The Science of Evaluation for Large Language Models

Lei Li Language Technologies Institute Carnegie Mellon University June 19, 2025

Large Language Models drive the Productivity

Translate Summarize

Editing

Write email



Answer questions Suggest names Write code Recommend restaurants

How good is LLM generation?



Evaluation

LLM output: The outbreak of the new crown crisis

Reference: The outbreak of the COVID-19 crisis

Prompt: Translate "新冠疫情危机爆发".



Metrics: comparing output against references, used for testing.



Reward / Quality estimation (QE) model.

Rule-based and Learned Metrics

Rule-based

Supervised Metric

- BLEU
- chrF
- TER

- BLEURT
- COMET
- MetricX

- **Unsupervised Metric**
 - SEScore
 - BERTScore
 - PRISM
 - BARTScore

- ROUGE Human rating is scarce Only surface form difference

LLM as evaluator?

Learning from Reward / Quality-Estimation Model(QE)



Challenges in Evaluating LLM

- BLEU/ROUGE will have significantly decreased correlations with human judgments.
- Comprehensive tasks instead of just one task (e.g. MT)
- Open-end generation tasks
- What if no ground truth is given?
 Source-based evaluation is difficult

Outline

- Can we trust LLM evaluators?
 Self-bias in LLM Evaluators (source-based)
 - Evaluating LLM Generation Quality

 InstructScore: Interpretable text generation score
 - Assessing Knowledge in LLMs (KaRR)

LLM as an Evaluator? (source-based)

Prompt: Translate " 新冠疫情危机爆发 ".

LLM output: The outbreak of the new crown crisis

ask LLM: how good is the above translation? (MQM scheme: major error=-5, minor error=-1) LLM output: -5



Aman Madaan, Niket Tandon ..., and Peter Clark. 2023. <u>Self-refine: Iterative refinement with self-feedback.</u> Xinyun Chen, Maxwell Lin, Nathanael Schärli, and Denny Zhou. 2023. <u>Teaching large language models to self-debug</u>.

LLM (GPT4) evaluator highly correlates with human evaluation



Liu et al. G-EVAL: NLG Evaluation using GPT-4 with Better Human Alignment. 2023. Chen et al. Exploring the Use of Large Language Models for Reference-Free Text Quality Evaluation: An Empirical Study. 2023.

But, are LLM evaluators fair? GPT4 evaluator gives higher scores to its generation!



Liu et al. G-EVAL: NLG Evaluation using GPT-4 with Better Human Alignment. 2023.

A Translation Example

Yoruba text: Ní bayii a ni àwon eku oloshu merin ti ko ni dayabetesi telele to ti ni ayabetesi," o she afikun.

GPT-4's translation: At this point, we have four rats without diabetes that have developed diabetes," he added.

Using LLM self-evaluate and refine

Human Post Edits: At this point, we have 4-month-old rats mice without diabetes that have developed diabetes that are non-diabetic that used to be diabetic," he added.

Major error (-5) Minor error (-1)

GPT-4's evaluation: At this point, we have four rats without diabetes that have developed diabetes," he added.

Self-refine

Human Score:
-11





LLM self-refine leads to inflated self-score!

Human Post Edits: Currently, we have 4-month-old healthy rats mice that have developed diabetes that are non-diabetic that used to be diabetic ," he clarified.

Major error (-5)

Minor error (-1)

GPT-4's evaluation: "Currently, we have four healthy rats that have developed diabetes," he clarified.



Self-refine







LLM self-refine leads to inflated self-score!



GPT-4's evaluation: Presently, we have four non-diabetic rats that have developed diabetes," he elaborated.



While <u>GPT-4</u> thinks it performed self-refine, humans observe all errors persist

LLM 1st generation: <u>At this point</u>, we have four <u>rats</u> without diabetes that have developed diabetes," he added.

LLM 2nd generation: "<u>Currently</u>, we have four <u>healthy rats</u> that have developed diabetes," he clarified.

LLM 3rd generation : Presently, we have four non-diabetic rats that have developed diabetes," he elaborated.

Defining bias in LLM Evaluators

Bias definition 1: Statistical Bias Estimation



Defining bias in LLM

Bias definition 2: Distance Skewness estimation



Self-Bias Amplifies in LLM Translation



What is the root cause of self-bias amplification?

GPT-4 and Gemini overestimate improvements in selfrefined outputs, compared to actual performance measured by BLEURT

Self-Bias Amplifies in LLM Data-to-Text and Math



Wenda Xu, Guanglei Zhu, Xuandong Zhao, Liangming Pan, Lei Li, William Yang Wang. Pride and Prejudice: LLM Amplifies Self-Bias in Self-Refinement. ACL 2024

What is improving at Self-refine if not quality?

