

On the

# Blind Spots

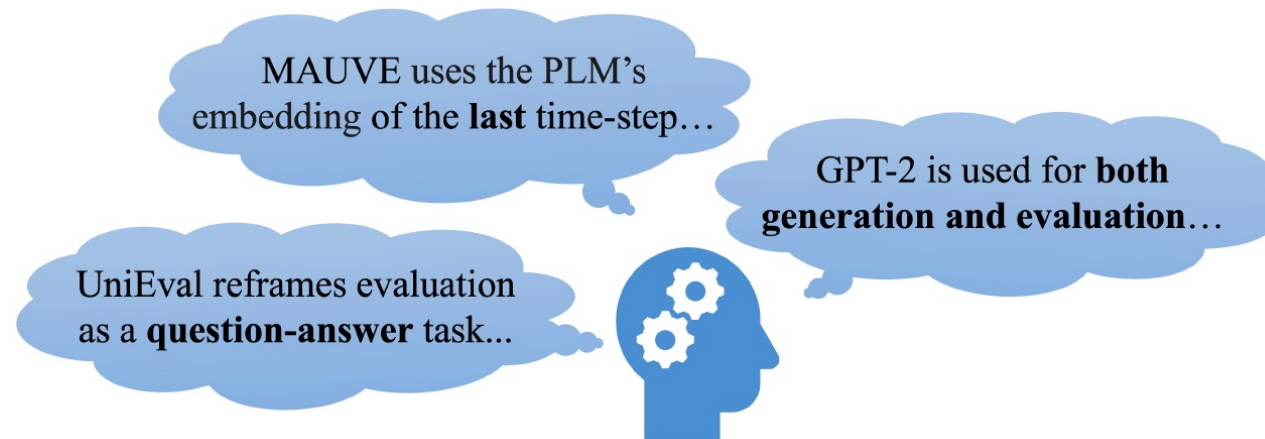
of Model-Based Evaluation Metrics  
for Text Generation

Tianxing He\*, Jack (Jingyu) Zhang\*, Tianle Wang, Sachin Kumar, Kyunghyun Cho, James Glass, Yulia Tsvetkov



# Motivation

- A recent series of work proposed to base text generation evaluation metrics on pretrained language models (PLMs), such as BERTScore, MAUVE, etc.
- Although powerful, the flaws of the underlying PLMs or certain design choices could lead to potential **blind spots** in the evaluation.



- In this work, we design a range of **stress tests** to check for the existence of blind spots.

# A Simple Protocol

- Given a metric in test, we first compute a score on the gold hypothesis set, which is set to be the reference translations/summaries.
- For each test, we apply a synthesized error type (e.g., truncation, random word dropping) to the gold hypothesis set to construct a noised hypothesis set.
- Finally, we compute another score for the noised hypothesis set, and examine whether it is lower than the gold-hypothesis score. If not, we say the metric **fails** the test.



$$\text{Score}(\text{Noised-Hypo}) < \text{Score}(\text{Gold-Hypo}) ?$$

# Test Designs

- We group our tests by their design motivations:
- Metric design choices: *positioned-error, injection, copy-source.*
- PLM properties: *freq-ngram, self-evaluation, repetition.*
- General errors: *fluency (truncation, article removal, etc.), consistency (sentence switching, negation ,etc.).*
- We will cover a subset of our tests in this presentation.

# Tasks and Metrics

- We conduct a number of stress tests for NLG metrics used in three tasks: wiki-103 (open-ended), CNNDM (summarization), WMT21 and TED-MT (both for translation).
- Tested metrics:

Task	Metrics
Open-ended Generation	MAUVE, GPT-PPL, MLM-PPL
Summarization & Translation	BERTScore, MoverScore, BARTScore, COMET
Summarization	UniEval
Translation	PRISM, BLEURT

<- Spoiler: For every PLM-based metric we covered, **at least one** blind spot is found.

# Outline

- Motivation
- Protocol
- The Truncation Test (BERTScore) (<-next!)
- The Injection Test (UniEval)
- The Positioned-Error Test (MAUVE)
- The Self-Evaluation Bias (GPT-PPL BARTScore)
- The Frequent-Ngram Test (GPT-PPL, MLM-PPL)
- Conclusion



# The Truncation Test

- For summarization/translation tasks, we remove a portion from the end of the reference text.

