



Measuring and Improving BERT's Mathematical Abilities by Predicting the Order of Reasoning

Piotr Piękos¹, Henryk Michalewski^{1,2}, Mateusz Malinowski³

¹University of Warsaw, ²Google, ³DeepMind

bert-math.github.io

Accepted to



ACL-IJCNLP 2021



Motivation

Improving mathematical tasks performance of language models

Weak BERTs results on mathematics

“If I have 2 apples and I get 3 more apples, then I have [MASK] apples”
BERT predicts uniform distribution over one digit-number, with “3” being a slight favourite.

Models are biased in different ways

BERT relies heavily on biases and learns shortcuts instead of doing proper computations to solve math problems.

Biases

Order bias

BERT chooses the answer in a closed test based on the order of prompted possible answers. We solve that problem by using retrieval networks.

Distinctive answer bias

BERT prefers answers with round numbers or integers when other possibilities are non-integer.

Poor representation of math related sentences

BERTs embeddings poorly differentiate between basic 4 operators (+, -, *, /)

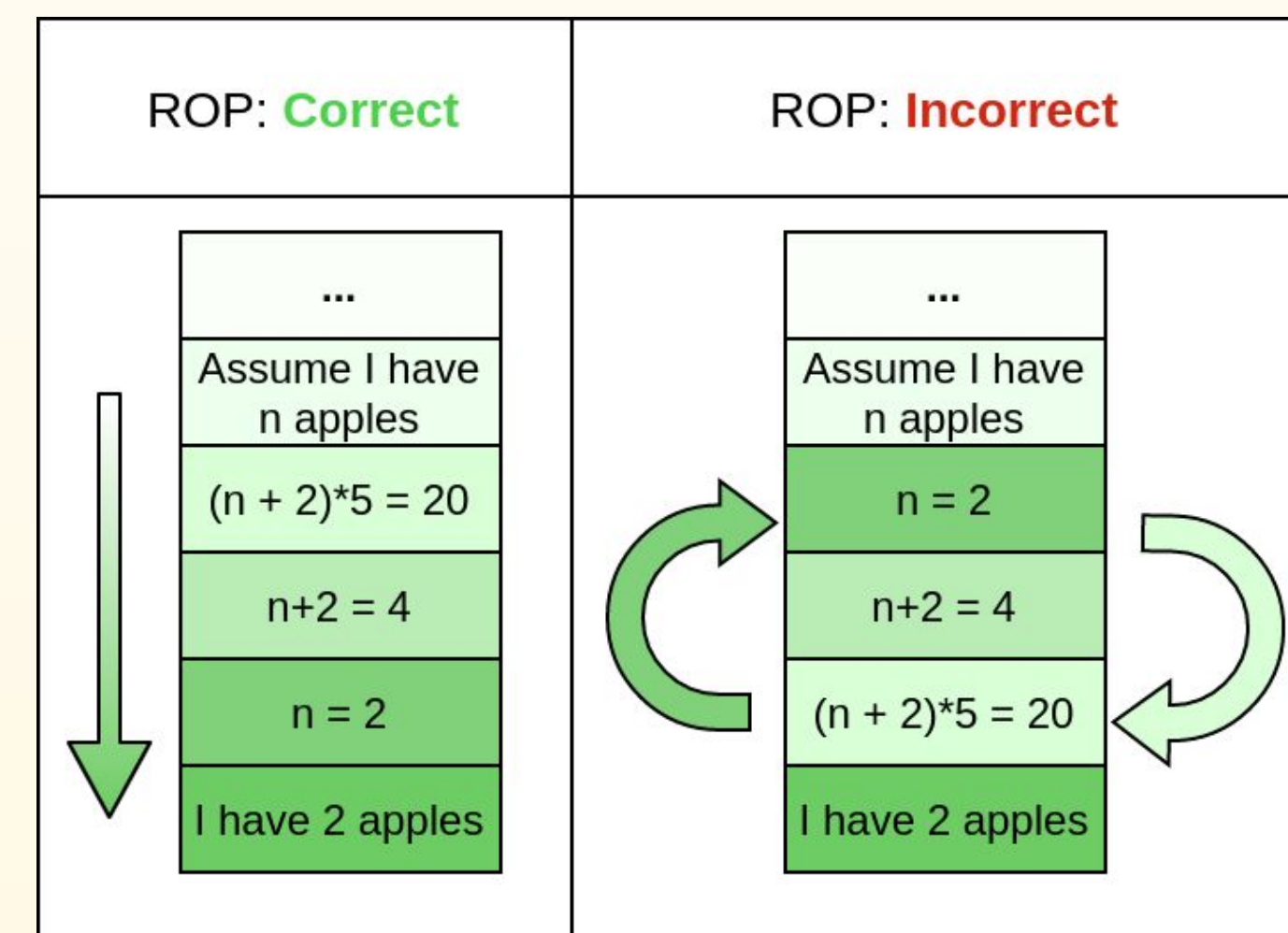
Rationales

Rationales are a step-by-step derivation of the answer. We train models on them along with mathematical questions. We hypothesize that rationales are useful for improving mathematical skills in language models because they mix natural language and formal mathematics. Hence, rationales are a bridge between them.