Self-refine improves understanding and fluency of the text



Key insights

- LLM evaluators have strong self-bias
- Self-bias is amplified during LLM self-refine/self-rewarding process
- Self-refine can improve fluency of text but not necessarily quality
- LLMs favor texts that follow their 'style'



Outline

- Can we trust LLM evaluators?
 Self-bias in LLM Evaluators (source-based)
- Evaluating LLM Generation Quality

 InstructScore: Interpretable text generation score
 - Assessing Knowledge in LLMs (KaRR)

When you made a mistake...

Teacher 1: You have a bad translation. You get score of 20/100

outbreak of

Teacher 2: 'New crown' is a major mistranslation error. The correct translation is 'COVID-19'. Score: 20/100

Evaluating Text Generation Quality – Existing metrics



Ideal Metric: Fine-grained Explanation

Error location: new crown

Reference: The outbreak of the COVID-19 crisis



Candidate: The outbreak of the new crown crisis **Error type:** Terminology is used inconsistently

Major/Minor: Major

Explanation: The term " new crown" is not the correct term for "Covid-19".

Direct Prompting ChatGPT

Raw text: "The art ... between providing enough detail to ... too much information."



Error type 1: Translation includes information not present in the correct translation

Major/minor: Major

Incorrect generation: [GPT4 fill in] Error location 1: [GPT4 fill in] Explanation for error 1: [GPT4 fill in]

But, failed explanation in GPT4



Error type 3: Missing information

Explanation for error 3: The incorrect translation adds the word "annual" to the phrase ...

Error type is inconsistent with explanation

But, failed explanation in GPT4



Evaluated text: The outbreak of the new crown crisis



Hallucination

But, failed explanation in GPT4



Explanation for error 1: The incorrect translation uses the word "annual" instead of "annual"

Explanation is illogical

Failures of GPT4 generated explanation

Fields	Failure Mode	Description (<mark>M is local failure mode</mark> , G is global failure mode)		
Error Type	Inconsistency to explanation	M1: Error type is inconsistent with explanation		
Error Location	Inconsistency to explanation	M2: Error locations are not consistent with the explanation		
	Hallucination	M3: Error locations are not referred in the output text		
Major/Minor	Major/Minor disagreement	M5: Major and minor labels are not correct		
Explanation	Hallucination	M4: Error locations are not referred in the output text		
	Explanation failure	M6: Explanation is illogical		
All 4 Fields	False negative error	G1: Error described in the explanation is not an error		
	Repetition	G2: One error is mentioned more than once among explanations		
	Phrase misalignment	G3: Incorrect phrase and correct phrase are not aligned		
	Mention multiple errors	G4: One error span mentions multiple errors		

Introducing InstructScore





InstructScore: Automatic Feedback

Reference		Error1	Error location	\checkmark
Candidate			Error type	 Image: A second s
Error location 1	/ \		Major/minor	\times
Error Type1	AS		Explanation	\checkmark
Major/Minor	20	Error2	Error location	\checkmark
Explanation1			Error type	\checkmark
Error location2			Major/minor	~
Major/Minor			Explanation	
Explanation2			Alignment Score: 7,	/8

InstructScore: Training with Refinement



InstructScore can judge machine translation!

WMT22 Chinese-to-English Translation



InstructScore can evaluate text generation!



InstructScore vs GPT4

Candidate: Or how long before the restaurant is ready? Can you ask for me? **Reference:** And how much longer does it take for the restaurant to prepare it? Could you help me ask them?



InstructScore vs GPT4

Error Type: Lexical Selection Major/Minor: Major Error Location: "Or" instead of "And"

Error Type: Lexical Selection/Omission **Major/Minor:** Major **Error Location:** "Can you ask for me?"

instead of "Could you help me ask them?"

Error Type: Omission

Major/Minor: Minor

Error Location: Missing "Could" in "Could you help me ask them?"

Error type 1: Incorrect translation is missing content from the correct translation Major/minor: Major Error location 1: how long before

Error type 2: Problems with grammar, other than orthography Major/minor: Minor Error location 2: help me ask





Highlights of InstructScore

- We develop a new model-based evaluation metric for Explainable text generation-based metric and leverage automatic feedback to align with human requirements!
 - 1. Fine-grained Explainability
 - 2. Highly Aligned with Human
 - 3. Generalizability (No human ratings are required!)



Outline

- Can we trust LLM evaluators?
 Self-bias in LLM Evaluators (source-based)
- Evaluating LLM Generation Quality

 InstructScore: Interpretable text generation score
- → Assessing Knowledge in LLMs (KaRR)

LLMs generates Unreliable Answers

• e.g. LLaMA-7B

When did Shakespeare die?