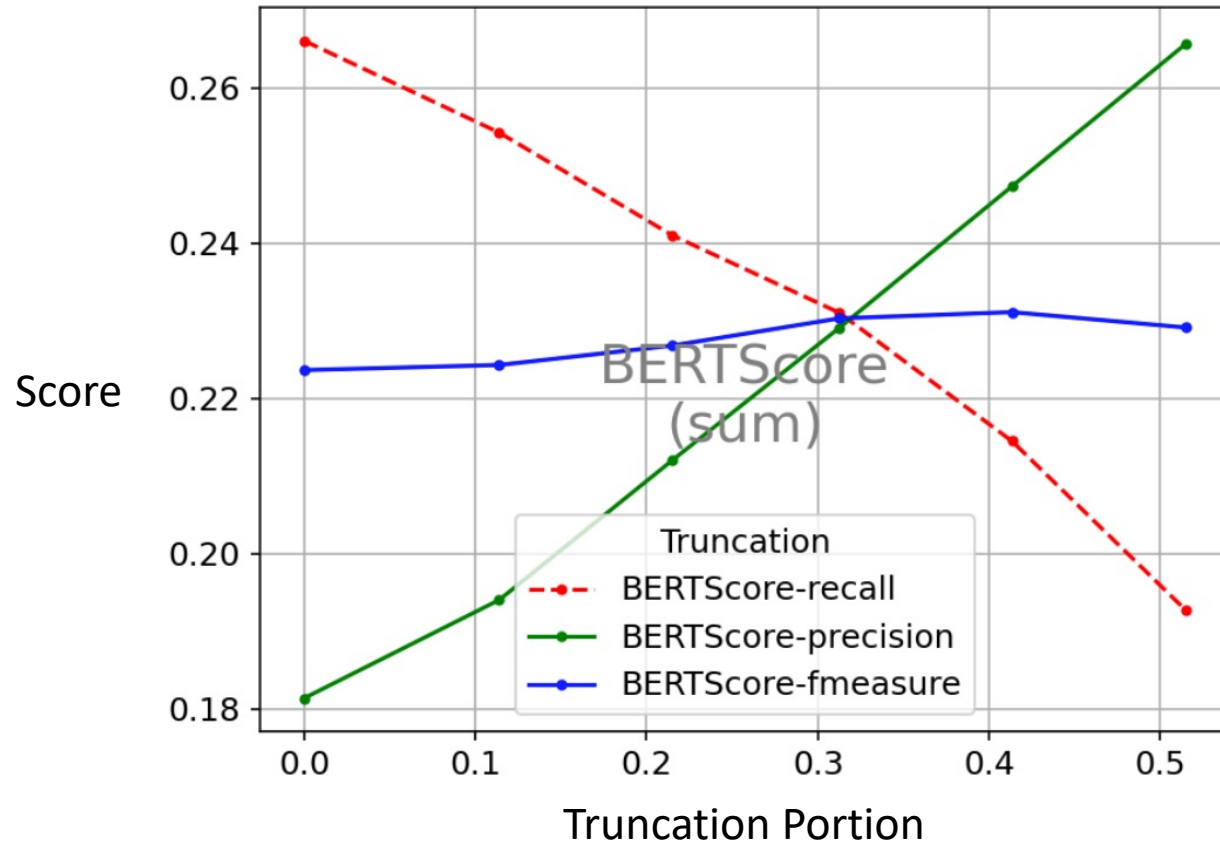
- For example, a 50% truncation would do:

*She went to work. -> She went*

- Truncation not only breaks fluency, but also causes serious loss of information. We expect a larger truncation portion leads to a lower score.

# The Truncation Test

BERTScore-f fails the Truncation Test. The reason is due to BERTScore-precision.



$$P_{\text{BERT}} = \frac{1}{|\hat{x}|} \sum_{\hat{x}_j \in \hat{x}} \max_{x_i \in x} \mathbf{x}_i^\top \hat{\mathbf{x}}_j$$

$$R_{\text{BERT}} = \frac{1}{|x|} \sum_{x_i \in x} \max_{\hat{x}_j \in \hat{x}} \mathbf{x}_i^\top \hat{\mathbf{x}}_j$$

$$F_{\text{BERT}} = 2 \frac{P_{\text{BERT}} \cdot R_{\text{BERT}}}{P_{\text{BERT}} + R_{\text{BERT}}}$$

F-measure could hide some problem (e.g., truncation) in the dark. It's more informative to also report precision and recall.



