

# Theories of Pre-Trained Language Models with Practical Implications

Xinyi Wang (Proposal)

Committee: William Wang, Kun Zhang, Shiyu Chang, Xifeng Yan

# Outline

- Background: the many ways of understanding large language models (LLMs)
  - Interpretability v.s. theory
  - Existing LLM theories
- My progress
  - A latent variable theory
  - A data composition theory
- Future research

**Background**

# The many ways of understanding large language models (LLMs)

## Interpretability

### Step 1: Generate explanation using GPT-4

The Avengers to the big screen, Joss Whedon has returned to reunite Marvel's gang of superheroes for their toughest challenge yet. Avengers: Age of Ultron pits the titular heroes against a sentient artificial intelligence, and smart money says that it could soar at the box office to be the highest-grossing film of the

introduction into the Marvel cinematic universe, it's possible, though Marvel Studios boss Kevin Feige told Entertainment Weekly that, "Tony is earthbound and facing earthbound villains. You will not find magic power rings firing ice and flame beams." Spoilsport! But he does hint that they have some use... STARK T

, which means this Nightwing movie is probably not about the guy who used to own that suit. So, unless new director Matt Reeves' The Batman is going to dig into some of this backstory or introduce the Dick Grayson character in his movie, the Nightwing movie is going to have a lot of work to do explaining

of Avengers who weren't in the movie and also Thor try to fight the infinitely powerful Magic Space Fire Bird. It ends up being completely pointless, an embarrassing loss, and I'm pretty sure Thor accidentally destroys a planet. That's right. In an effort to save Earth, one of the heroes inadvertently blows up an

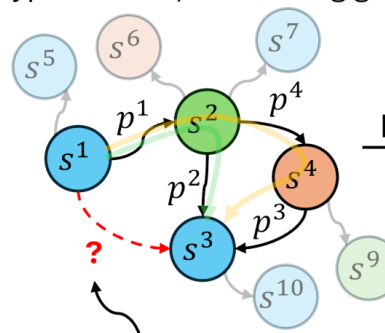
Given a GPT-2 neuron, generate an explanation of its behavior by showing relevant text sequences and activations to GPT-4.

**Model-generated explanation:**  
references to movies, characters, and entertainment.

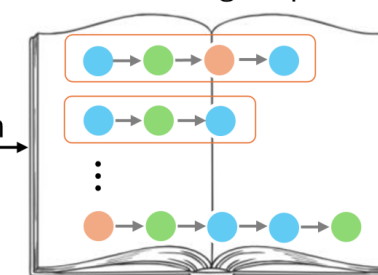
(Source: OpenAI 2023)

## Theory

### (Hypothetical) Reasoning graph G



### Pre-training corpus D



Random walk

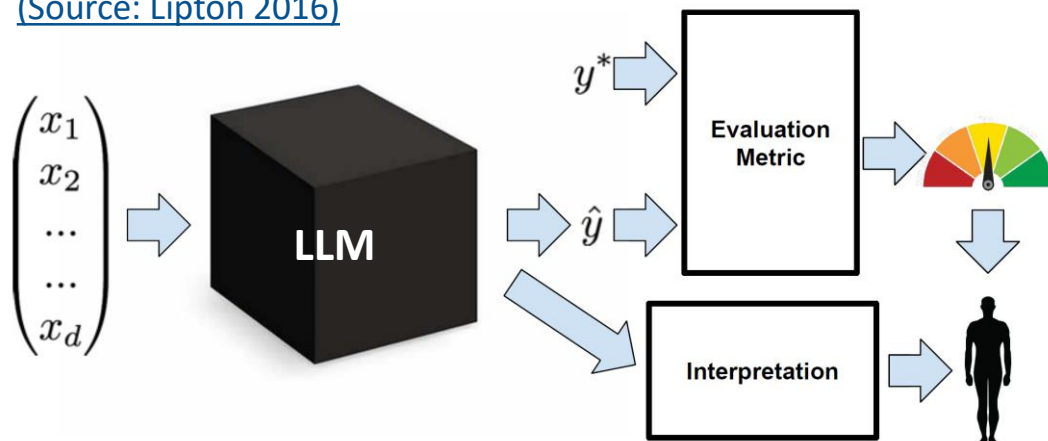
$$P_{LM}(s^1 \xrightarrow{p^i} s^3) \propto \exp[w_i^1 P_D(s^1 \xrightarrow{p^1} s^2 \xrightarrow{p^2} s^3) + w_i^2 P_D(s^1 \xrightarrow{p^1} s^2 \xrightarrow{p^4} s^4 \xrightarrow{p^3} s^3)]$$

(Source: Ours 2024)

# Interpretability v.s. Theory

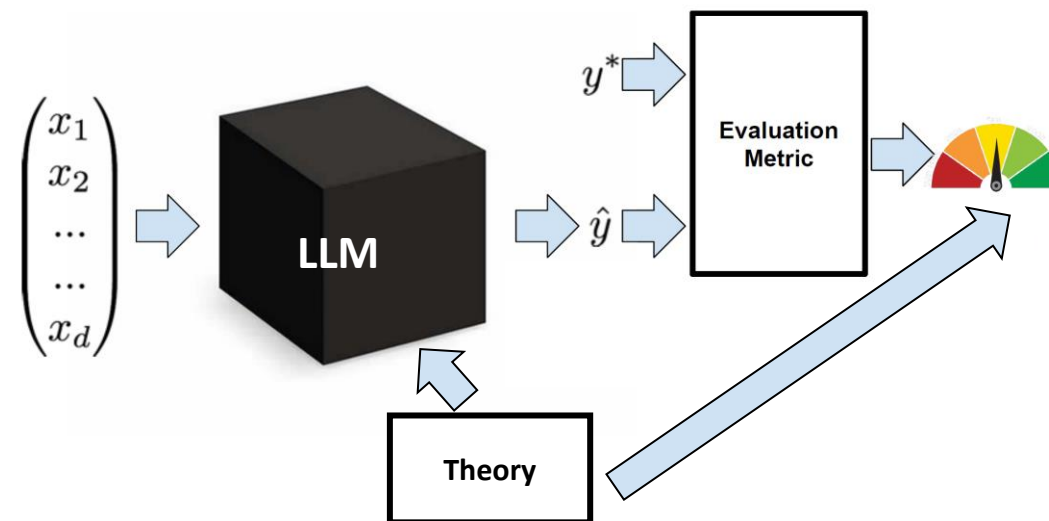
## Interpretability

(Source: Lipton 2016)



⊃

## Theory



### Goal

Make the LM prediction more transparent to humans, thus assist the final decision making progress or improve user experience.\*

Propose a self-consistent theory to explain the LM behavior/learning process as it is, which can be applied to improve the LM performance.

### Difference

Not necessarily corresponding to the underlying mechanism of LM learning/inference.

Must revealing the underlying mechanism of LM learning/inference.

### Common

Make the LM behavior more predictable. Reduce the risk of unwanted behavior of LMs.

\* [Mechanistic interpretability](#) focus on revert LM to human understandable programs. (detailed in next slides.)

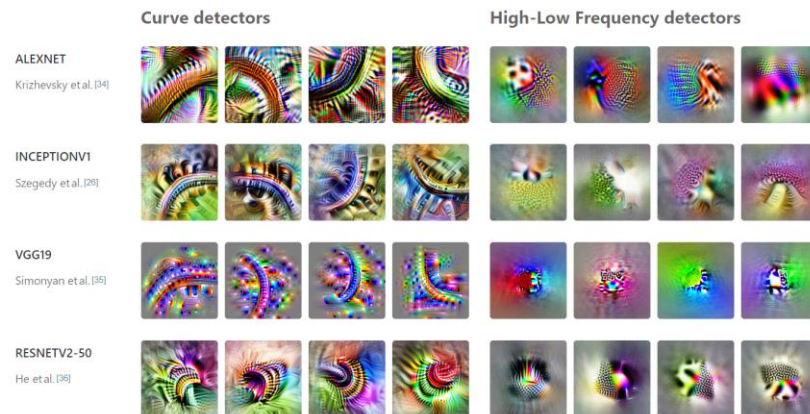
# Zoom in or zoom out?

**Universality:** Analogous features and circuits form across models and tasks.

**Circuit:** subgraph of connected neurons.

**Feature:** neuron.

**Zoom in:** the *circuit* view from mechanistic interpretability.



A car detector (4c:447) is assembled from earlier units.

**Windows** (4b:237) excite the car detector at the top and inhibit at the bottom.



**Car Body** (4b:491) excites the car detector, especially at the bottom.



**Wheels** (4b:373) excite the car detector at the bottom and inhibit at the top.



**Theory:** universal principle governing all models.

**Mechanism:** how theoretical principles are implemented as algorithm.

**Parameter:** pinpoint model parameters corresponding to a mechanism.



**Zoom out:** verify a hypothesis via theoretical analysis/experiments.

(Source: Olah et al. 2020)

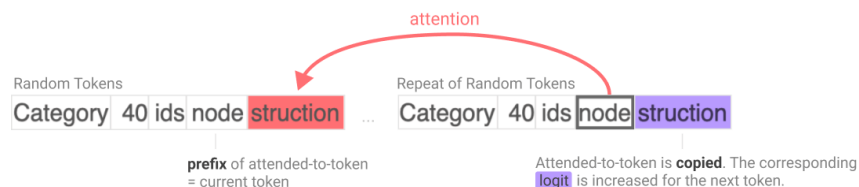
# Example In-Context Learning Theories

## Zoom in: Induction head

(Source: Olssen et al. 2022)

**Universality:** Induction heads might constitute the mechanism for the actual majority of all in-context learning (ICL) in large transformer models.