Llama-7B : 23rd April 1616.

LLMs generates Unreliable Answers

• e.g. LLaMA-7B

On what date did William Shakespeare's death occur?

Llama-7B : It was on 23 august 1616.



Knowing versus Guessing

1. Distinguish if text generation stems from genuine knowledge or just high co-occurrence with given text.

William Shakespeare's job is a writer.

John Smith's job is a writer.

Unreliable Factual Knowledge in LLMs

- LLMs often generate unreliable answers given varying prompts.
- Example1: Alpaca-7B

• Example2: ChatGPT

William Shakespeare's job is? Solution: Solu Is William Shakespeare a teacher?

Assessing LLM's Knowledge

 Given varying prompts regarding a factoid question, can a LLM reliably generate factually correct answers?



Challenges in LLM Knowledge Assessment

• Knowledge irrelevant generation: The freely generated results of generative models might be irrelevant to factual knowledge.





Risk Ratio

- In statistics, **risk ratio** estimate the strength of the association between exposures (treatments or risk factors) and outcomes.
- Example: a disease noted by *D*, and no disease noted by ¬*D*, exposure noted by *E*, and no exposure noted by ¬*E*. The risk ratio can be written as:

• Risk Ratio =
$$\frac{P(D|E)}{P(D|\neg E)}$$

	E (exposure)	$\neg E$ (no exposure)
D (disease)	P(D E)	P(D ¬E)
¬D (no disease)	P(¬D E)	P(¬D ¬E)

Knowledge Assessment Risk Ratio (KaRR)

• Assesses the joint impact of subject and relation symbols on the LLM's ability to generate the object symbol.

Qingxiu Dong, Jingjing Xu, Lingpeng Kong, Zhifang Sui, Lei Li. Statistical Knowledge Assessment for LLMs. Neurips 2023

KaRR Dataset

Broad coverage

 1million entities
 600 relations

Method	Subj. Alias	Obj. Alias	Rel. Alias	Rel. Cvg.
LAMA@1	×	×	×	6.83%
LAMA@10	×	×	×	6.83%
ParaRel	×	×	 Image: A second s	6.33%
KaRR	√	√	√	100%

"P36": {

```
"capital city": "[X] is the capital city of [Y].",
```

"administrative capital": "[X] is the administrative capital of [Y].",...

```
"P19": {
```

```
"birthplace": "[X]'s birthplace is [Y].",
```

"born in": "[X] was born in [Y].",

"POB": "The POB of [X] is [Y].",

"birth place": "The birth place of [X] is [Y].",

"location of birth": "The location of birth of [X] is [Y].", ...

Qingxiu Dong, Jingjing Xu, Lingpeng Kong, Zhifang Sui, Lei Li. Statistical Knowledge Assessment for LLMs. Neurips 2023

KaRR Scores for 20 LLMs

- Small and mediumsized LLMs struggle with generating correct facts consistently.
- Finetuning LLMs with data from more knowledgeable models can enhance knowledge.

71

Summary of LLM Knowledge Assessment

- Graphical model for knowledge Assessment
- New metric -- KaRR Score
- High human correlation
- Less evaluation bias

Code and data:

73

Summary

- Can we trust LLM evaluators?

 LLM Evaluators exhibit strong bias towards itself
 Self-bias is amplified in LLM self-refine
- Evaluating LLM Generation Quality

 InstructScore: Interpretable text generation score
- Assessing knowledge in LLMs (KaRR)

 KaRR measures how reliable are LLM in generating fact-related answers

Future thoughts

- Evaluating
 - o complex knowledge
 o LLM RAG
 o LLM Agent
- Evaluation for open-end generation

 PerSE at EMNLP 2024

Reference

- Wenda Xu, Guanglei Zhu, Xuandong Zhao, Liangming Pan, Lei Li, William Yang Wang. Pride and Prejudice: LLM Amplifies Self-Bias in Self-Refinement. ACL 2024.
- Xu, Wang, Pan, Song, Freitag, Wang, Li. INSTRUCTSCORE: Explainable Text Generation Evaluation with Finegrained Feedback. EMNLP 2023.
- Dong, Xu, Kong, Sui, Li. Statistical Knowledge Assessment for Large Language Models. NeurIPS 2023.
- Xu, Qian, M. Wang, Li, W. Y. Wang. SESCORE2: Learning Text Generation Evaluation via Synthesizing Realistic Mistakes. ACL 2023.
- Xu, Tuan, Lu, Saxon, Li, Wang. Not All Errors are Equal: Learning Text Generation Metrics using Stratified Error Synthesis (SEScore). EMNLP 2022.