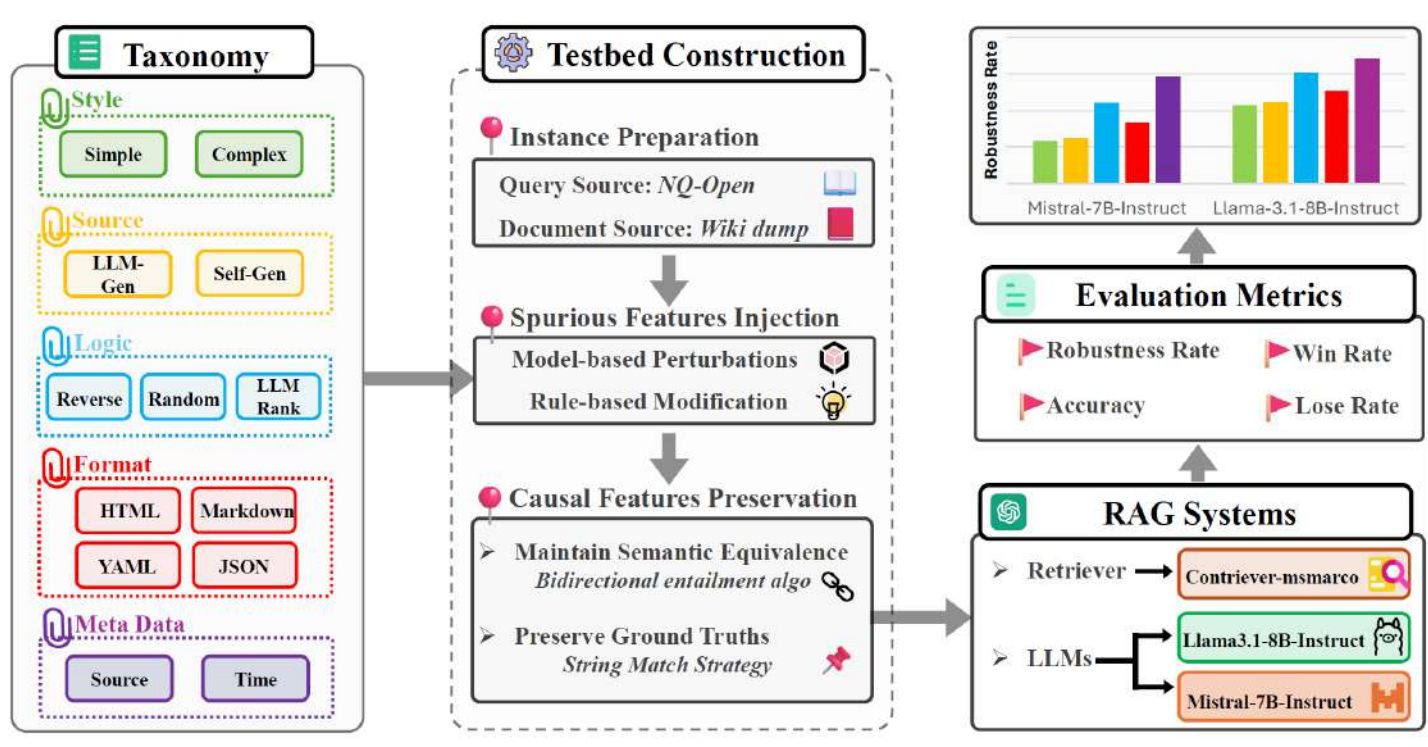
RAG mitigates hallucinations by grounding answers in retrieved documents. Prior work has shown that RALMs are sensitive to *irrelevant or counterfactual* grounding data — a well-known *RAG robustness* issue. However, these studies all examine *explicit noise* that alters semantics.

In practice, the Internet is saturated with documents that *share the same meaning* while differing in *style, format, source, ordering, or metadata*. We instead study *implicit noise: spurious features*, i.e., semantic-agnostic perturbations, and ask whether they likewise affect RALM responses.

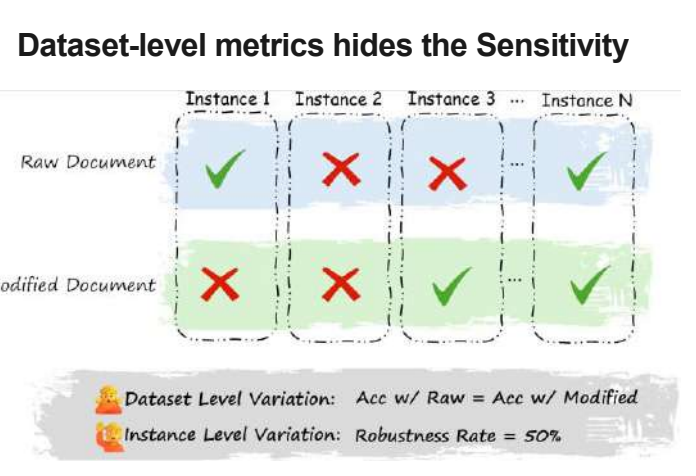
## The Phenomenon: Sensitive to Spurious Features