**Circuit:** Induction heads “complete the pattern” by copying and completing sequences that have occurred before.



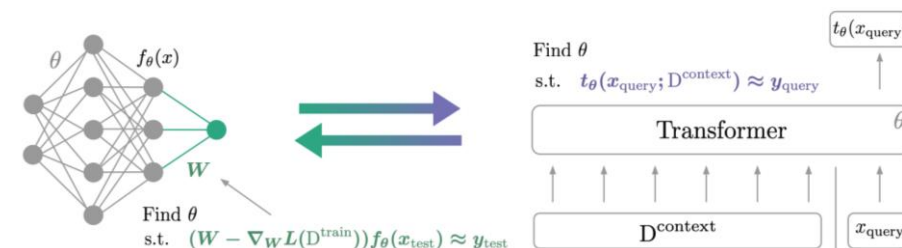
**Feature:** An attention head *copies information* from the previous token into each token, which enables the *induction head* to attend to tokens based on what happened before them, rather than their own content.

**In-context learning score:** the loss of the 500th token in the context minus the loss of the 50th token in the context.

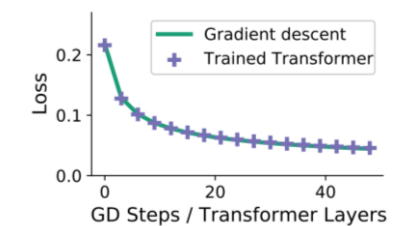
## Zoom out: Gradient descent

(Source: Akyurek et al. 2022)

**Theory:** In-context learning (ICL) is implemented by *gradient descent* on given demonstrations in Transformers.



**Mechanism:** An ordinary gradient-based *least squares algorithm* is implemented for the linear regression task.



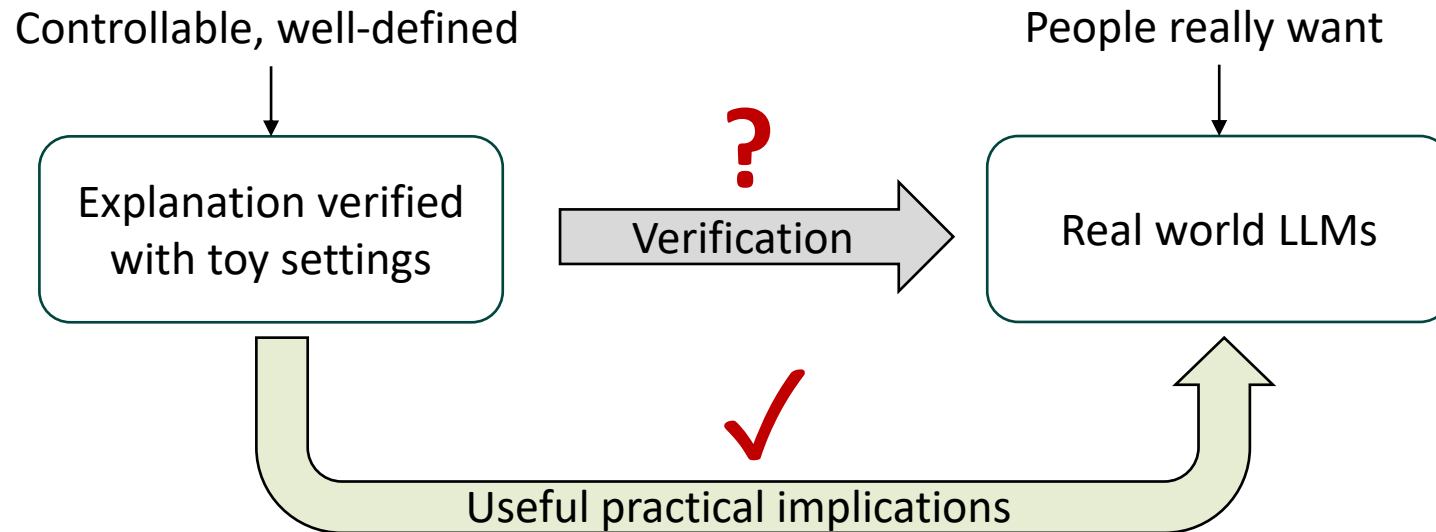
**Parameter:** There is a Transformer construction to exactly implement least square algorithm. It is unknown how it is *actually* implemented in a Transformer.

