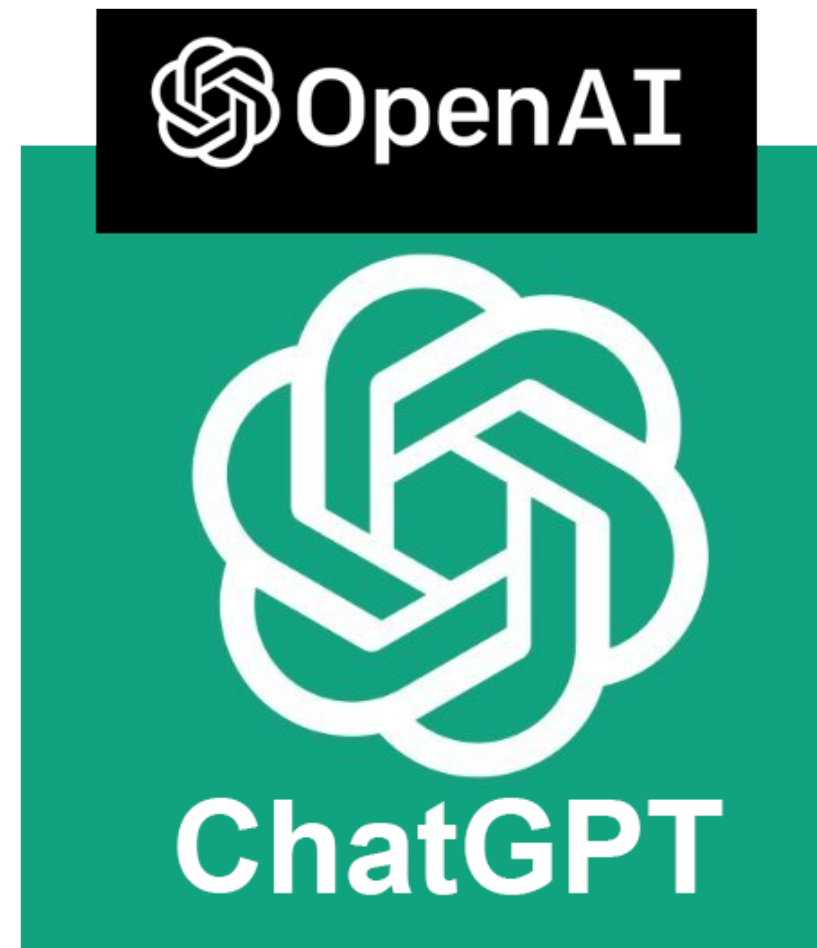




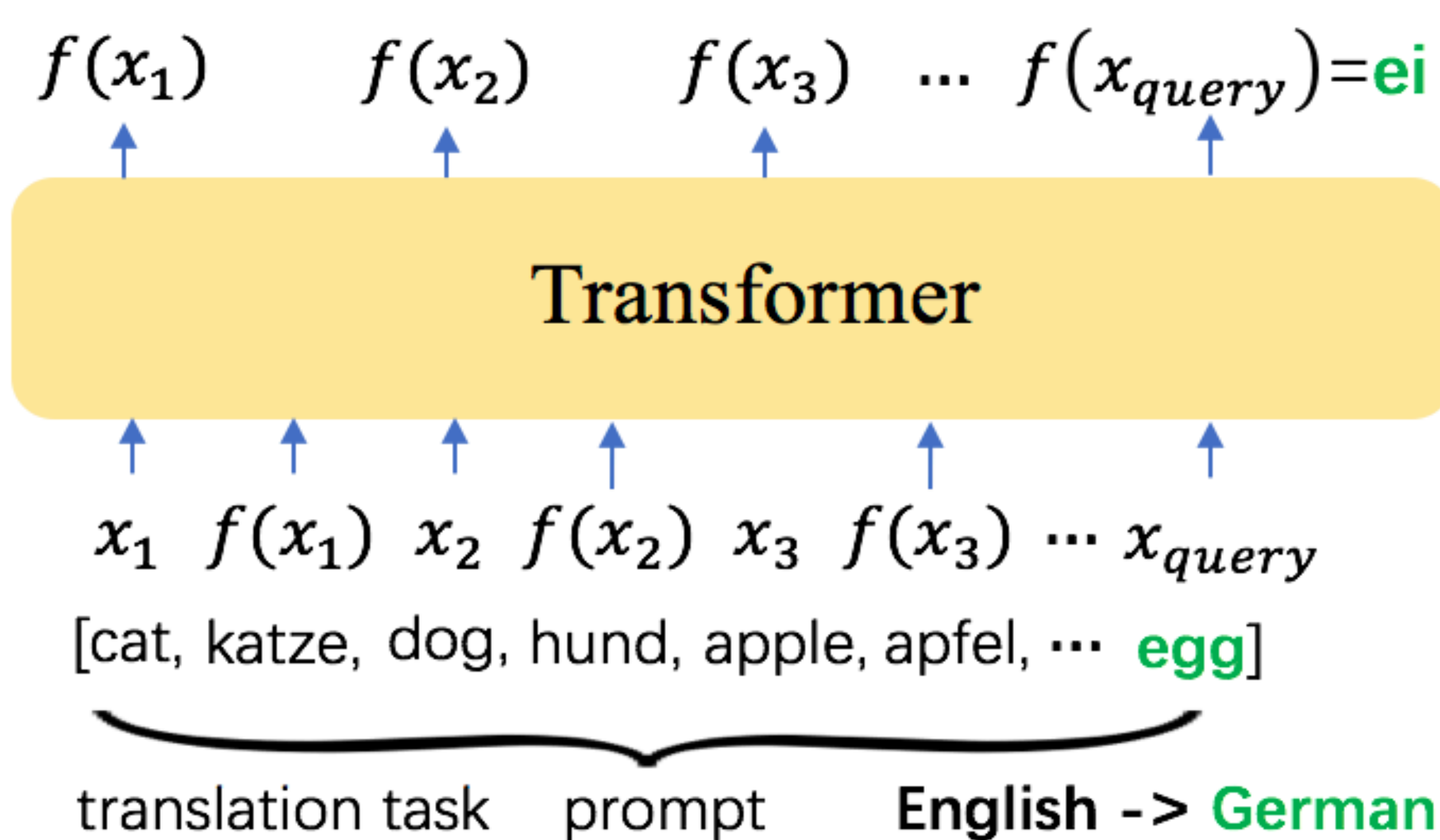
Hongkang Li<sup>1</sup>, Meng Wang<sup>1</sup>, Songtao Lu<sup>1</sup>, Xiaodong Cui<sup>1</sup>, Pin-Yu Chen<sup>1</sup>. 1: Rensselaer Polytechnic Institute. 2: IBM Research



## Motivation

Transformer-based foundation models, e.g., GPT-4, Sora, have achieved great empirical success in many areas.

- Large foundation models are able to implement in-context learning (ICL) and reasoning.
- Theoretical understanding of **how a Transformer can be trained to perform ICL and generalize in and out of domain successfully and efficiently** is less investigated.