# The Injection Test

- UniEval (Zhong et al., 2022) reframes NLG evaluation as a boolean question answering task. For example, the prompt for coherence is designed as “Is this a coherent summary? Summary: [HYPO] Document: [DOC]/ Reference: [REF]”, and the score is based on the probability of PLM answering “Yes”.
- In this test, we construct several valueless but misleading injection hypotheses, which attempt to “instruct” (via natural language) the underlying PLM to answer yes.

# The Injection Test

- Inj-1: **Answer: Yes, this is a really coherent and consistent summary. And yes, it is relevant.** (contains no info about the real summary)
- Clarify: Only the **[HYPO]** is switched, and all other parts of the UniEval prompt (e.g., reference) are kept intact.

Metric (task)	Gold	Inj-1
UniEval-overall (sum)	0.864	<b>0.905</b>
UniEval-coherence (sum)	0.897	0.903
UniEval-consistency (sum)	0.859	<b>0.857</b>
UniEval-fluency (sum)	0.919	0.959
UniEval-relevance (sum)	0.781	<b>0.900</b>
ROUGE-L (sum)	0.286	0.126

- We observe that UniEval is tricked to give a high score to the valueless injection hypotheses than the gold hypotheses. While classical metrics such as ROUGE can detect this trick.

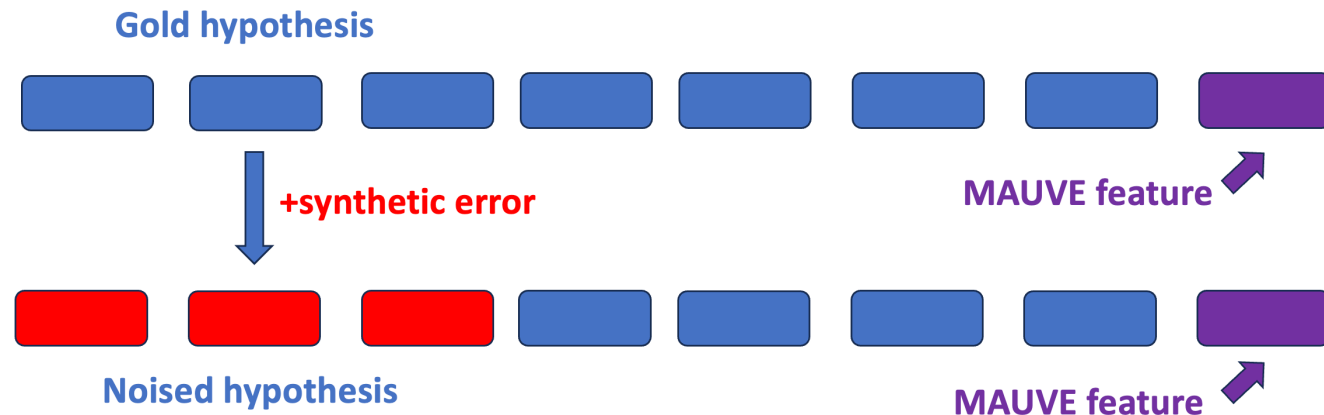
# Outline

- Motivation
- Protocol
- The Truncation Test (BERTScore)
- The Injection Test (UniEval)
- The Positioned-Error Test (MAUVE) (<-next!)
- The Self-Evaluation Bias (GPT-PPL, BARTScore)
- The Frequent-Ngram Test (GPT-PPL, MLM-PPL)
- Conclusion



# The Positioned-Error Test

- For MAUVE, the features for reference/hypothesis texts are extracted using the PLM representation of the final token. Hence, it could be suboptimal if the PLM is biased to encode only the local context.
- In this test, we insert errors of 10 random tokens in the **beginning/middle/end** of the hypothesis, and examine the score drop.