**In-context learning task:** linear regression

**Artificial!**

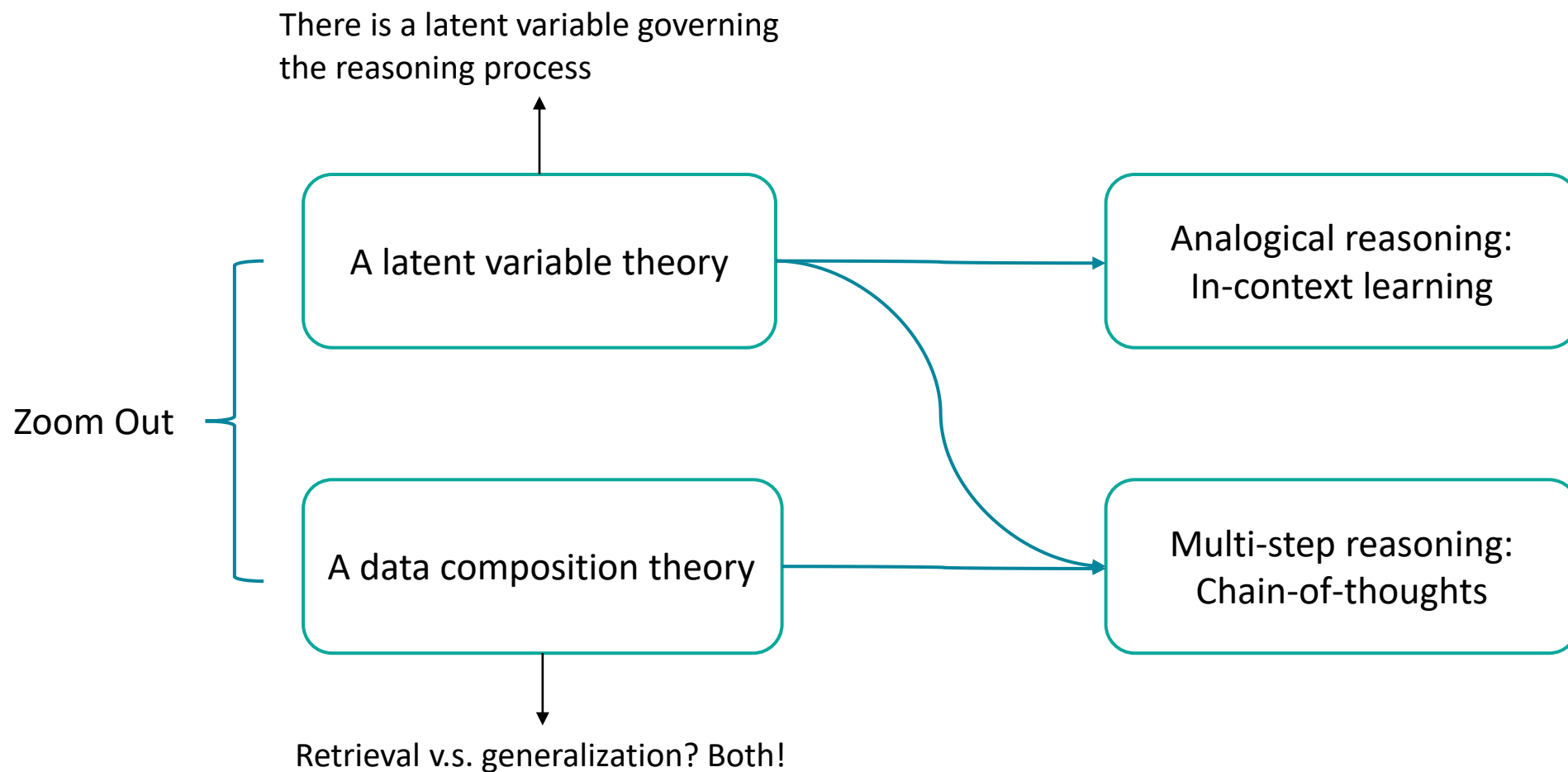
# The Common Issue

*The big gap between real world LLMs and the proposed explanations.*



**My progress**

# Our Proposed LLM Theories



UC SANTA BARBARA

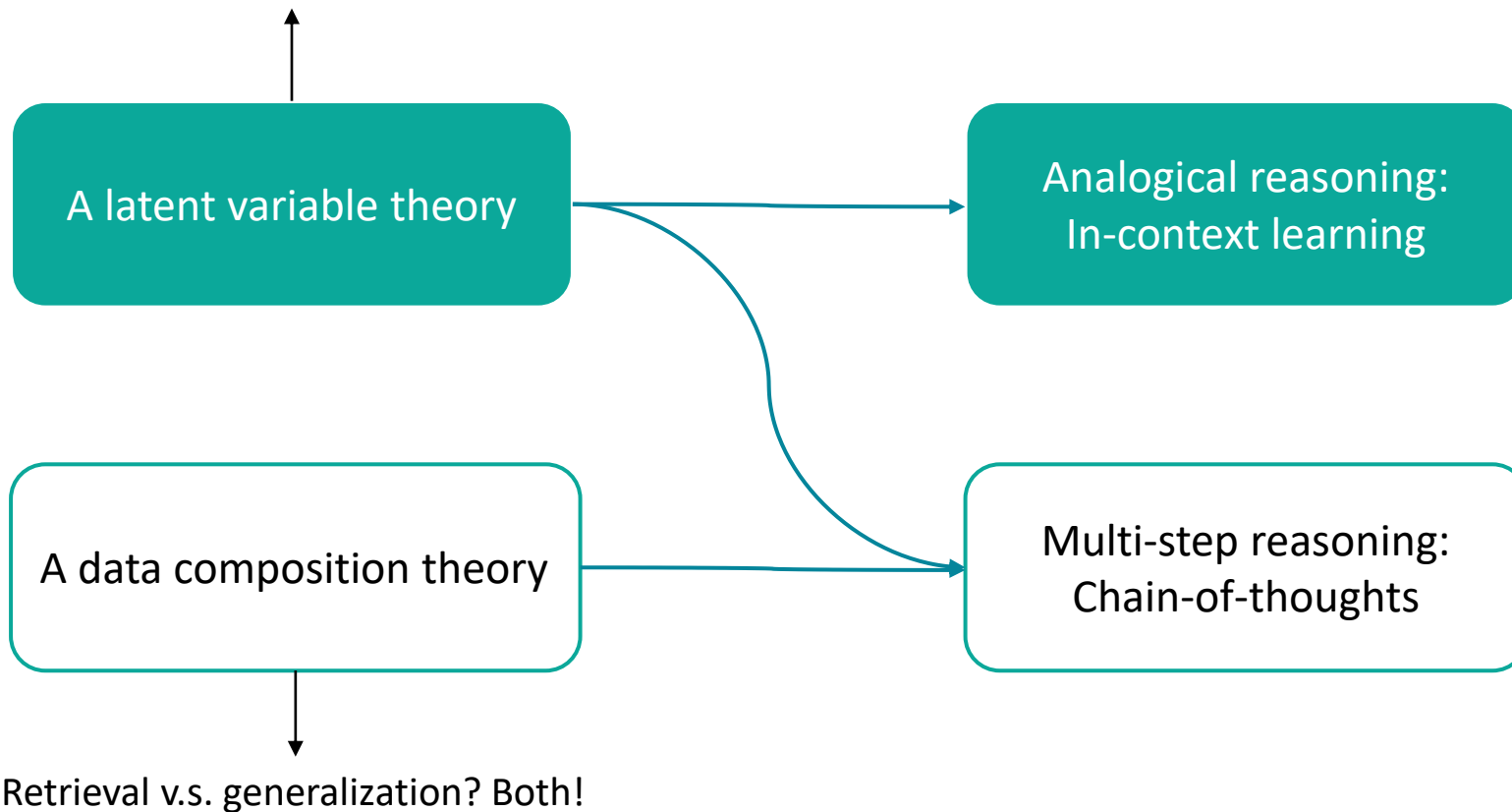
University of California, Irvine

# Large Language Models are Latent Variable Models: Explaining and Finding Good Demonstrations for In- Context Learning

Xinyi Wang, Wanrong Zhu, Michael Saxon, Mark Steyvers, William Yang  
Wang ([NeurIPS 2023](#))

# Our Proposed LLM Theories

There is a latent variable governing the reasoning process.



# LLMs are latent variabel models

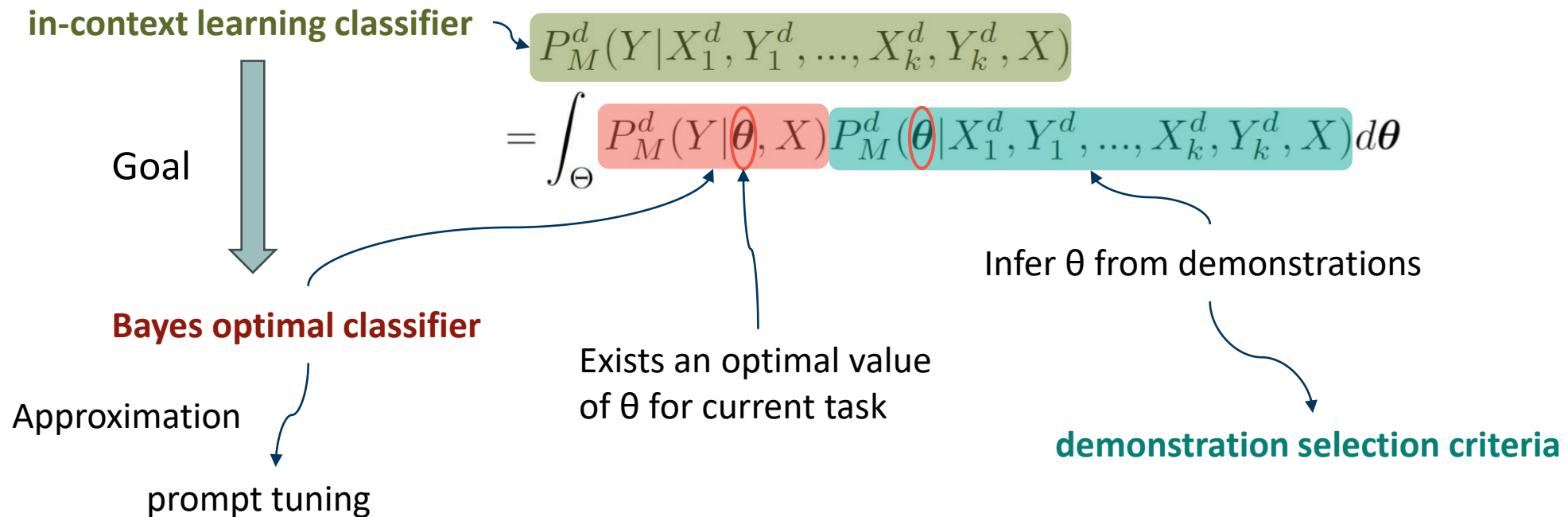
**LLM:**  $P(w_{1:T}) = \prod_{i=1}^T P(w_i | w_{i-1}, \dots, w_1)$     **Latent variabel model:**  $P(w_{1:T}) = \int_{\Theta} P(w_{1:T} | \theta) P(\theta) d\theta$

**Our assumption:**  $P_M(w_{t+1:T} | w_{1:t}) = \int_{\Theta} P_M(w_{t+1:T} | \theta) P_M(\theta | w_{1:t}) d\theta$

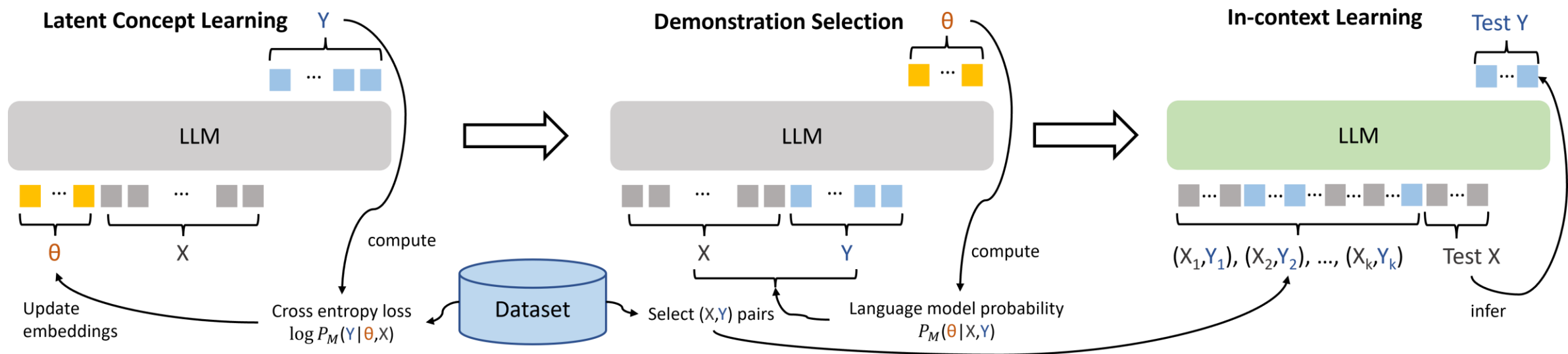
Language model probability output by an LLM

Generated continuation    Prompt    2. Generate the continuation exclusively based on the inferred concept variable  $\theta$     1. Implicitly infer a latent concept variable  $\theta$  from the prompt

# Analysis in-context learning classifier



# Algorithm overview

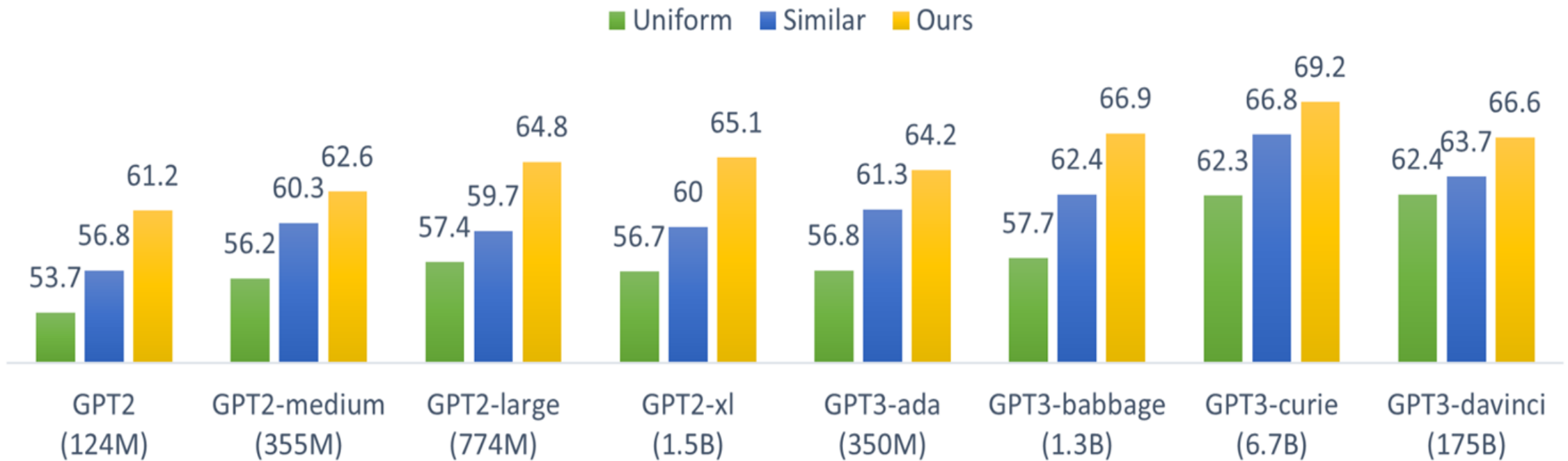


prompt tuning to obtain  
the latent  $\theta$

Scoring examples by  
 $P_M^d(\theta|X_1^d, Y_1^d, \dots, X_k^d, Y_k^d, X)$

Use the top examples for in-  
context learning with any LMs

# Text classification results



- Results are averaged over 8 text classification datasets, each experiment is repeated by 5 runs.
- We select the optimal demonstrations by GPT2-large, and use the same set of demonstrations for all other LLMs.