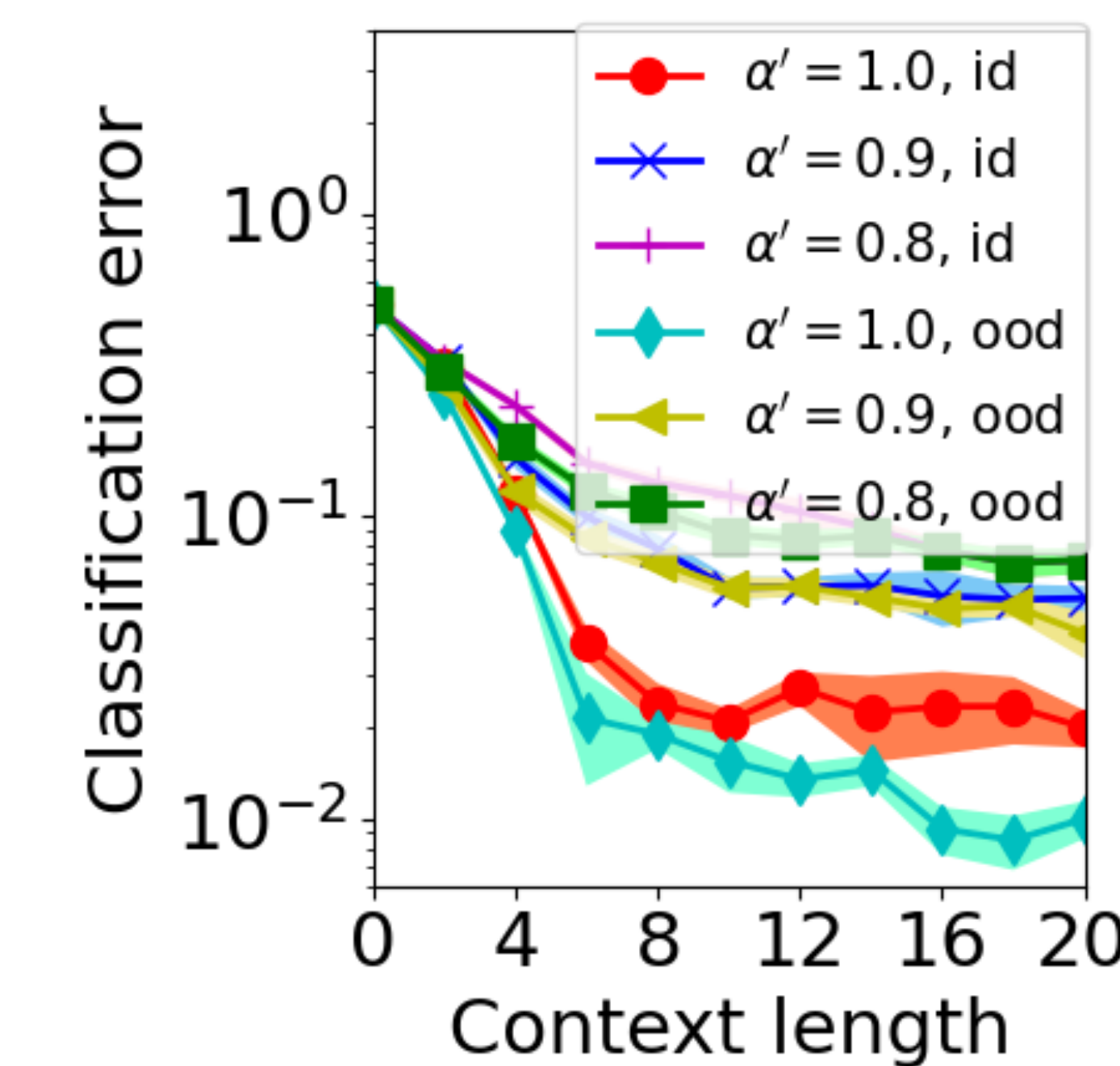


## Current Progress

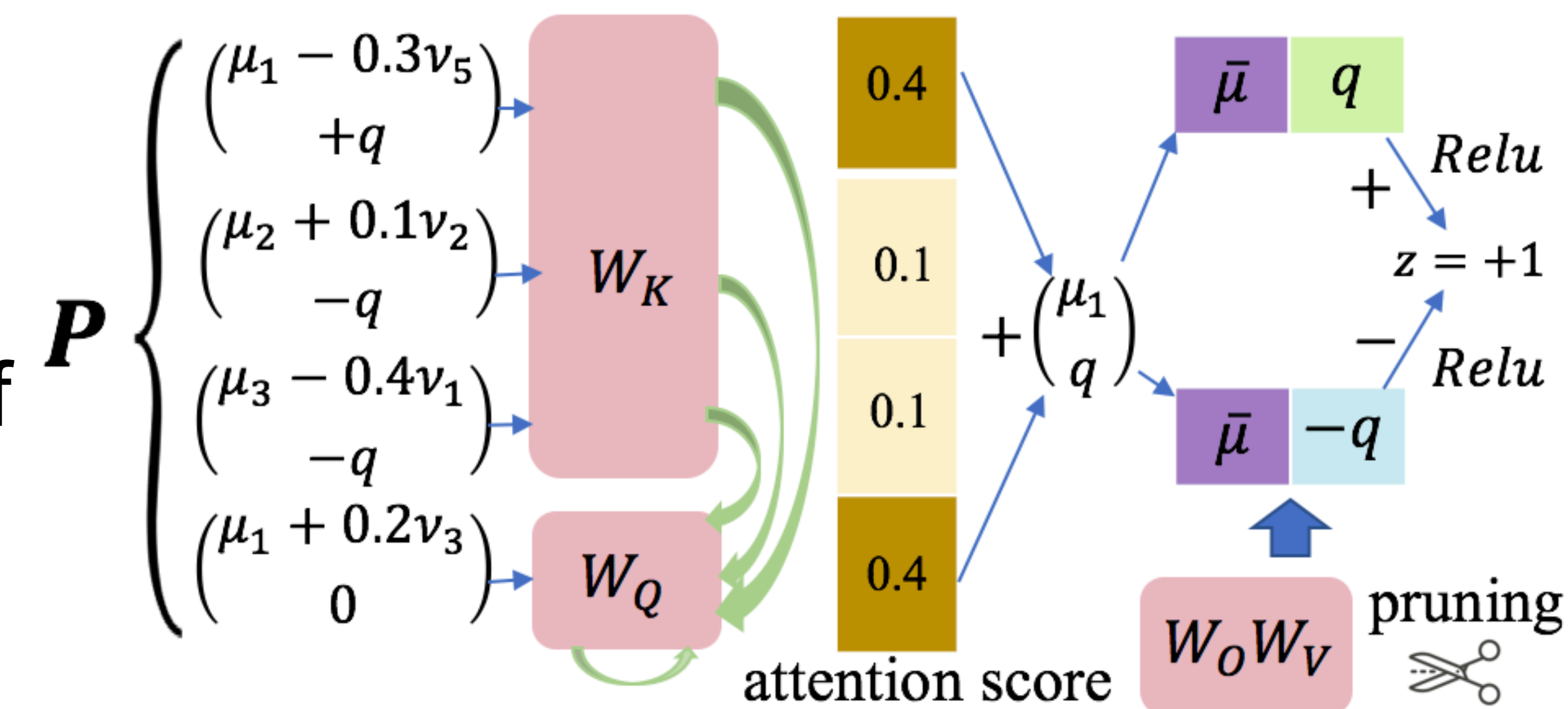
- We provide a theoretical characterization of how to train nonlinear Transformers to enhance their ICL capability on classification tasks. .

*Theorem 1 (informal): Given enough neurons and a large batch, and prompt lengths inverse in the fraction of relevant tokens  $\alpha$ , then after training with  $\Theta(\alpha^{-1})$  steps,*