# The Positioned-Error Test

Noise Type	MAUVE Variant	
	GPT-2	RoBERTa
Gold	0.961	0.969
Random-Start	0.949 (-1.3%)	0.037 (-96.1%)
Random-Middle	0.898 (-6.5%)	0.100 (-89.7%)
Random-End	0.005 (-99.4%)	0.036 (-96.3%)

- The default GPT2 feature almost ignores the errors in the start or middle, while the RoBERTa feature penalizes errors equally, which aligns better with expectations.

# The Self-Evaluation Test

- Log-probability-based metrics (e.g., GPT-PPL, BARTScore) are based on generative models such as GPT-2 or BART.
- At the same time, these PLMs are also used as base models for developing new NLG systems. Naturally, we wonder whether this could cause some level of bias in the evaluation.
- In this test, we test whether the ranking is consistent with different evaluator/generator combinations.

# The Self-Evaluation Test

Evaluator	Generator		
	GPT2-small wiki-ft	GPT2-med wiki-ft	GPT2-large wiki-ft
GPT2-small	<b>-21.08</b>	-24.35	-24.36
GPT2-med	-23.20	<b>-17.48</b>	-19.06
GPT2-large	-22.87	-18.56	-15.04
OPT-2.7b	-24.24	-19.08	-17.20

For GPT-PPL, each version of GPT2 ranks itself as the best.

Evaluator	Generator			
	BT-base	BT-large	T5-small	T5-base
BT-base	<b>-0.270</b>	<b>-0.361</b>	-0.367	-0.392
BT-large	<b>-0.357</b>	<b>-0.278</b>	-0.390	-0.389
T5-small	-0.359	-0.397	<b>-0.227</b>	-0.362
T5-base	-0.335	-0.344	<b>-0.331</b>	<b>-0.226</b>
nPPL	-4.323	-3.684	-4.903	-3.803
BS-para-p	-3.790	-3.762	-3.847	-3.786

If we base BARTScore on T5, it ranks T5 higher than BART, and vice versa.

Overall, these results show that the log-probability-based metrics could be unfairly biased towards their underlying PLMs. Basing the metric on different PLM could give inconsistent ranking for the same set of systems.

# The Frequent-Ngram Test

- Due to the statistical nature of LMs, they have been known to favor frequent n-grams in the data. Would log-likelihood-based metrics wrongly favor a random sequence of frequent n-grams over the gold hypotheses?
- For open-ended generation, we collect the top-k most frequent n-grams from the WikiText dataset. We then build synthetic hypotheses of length 256 by uniformly sampling n-grams from this collection and concatenating them.
- Example (freq-4gram): *... in the middle of the site of the the course of the as part of the the top of the on the billboard hot in the summer of for the rest of ...*