# GSM8K results

|                         | Uniform | Similar | Ours w/ Llama 2 (7B) | Ours w/ GPT2-XL (1.5B) |
|-------------------------|---------|---------|----------------------|------------------------|
| Prompt tuning           | -       | -       | 15.2                 | 7.3                    |
| Llama 2 (7B)            | 11.4    | 13.1    | 19.3                 | 15.9                   |
| Llama 2 (13B)           | 17.0    | 18.3    | 21.6                 | 20.5                   |
| Llama 2 (70B)           | 50.2    | 53.5    | 54.3                 | 52.9                   |
| ChatGPT (gpt-3.5-turbo) | 76.5    | 78.1    | 81.2                 | 80.4                   |

Table 1: Prompt tuning and 4-shot in-context learning accuracy on a subset of GSM8K test set. Our demonstrations are selected with either 7B Llama 2 or GPT2-XL

- We select the optimal demonstrations by Llama 2 (7B)/ GPT2-XL, and use the same set of demonstrations for all other LLMs.



UC SANTA BARBARA



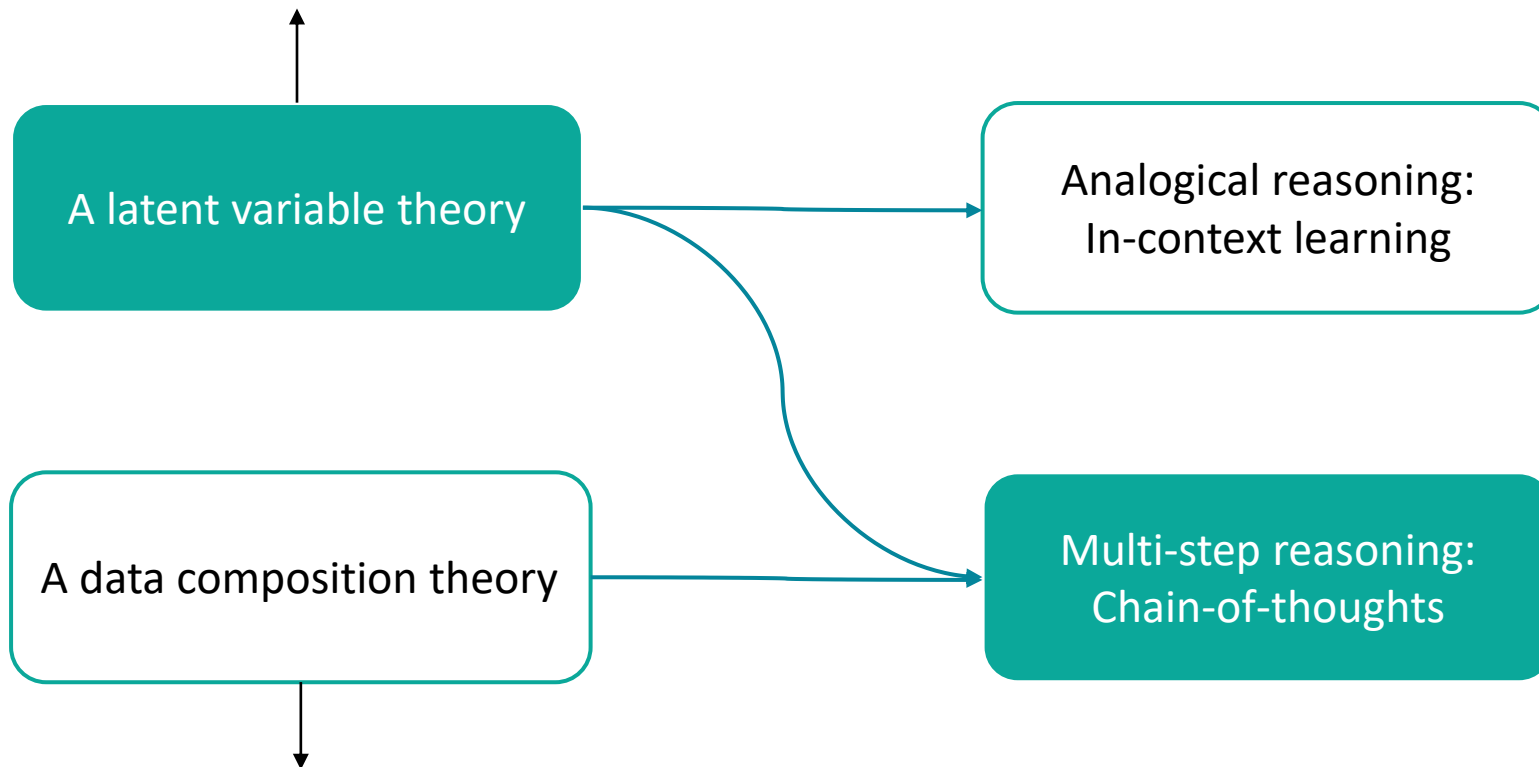
Mila

# Guiding Language Model Math Reasoning with Planning Tokens

Xinyi Wang, Lucas Caccia, Oleksiy Ostapenko, Xingdi Yuan, Alessandro Sordoni ([Arxiv](#))

# Our Proposed LLM Theories

There is a latent variable governing the reasoning process.



Retrieval v.s. generalization? Both!

# LLM fine-tuned with chain-of-thoughts data

Question: Every day, Wendi feeds each of her chickens three cups of mixed chicken feed, containing seeds, mealworms and vegetables to help keep them healthy. She gives the chickens their feed in three separate meals. In the morning, she gives her flock of chickens 15 cups of feed. In the afternoon, she gives her chickens another 25 cups of feed. **How many cups of feed does she need to give her chickens in the final meal of the day** if the size of Wendi's flock is 20 chickens?

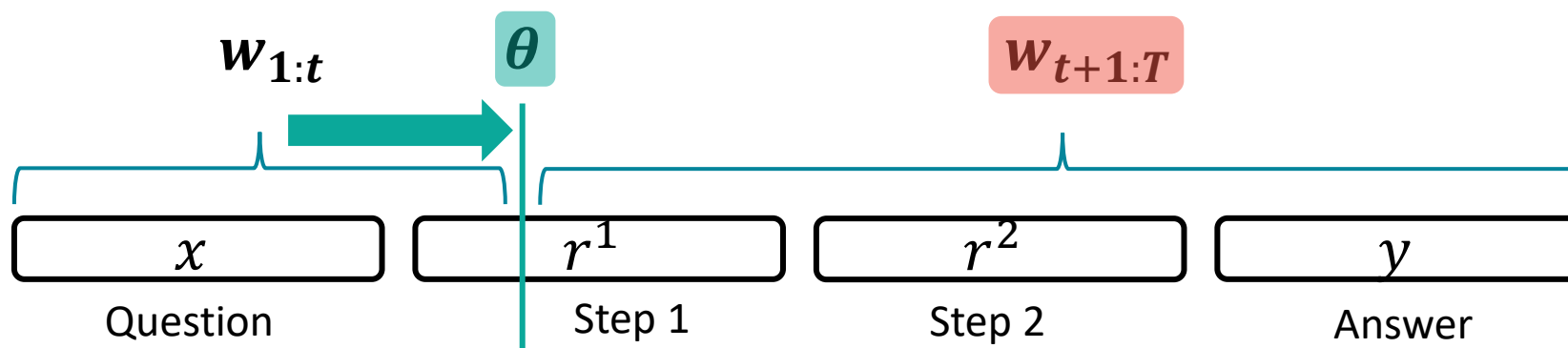
Every day, Wendi gives her chickens 15 cups of feed in the morning + 25 cups of feed in the afternoon =  $\langle\langle 15+25=40 \rangle\rangle$  40 cups of feed.

If she has 20 chickens and she feeds them 40 cups of feed every day, then each chicken gets  $40/20 = \langle\langle 40/20=2 \rangle\rangle$  2 cups of feed per chicken.

The answer is: 2

# A Bayesian view of chain-of-thoughts (CoTs)

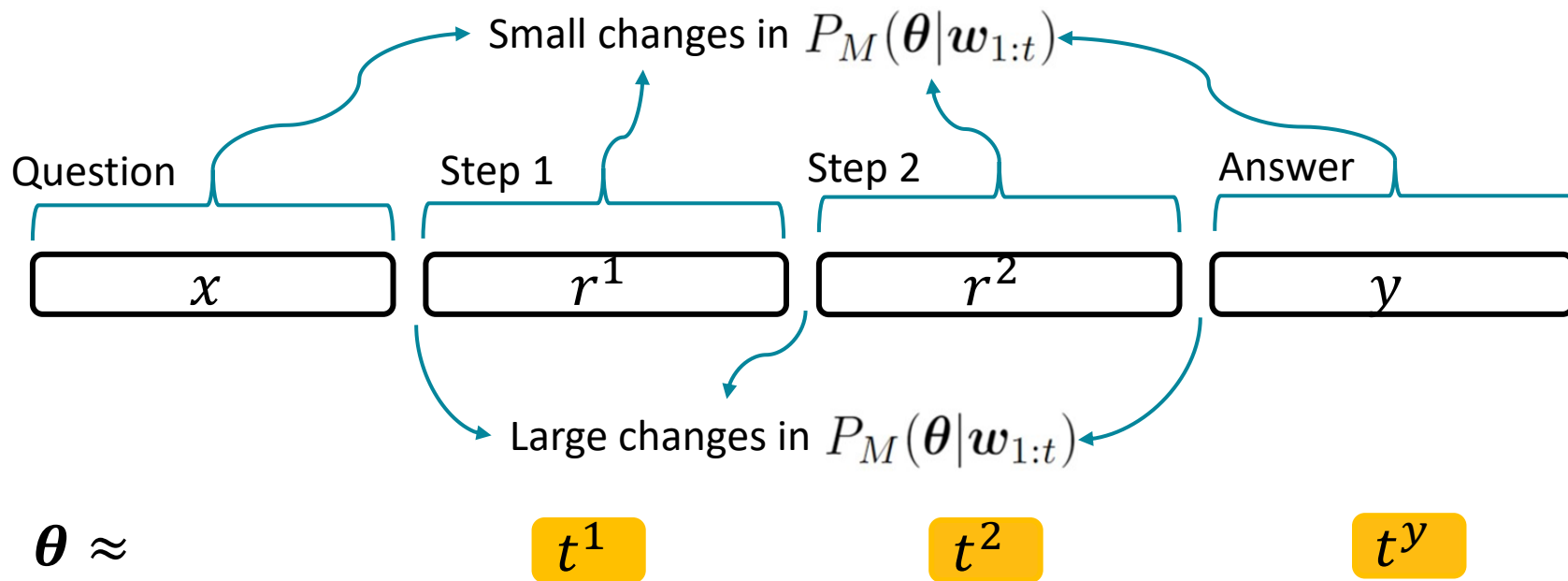
$$P_M(w_{t+1:T}|w_{1:t}) = \int_{\Theta} P_M(w_{t+1:T}|\theta) P_M(\theta|w_{1:t}) d\theta$$