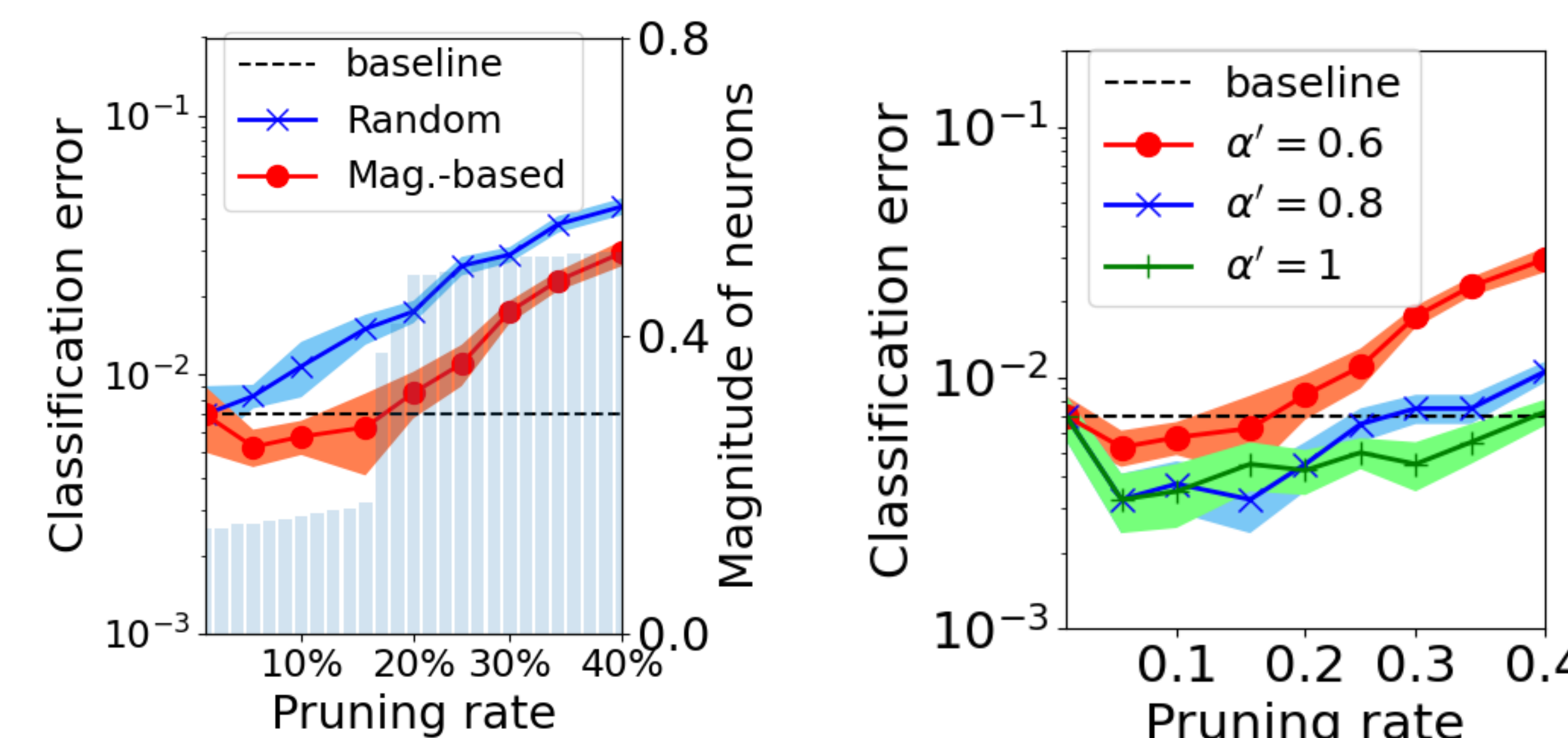
- the returned one-layer Transformer model achieves an in-domain generalization error no larger than  $\epsilon$ .
- If the testing relevant patterns are linear combinations of the trained ones with coefficient summation no larger than 1, the out-of-domain generalization error is no larger than  $\epsilon$ .



- We expand the theoretical understanding of the mechanism of the ICL capability of Transformers.