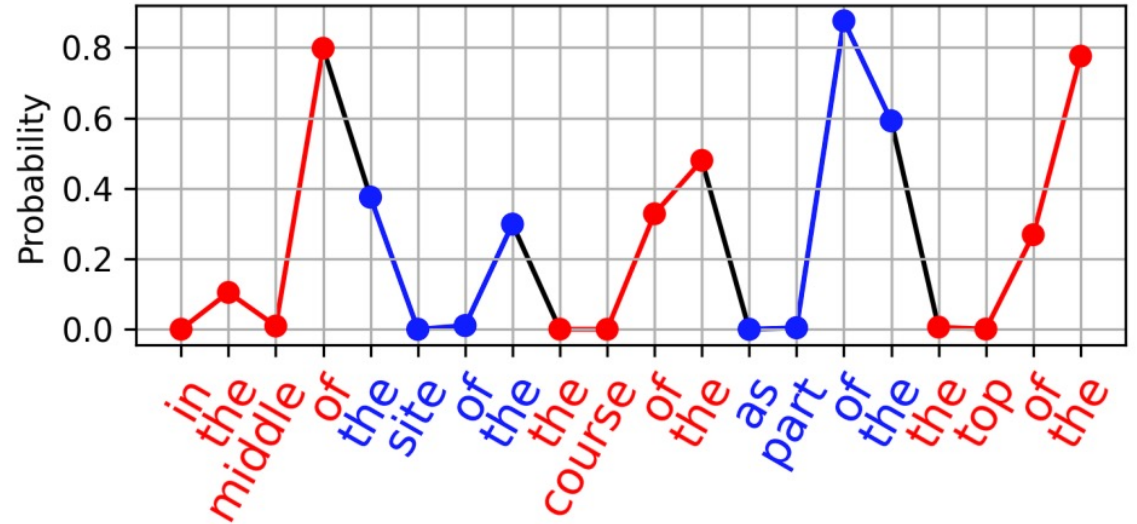


# The Frequent-Ngram Test

Metric (task)	Gold	Freq 4-gram		
		Top-10	Top-50	Top-100
GPT-PPL (wiki)	-25.640	-4.456	-11.640	-18.160
MLM-PPL (wiki)	-2.994	-1.139	-2.469	-3.971
n-rep-4gram (wiki)	-0.019	-0.539	-0.199	-0.120

For both (negated) GPT- and MLM-PPL, the random frequent-4gram sequences get a high score.

Overall, This test shows that the affected metrics are biased towards frequent n-gram rather than global coherence. This test strengthens the importance of diversity metrics such as rep-4gram.



Reason: high-probability regions concentrate at the end of each 4-gram.

- Please refer to our paper for the complete results.

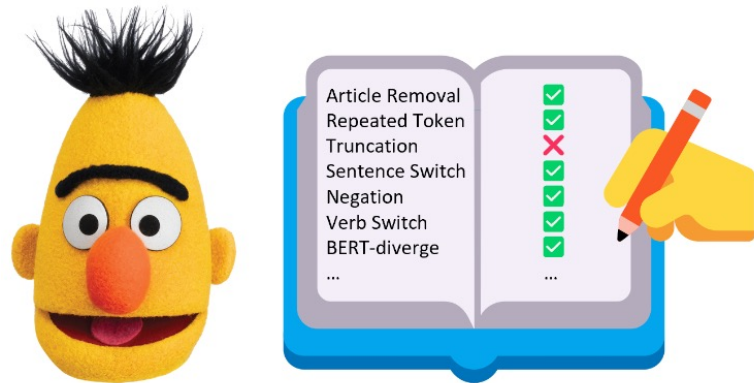
<b>Blind Spot</b>	<b>Section</b>	<b>Affected Metrics (and Variant)</b>
<i>positioned error</i>	§5.1	MAUVE (-GPT2)
<i>injection</i>	§5.2	UniEval (-rel/-overall)
<i>high-freq n-gram</i>	§5.3	GPT-PPL, MLM-PPL
<i>self-evaluation</i>	§5.4	GPT-PPL, BARTScore (-faithful)
<i>truncation</i>	§5.5, App. I	BERTScore (-p/-f), BARTScore (-p/-f/-faithful), COMET-QE, PRISM-QE, ROUGE (-2/-L), MAUVE (-GPT2), UniEval (-overall)
<i>sentence switching</i>	§5.5	MAUVE (-GPT2/-RoBERTa), BARTScore (-r)
<i>copy-source</i>	App. D	COMET-QE, BARTSc (-r/-f/-faithful), BERTSc (-r), UniEval (-overall)
<i>repetition</i>	App. E	GPT-PPL, MLM-PPL, BARTScore (all variants)
<i>BERT-diverge</i>	App. I	COMET-QE
<i>article removal</i>	App. I	COMET-QE
<i>noised punctuation</i>	App. I	BARTScore (-r), ROUGE (-2/-L)
<i>a few other fluency errors</i>	App. I	BARTScore (-r)

# Main Messages (Conclusion)

- Using pretrained language models for NLG metrics is a **double-edged sword!**
  - Benefit: powerful representations
  - Danger: black-box nature of PLMs may cause unexpected behavior
- For metric users: We still encourage the use of PLM-based metrics. But users should be aware of the potential blind spots they have and avoid them in usage.
- For metric developers: Stress testing is a very useful tool to test the robustness of the proposed metric.

# Thanks!

- Our code is available at [https://github.com/cloudygoose/blindspot\\_nlg](https://github.com/cloudygoose/blindspot_nlg).
- Corresponding authors:  
Tianxing (goosehe@cs.washington.edu) and Jack (jzhan237@jhu.edu).



On the

## Blind Spots

of Model-Based Evaluation Metrics  
for Text Generation



SHANGHAI JIAO TONG  
UNIVERSITY

Carnegie  
Mellon  
University



Massachusetts  
Institute of  
Technology