**Bayesian assumption:** there is a latent variable  $\theta$  governing the generation of the whole CoT sequence.

# A Bayesian view of chain-of-thoughts (CoTs)

$$P_M(\mathbf{w}_{t+1:T}|\mathbf{w}_{1:t}) = \int_{\Theta} P_M(\mathbf{w}_{t+1:T}|\boldsymbol{\theta})P_M(\boldsymbol{\theta}|\mathbf{w}_{1:t})d\boldsymbol{\theta}$$



**Simplified Bayesian assumption:** there is a discrete planning variable  $t$  governing the generation of each chain-of-thoughts step.

# LLM fine-tuning with planning tokens

Question: Every day, Wendi feeds each of her chickens three cups of mixed chicken feed, containing seeds, mealworms and vegetables to help keep them healthy. She gives the chickens their feed in three separate meals. In the morning, she gives her flock of chickens 15 cups of feed. In the afternoon, she gives her chickens another 25 cups of feed. **How many cups of feed does she need to give her chickens in the final meal of the day** if the size of Wendi's flock is 20 chickens?

General  
planning  
tokens

Specialized planning tokens

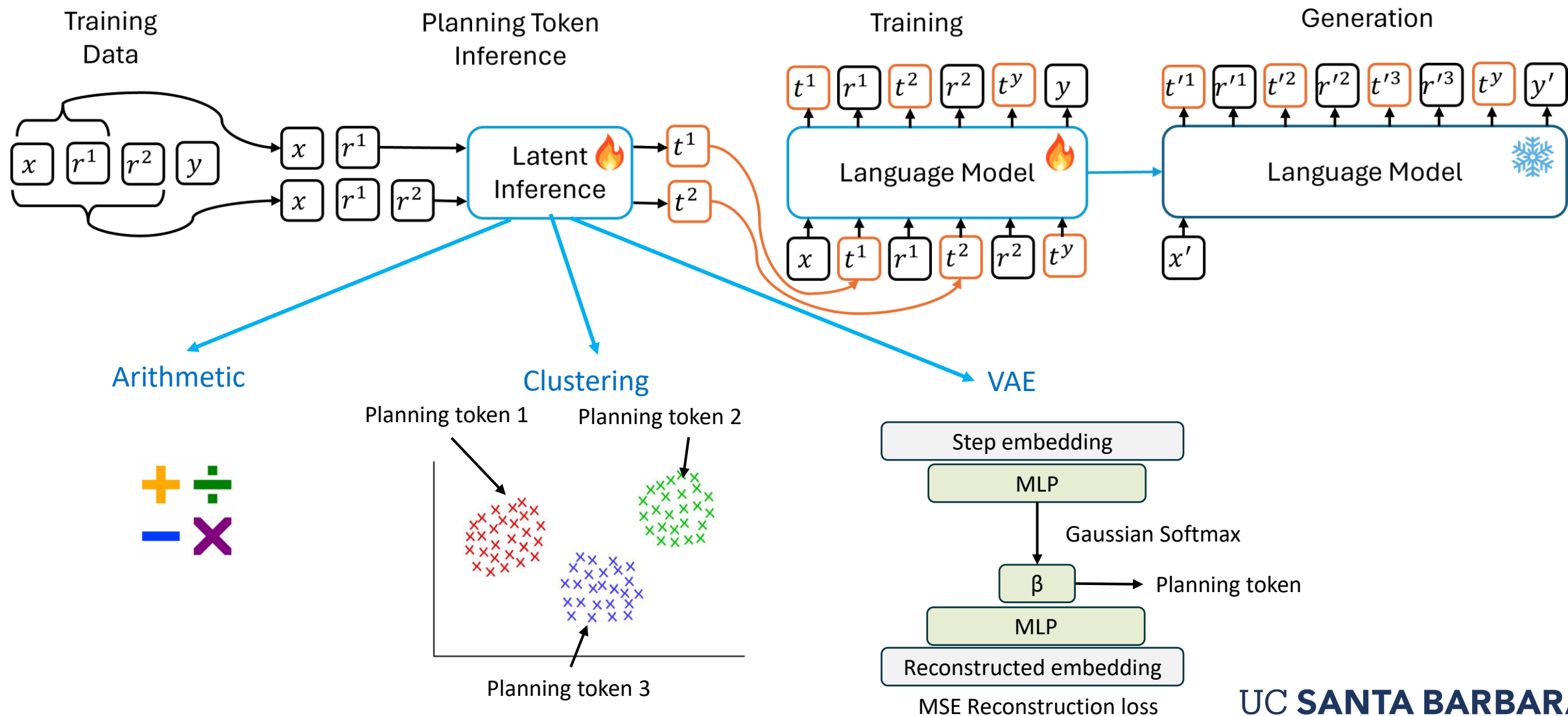
`<prefix>` `<+>` Wendi gives her flock 15 cups of feed in the morning and another 25 cups in the afternoon, for a total of  $15+25 = \ll 15+25=40 \gg 40$  cups of feed.

`<prefix>` `<*>` If Wendi has 20 chickens, then she needs  $20*3 = \ll 20*3=60 \gg 60$  cups of feed to feed her flock.

`<prefix>` `<->` If Wendi has already given her flock 40 cups of feed, then she needs to give her flock  $60-40 = \ll 60-40=20 \gg 20$  more cups of feed.

`<prefix>` `<answer>` The answer is: 20

# Algorithm overview



# Results on Math Word Datasets

| LM                | Method                      | #clusters | #trainable | GSM8K       | MATH       | AQUA        | Avg         |
|-------------------|-----------------------------|-----------|------------|-------------|------------|-------------|-------------|
| Phi 1.5<br>(1.3B) | Full-FT                     | 0         | 100%       | 12.5        | 1.3        | 27.2        | 13.5        |
|                   | Full-FT + General           | 1         | 100%       | 15.4        | 2.0        | 35.4        | 17.6        |
|                   | Full-FT + <b>Arithmetic</b> | 4         | 100%       | 15.0        | 2.3        | 33.1        | 16.8        |
|                   | Full-FT + <b>K-Means</b>    | 5         | 100%       | 14.5        | 2.7        | <b>36.5</b> | 17.7        |
|                   | Full-FT + <b>SQ-VAE</b>     | 5         | 100%       | <b>15.8</b> | <b>3.3</b> | 34.3        | <b>17.8</b> |
| Llama2<br>(7B)    | LoRA                        | 0         | 0.343%     | 38.2        | 6.5        | 36.6        | 27.1        |
|                   | LoRA + General              | 1         | 0.344%     | 38.5        | 6.7        | 37.8        | 27.7        |
|                   | LoRA + <b>Arithmetic</b>    | 4         | 0.344%     | 39.5        | 5.6        | 38.2        | 27.8        |
|                   | LoRA + <b>K-Means</b>       | 5         | 0.344%     | 39.1        | 6.7        | 40.5        | 28.8        |
|                   | LoRA + <b>SQ-VAE</b>        | 5         | 0.344%     | <b>40.0</b> | <b>7.0</b> | <b>41.3</b> | <b>29.4</b> |
| Llama2<br>(13B)   | LoRA                        | 0         | 0.279%     | 44.6        | 7.2        | 41.3        | 31.0        |
|                   | LoRA + General              | 1         | 0.280%     | 47.9        | 7.9        | 42.5        | 32.8        |
|                   | LoRA + <b>Arithmetic</b>    | 4         | 0.280%     | 41.9        | 4.6        | 35.8        | 27.4        |
|                   | LoRA + <b>K-Means</b>       | 5         | 0.280%     | 49.6        | 8.4        | <b>44.1</b> | 34.0        |
|                   | LoRA + <b>SQ-VAE</b>        | 5         | 0.280%     | <b>50.6</b> | <b>8.5</b> | 43.9        | <b>34.3</b> |