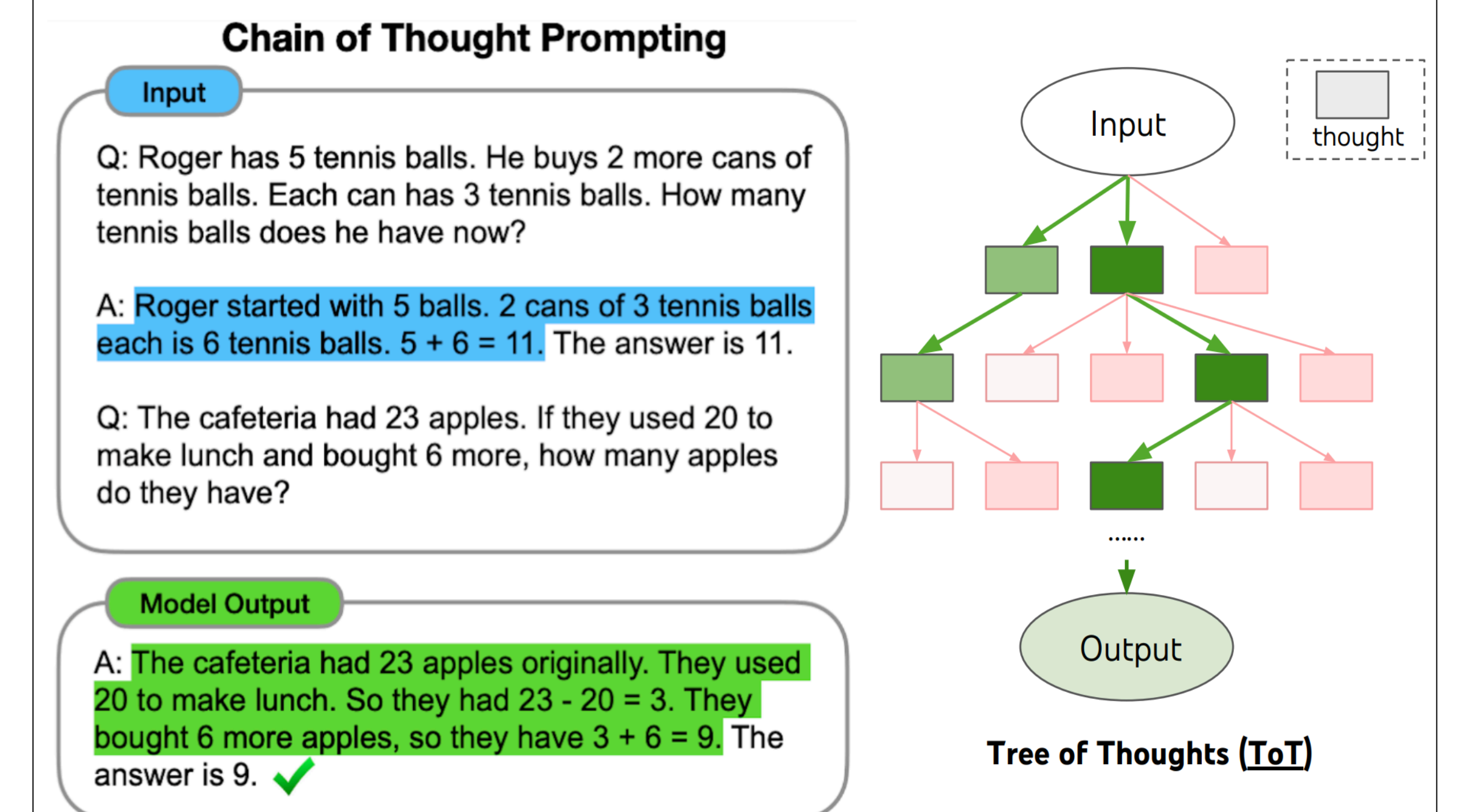


- We theoretically justify the Magnitude-based Pruning in preserving ICL.



## Future Plan

### LLM reasoning



### Problems to solve

- How can a Transformer be trained to learn different hidden causal structure?
- Why does adding intermediate steps help the reasoning in theory?
- What is the mechanism of a Transformer implementing reasoning in context?

### Theoretical contributions

- Hidden Markov chain modeling.
- Next token prediction beyond classification and regression.

### Experiments

- Evaluate the results on the arithmetic reasoning dataset GSM8K and the commonsense reasoning dataset CSQA.







## Math Word Problems (free response)

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls.  $5 + 6 = 11$ . The answer is 11.



Education Assistant

## StrategyQA

Q: Yes or no: Would a pear sink in water?

A: The density of a pear is about  $0.6 \text{ g/cm}^3$ , which is less than water. Thus, a pear would float. So the answer is no.



Smart Navigation