Training on rationales

We train on rationales because they contain a combination of natural language and mathematics, so the model can smoothly incorporate mathematics into its NLP repertoire.

```
Assume that C was there in the business for x months
A:B:C = 40000*12 : 60000*10 : 120000*x
= 40*12 : 60*10 : 120x = 40 : 5*10 : 10x
= 8 : 10 : 2x
= 4 : 5 : x
C's share = 375000*x/(9+x) = 150000
=> 375x/(9+x) = 150
=> 15x = 6(9+x)
=> 5x = 18 + 2x
=> 3x = 18
=> x = 18/3 = 6
It means C was there in the business for 6 months. Given
that B joined the business
after 2 months. Hence C joined after 4 months after B
joined
Answer is B
```



NROP

Extension of the ROP task, where only neighbor rows are swapped. This requires more subtlety from the model, therefore forces it to focus on rationales more and create richer representations.

Additional losses for better rationale utilization

In Reasoning Order Prediction (ROP) model predicts whether the order of rationale steps is correct.

For negative examples, two random steps are swapped. For positive examples the order is unchanged.

Experiments

Impact of new losses

Our methods improve the BERTs performance, even on-par with crafted models.

Measuring impact of the rationales

Training on rationales gives better results than training on more questions.

Permutation consistency tests

We investigate biases of the model by specially crafted tests.

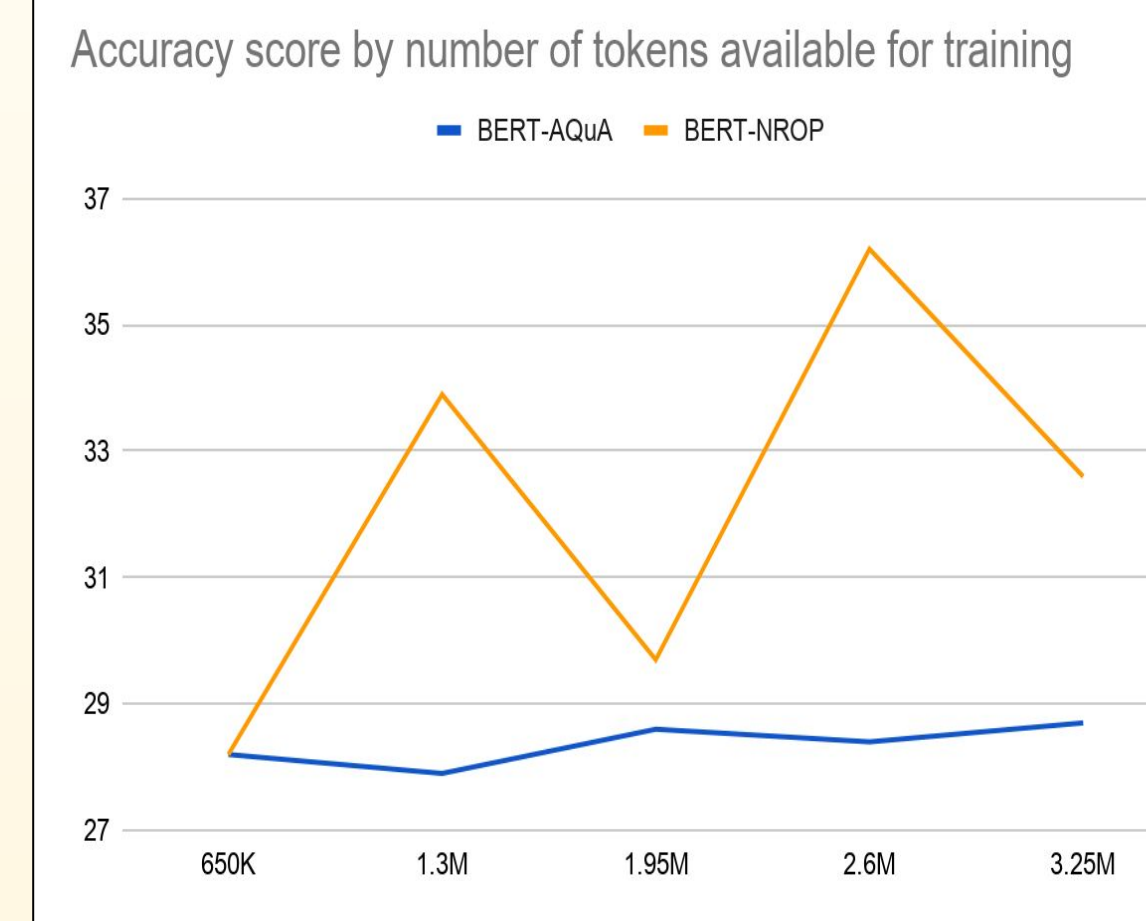
Question difficulty

We group questions by difficulty perceived by the model. A manual inspection of the clusters has shown that they form thematic groups like linear equations, etc.

Human study

Human study shows similarities between difficulty levels “perceived” by our models and humans.

Model	Accuracy
Random chance	20%
LSTM	20.8%
BERT-base	28.3%
BERT-NROP	37.0%
AQuA-RAT	36.4%
MathQA	37.9%



Conclusions

We showed that explanations are useful for learning better representations. Additionally, we proposed novel tasks and losses for utilizing rationales. We showed that they significantly improve mathematical reasoning in BERT. We showed that BERTs representation is biased in mathematics. In the end, we proposed permutation invariant losses for reducing that bias.