# Reasoning length effect

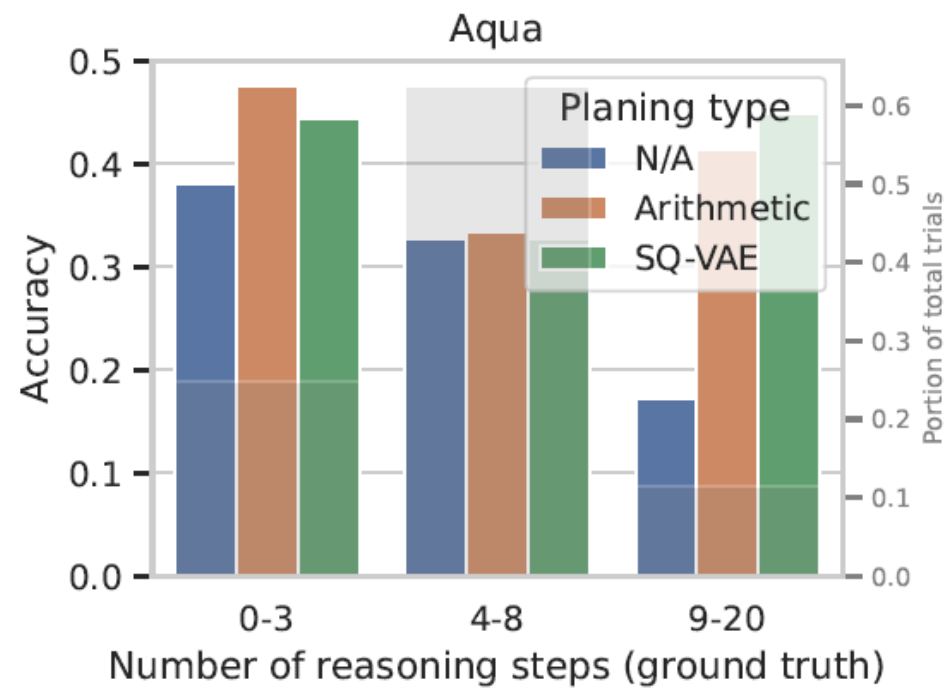
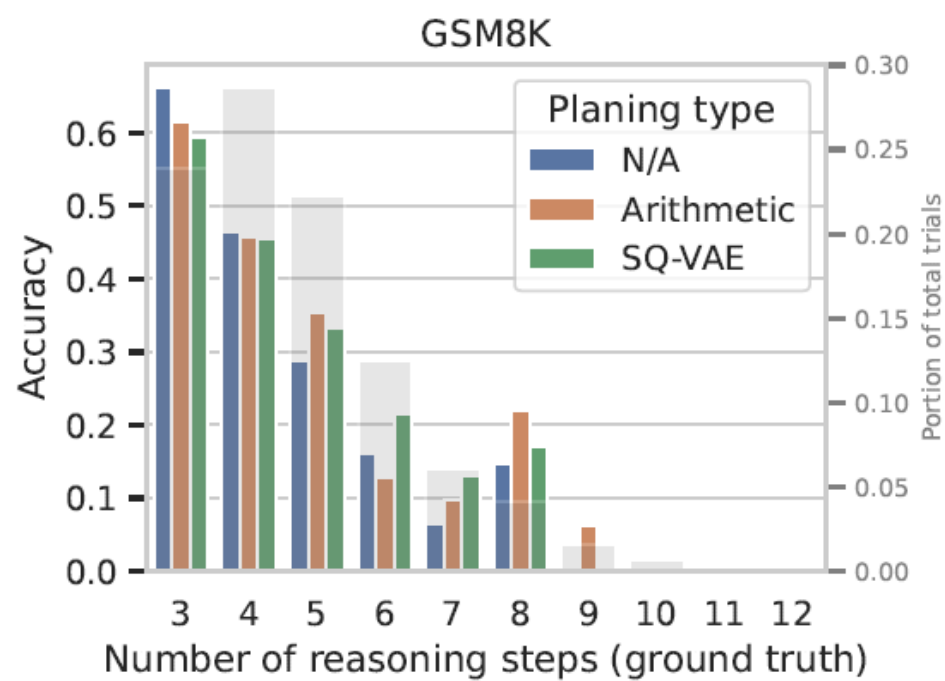


Figure 2: Accuracy on GSM8K (**left**) and Aqua (**right**) on test examples by their number of ground-truth reasoning steps. SQ-VAE consistently increases performance for test examples that require more steps of reasoning to be solved.



Carnegie  
Mellon  
University

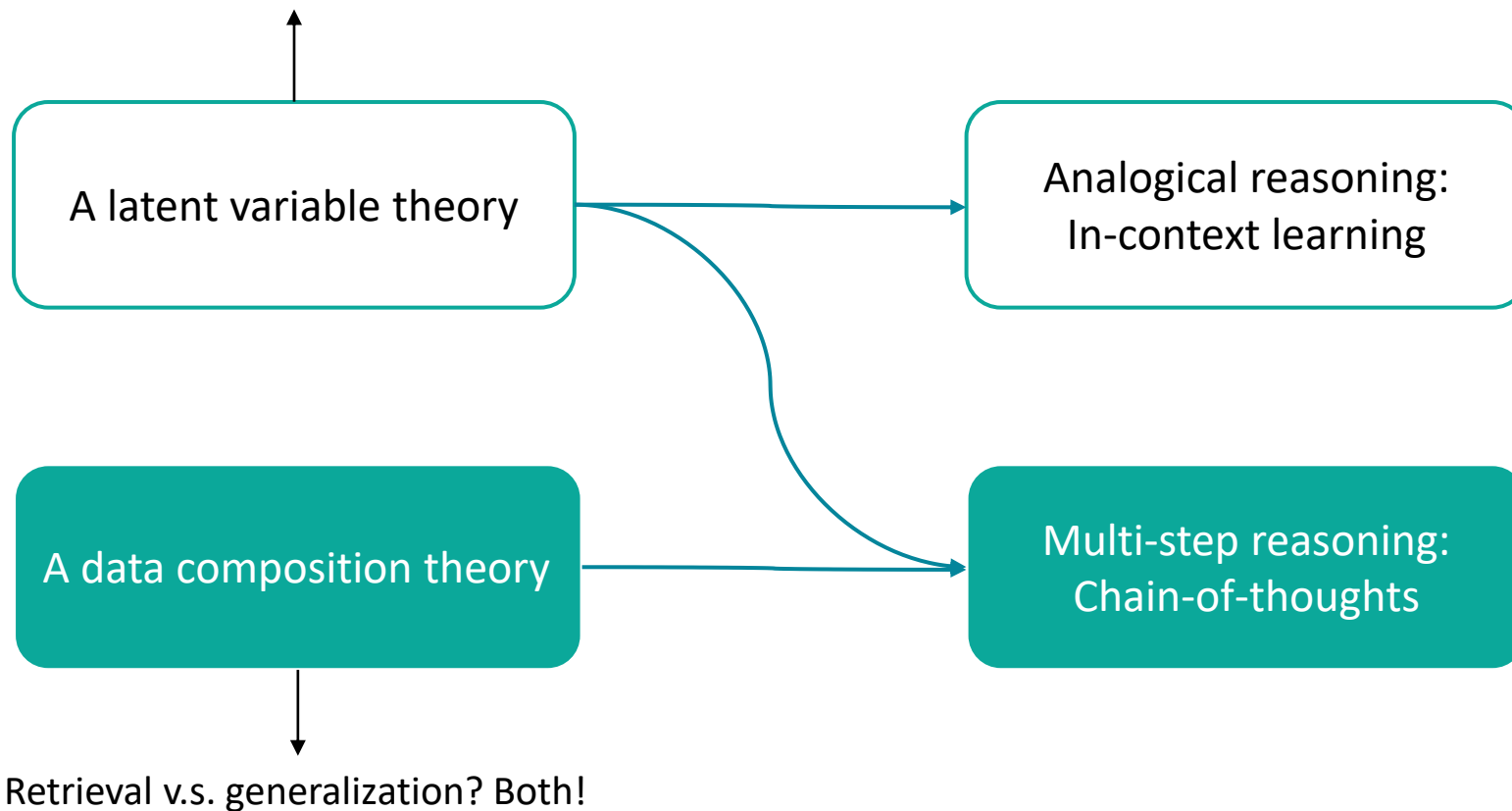
UC SANTA BARBARA

# Understanding the Reasoning Ability of Language Models From the Perspective of Reasoning Paths Aggregation

Xinyi Wang, Alfonso Amayuelas, Kexun Zhang, Liangmin Pan, Wenhui Chen, William Yang Wang ([Arxiv](#))

# Our Proposed LLM Theories

There is a latent variable governing the reasoning process.



# Reasoning with LLMs

## Chain-of-Thought Prompting

### Model Input

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.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

### Model Output

A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had  $23 - 20 = 3$ . They bought 6 more apples, so they have  $3 + 6 = 9$ . The answer is 9. ✓

## (d) Zero-shot-CoT (Ours)

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A: **Let's think step by step.**