## SURE Framework: Perturb then Evaluate



## How to Quantify Robustness

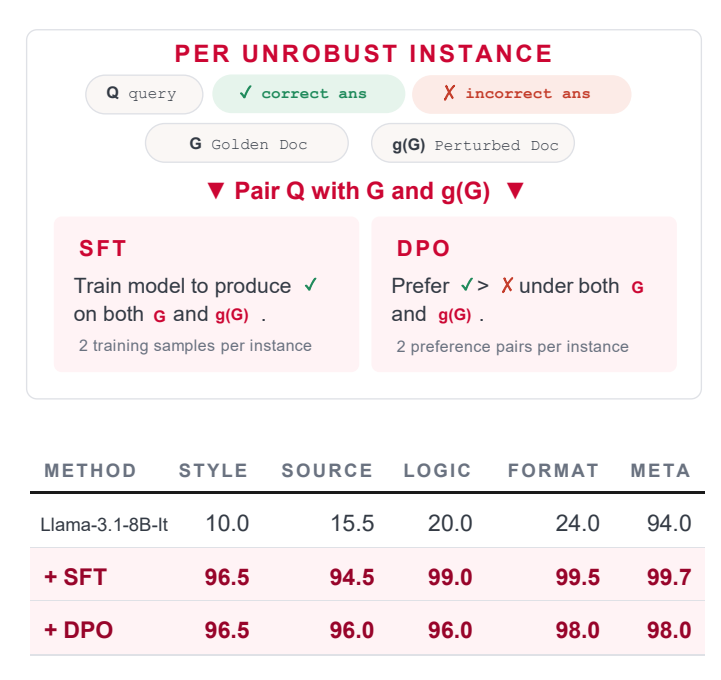


Robustness Rate	Win Rate (WR)	Lose Rate (LR)
Prediction unchanged	Incorrect to correct	Correct to incorrect

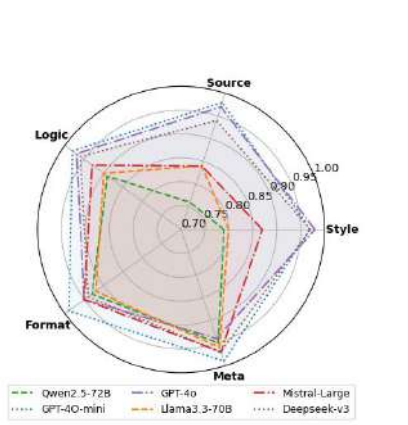
## Main Results

Taxonomy	Perturbations	Known-Golden					Known-Noise					Unknown-Golden					U-N
		LR	RR	WR	Org	Acc	LR	RR	WR	Org	Acc	LR	RR	WR	Org	Acc	
Style	Simple	7.79	83.04	<b>9.18</b>	66.03	67.42	1.70	95.80	<b>2.50</b>	4.12	4.92	8.43	82.88	<b>8.69</b>	51.42	51.68	99.45
	Complex	6.00	85.60	<b>8.40</b>	66.03	68.43	1.91	96.59	1.50	4.12	3.71	6.71	84.86	<b>8.43</b>	51.42	53.13	99.57
Source	LLM-Generated	5.89	86.43	<b>7.69</b>	65.62	67.43	1.43	96.83	<b>1.74</b>	4.13	4.45	6.20	85.71	<b>8.09</b>	49.15	51.04	99.56
	Self-Generated	6.55	85.01	<b>8.44</b>	65.62	67.52	1.55	96.37	<b>2.09</b>	4.13	4.67	6.52	86.36	<b>7.12</b>	49.15	49.74	99.57
Logic	Reverse	5.06	90.82	4.12	62.01	62.01	1.13	97.82	1.06	4.36	5.73	89.71	4.56	44.66	44.66	99.67	
	Random	3.91	93.16	2.93	62.95	61.97	0.86	98.31	0.83	4.43	4.40	4.21	91.67	4.12	45.84	45.74	99.72
	LLM-Ranked	3.24	93.93	2.83	62.54	62.54	0.82	98.43	0.74	4.36	3.58	93.36	3.06	45.32	45.32	99.76	
Format	JSON	7.01	88.25	4.74	63.91	61.64	1.70	97.25	1.05	3.21	5.92	89.63	4.45	47.88	47.88	99.61	
	HTML	11.85	84.46	3.69	63.91	55.75	2.70	96.90	0.40	3.87	1.56	9.33	86.78	3.90	49.35	43.92	99.61
	YAML	5.26	89.94	4.80	63.91	63.45	1.26	97.41	1.33	3.94	4.79	90.80	4.41	49.35	48.97	99.67	
	Markdown	2.32	92.23	<b>5.45</b>	63.91	67.04	0.60	96.89	<b>2.51</b>	3.87	5.77	2.34	93.46	<b>4.19</b>	51.20	51.20	99.61
Metadata	Timestamp (pre)	2.08	95.81	<b>2.11</b>	55.77	55.80	0.28	99.42	<b>0.29</b>	1.59	2.54	95.56	1.90	42.66	42.66	99.95	
	Timestamp (post)	2.04	95.86	<b>2.10</b>	55.77	55.84	0.25	99.43	<b>0.32</b>	1.58	1.64	2.81	95.56	1.63	43.31	42.12	99.95
	Datasource (wiki)	2.11	93.45	<b>4.44</b>	55.77	58.10	0.23	98.96	<b>0.81</b>	1.58	2.17	3.25	92.47	<b>4.27</b>	44.33	44.33	99.86
	Datasource (twitter)	2.27	94.11	<b>3.62</b>	55.77	57.11	0.31	99.25	<b>0.43</b>	1.58	1.70	2.77	93.97	<b>3.25</b>	43.79	43.79	99.91

## How to Improve Robustness



## Further Analysis



### TAKEAWAY 1

Spurious features are a **widespread** robustness gap in RAG — independent of **internal knowledge, architecture, and scale**.

### TAKEAWAY 2

Not every spurious feature is harmful — **WR > LR** on some perturbations shows certain features can even help.

### TAKEAWAY 3

Traditional strategies for improving robustness against **explicit noise** do not work for spurious features; training on SURE-synthesized pairs (SFT / DPO) effectively mitigates this sensitivity.

## Contact Information

Shiping Yang  
Email: yangshipingnl@gmail.com

PAPER CODE Personal Webpage