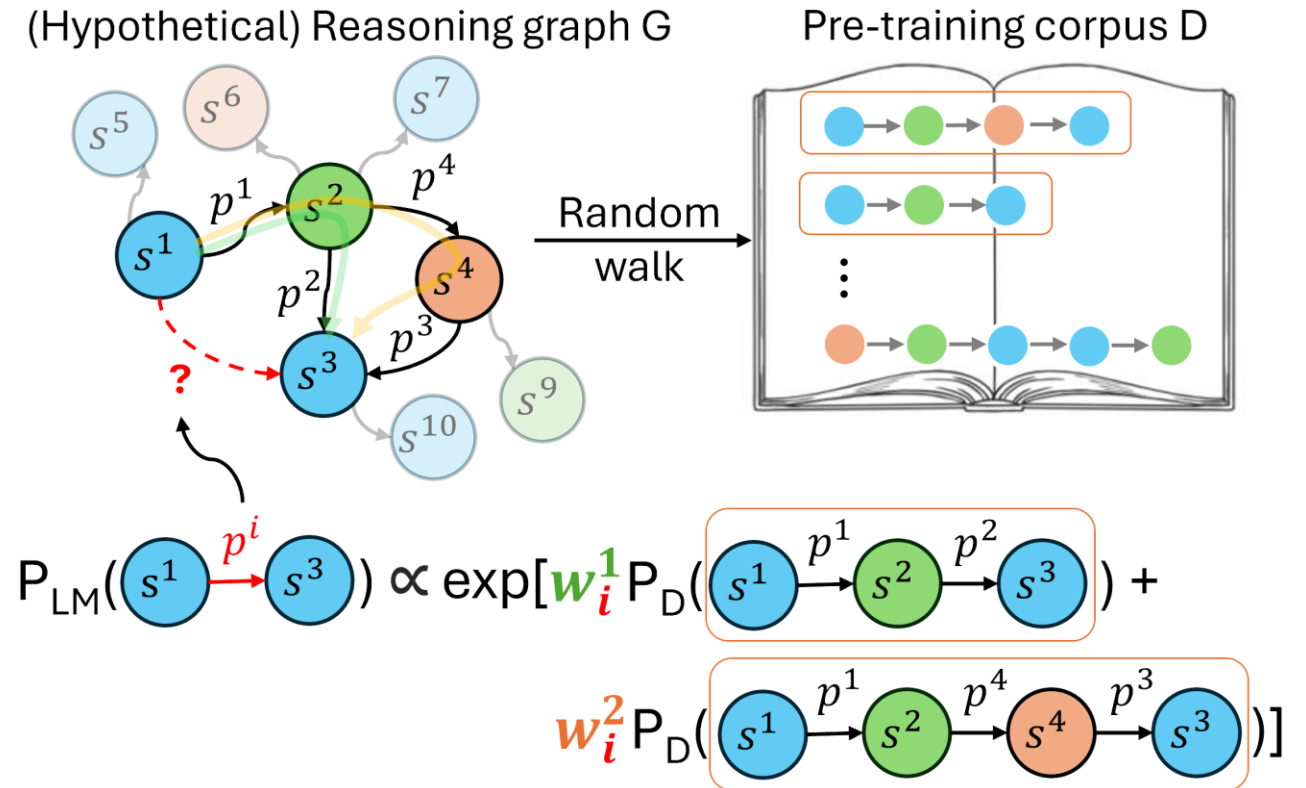
(Output) *There are 16 balls in total. Half of the balls are golf balls. That means that there are 8 golf balls. Half of the golf balls are blue. That means that there are 4 blue golf balls. ✓*

(Source: [Wei et al. 2022](#); [Kojima et al. 2022](#))

- **Definition** of reasoning: deriving new conclusions with novel conditions from the known facts.
- **Observation:** Pre-trained-only base LLMs exhibit impressive reasoning capability without any fine-tuning.

# Understand reasoning ability of LLMs

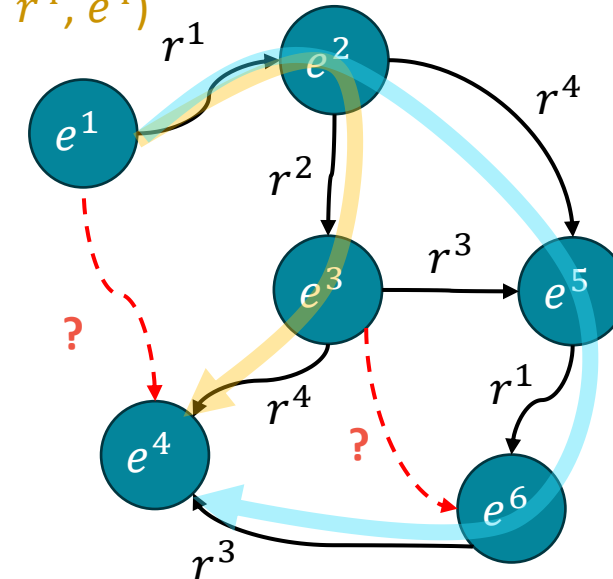
- **Hypothesis:** LLMs can retrieve and aggregate (random walk) reasoning paths seen at pre-training time to do complex reasonings at inference time.
- **Approach:** We study two specific cases of reasoning:
  - **logical**/knowledge graph (KG) reasoning: pre-train a toy transformer on KGs
  - **mathematical** reasoning: continue (pre-)train a pre-trained LM on more unlabeled augmented reasoning paths.



# Logical reasoning with knowledge graph

- **Knowledge graph (seen triples):**  $(e^1, r^1, e^2)$ ,  $(e^2, r^2, e^3)$ ,  $(e^2, r^4, e^5)$ ,  $(e^3, r^3, e^5)$ ,  $(e^3, r^4, e^4)$ ,  $(e^5, r^1, e^6)$ ,  $(e^6, r^3, e^4)$
- **Unseen triples:**  $(e^1, r^3, e^4)$ ,  $(e^3, r^4, e^6)$
- **Task:** how to infer unseen triples from the seen ones?
  - We propose to look at  $P(\text{tail} \mid \text{head}, \text{relation})$ .
  - i.e.  $P(e^4 \mid e^1, r^3)$ ,  $P(e^6 \mid e^3, r^4)$ .
  - The sample space is all entities in the graph.
- **Language model training:**
  - Translate each entity and relation into a new token.
  - Sample random walk paths from the knowledge graph to form the pre-training data.
  - Use the next-token-prediction objective to pre-train a small transformer based LM from scratch.
- **Language model inference:**
  - prediction of the tail entity: prompt the LM with head entity and relation.

path1 =  $(e^1, r^1, e^2), (e^2, r^2, e^3), (e^3, r^4, e^4)$



path2 =  $(e^1, r^1, e^2), (e^2, r^4, e^5), (e^5, r^1, e^6), (e^6, r^3, e^4)$

# Distributions

- Language model:

$$P_{\text{LM}}(e_2|e_1, r) = \frac{\exp(f_{\theta}(e_2|e_1, r))}{\sum_{e \in \mathcal{E}} \exp(f_{\theta}(e|e_1, r))} \quad (2)$$

logits

All Entities

- Unweighted aggregation:

A simplified version of  $P_w$  would be letting  $w_r(h) = 1$  for all  $h$  and  $r$ . And **we define this unweighted aggregation distribution to be  $P_s$** :

$$P_s(e_2|e_1, r) = \frac{\exp(\sum_{h \in \mathcal{H}_r} P(e_2|e_1, h)/T)}{\sum_{e \in \mathcal{E}} \exp(\sum_{h \in \mathcal{H}_r} P(e|e_1, h)/T)} \quad (4)$$

Rules related to relation  $r$

Sum of Random walk paths probability

- Weighted aggregation:

Path ranking algorithm (PRA) (Lao et. al. 2011)

$$P_w(e_2|e_1, r) = \frac{\exp(S_w(e_2|e_1, r)/T)}{\sum_{e \in \mathcal{E}} \exp(S_w(e|e_1, r)/T)} \quad (3)$$

$$S_w(e_2|e_1, r) = \sum_{h \in \mathcal{H}} w_r(h) P(e_2|e_1, h)$$

All possible rules

Rule weight learned by logistic regression

Sum of Random walk paths probability

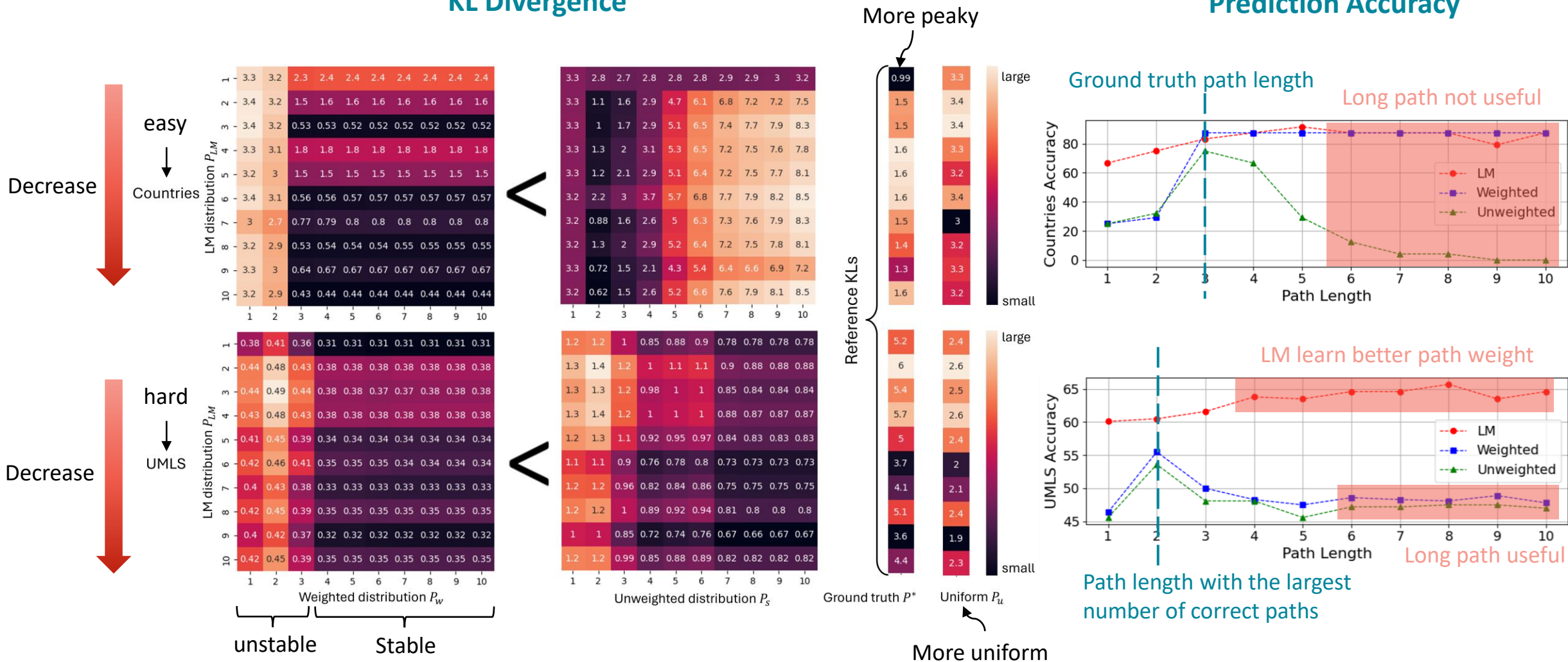
$$P(e_n|e_0, h) = \sum_{(e_0, r_1, e_1) \dots (e_{n-1}, r_n, e_n) \in \mathcal{P}_h} \prod_{i=1}^n P(e_i|e_{i-1}, r_i)$$

Uniform for random walk

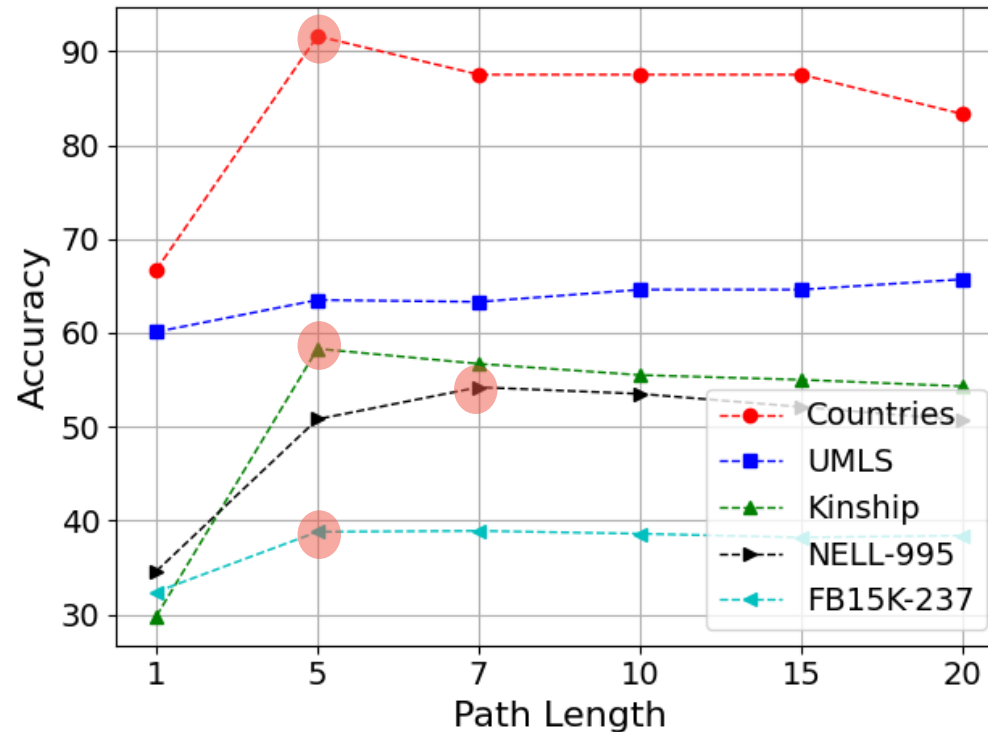
# Comparing LM with path aggregation

## KL Divergence

## Prediction Accuracy



# Ablation on random walk path length

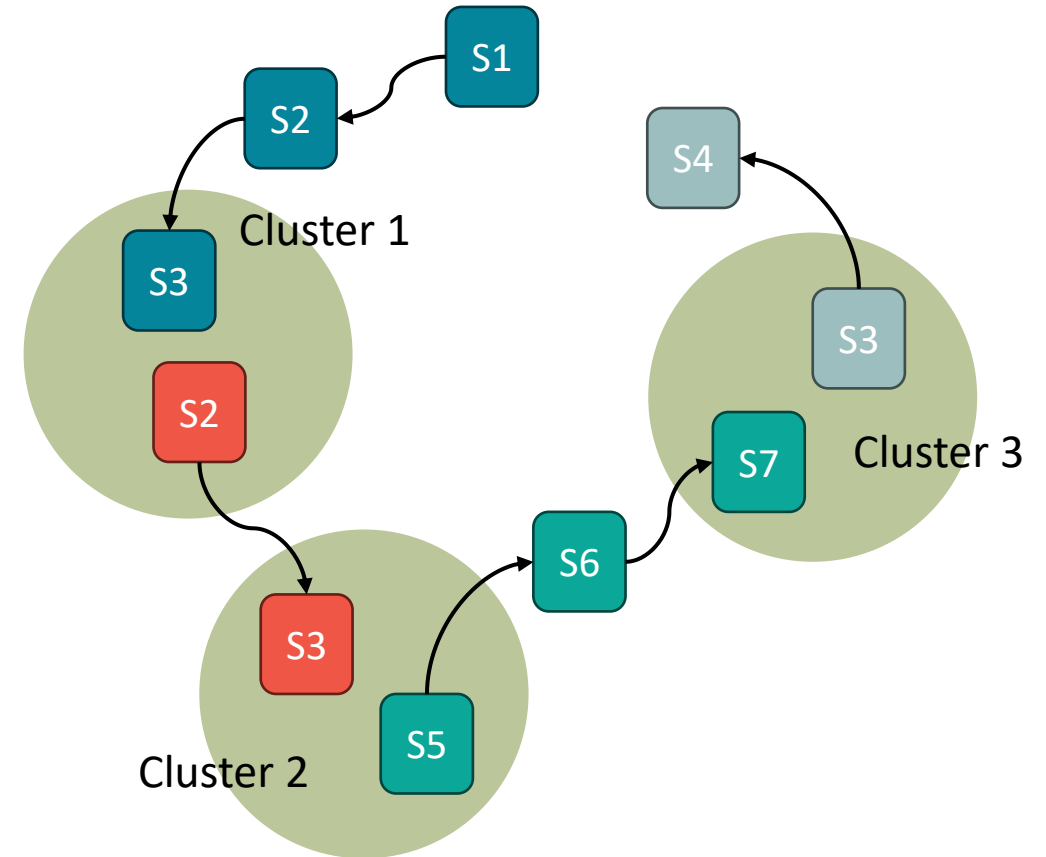


\* Exists an optimal path length

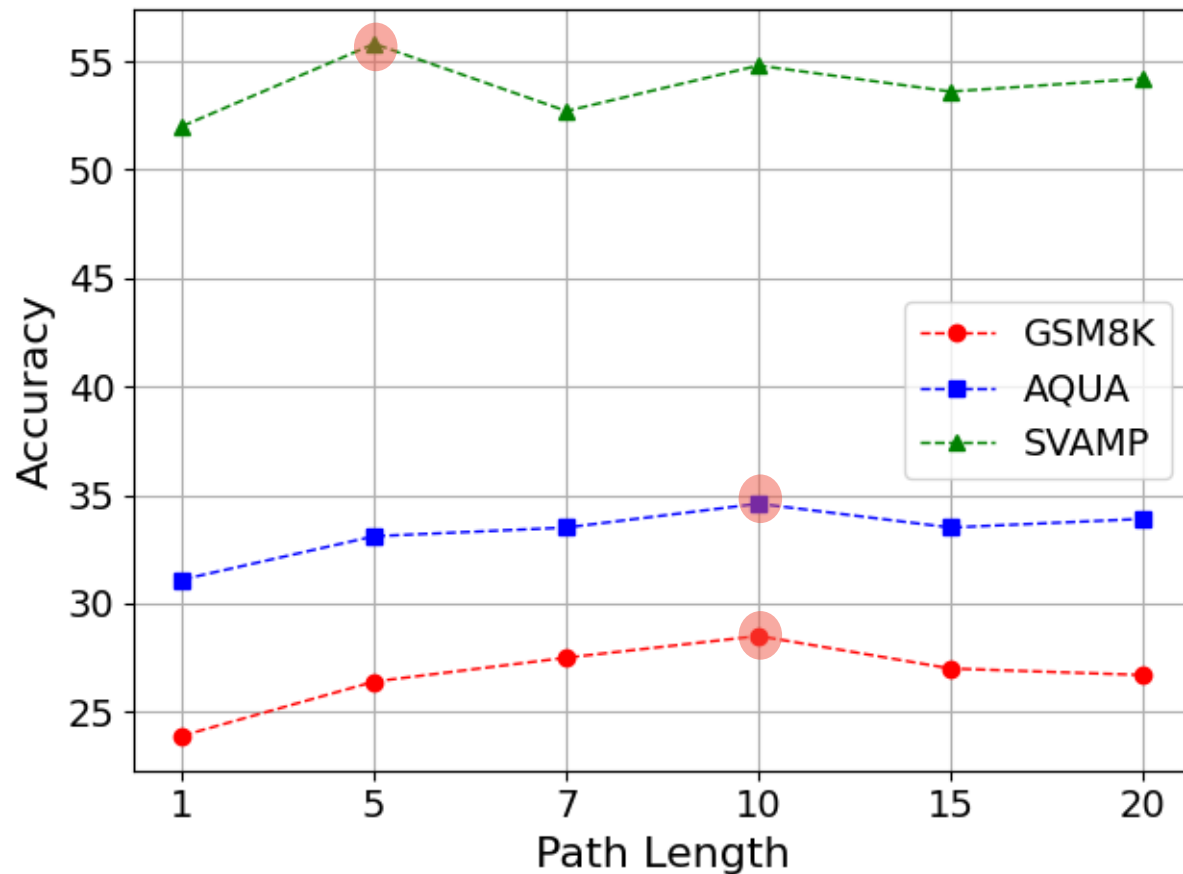
- Test accuracy (%) of GPT2 pre-trained on different length random walk paths.
- All entities and relationships are translated as new tokens. i.e. no natural language involved.

# Math reasoning with Chain-of-thoughts (CoT)

- **Reasoning graph:** CoTs can be regarded as walking on a graph whose nodes represent the current reasoning state.
- **Encode reasoning state:** cumulatively embed each CoT step with a pre-trained LM.
- **Construct nodes:** cluster the CoT steps. Each cluster represents a node in the reasoning graph.
- **Random walk** on this graph:
  - Step 1. Randomly select a starting step
  - Step 2. Follow the original CoT for m steps
  - Step 3. In the cluster of the end step, randomly select another step as the next step.
  - Step 4. Go back to step 2.
- This can be regarded as a light-weight **data augmentation** method for CoT reasoning.



# Ablation on random walk path length



\* Exists an optimal path length

- GSM8K, AQUA, SVAMP are three math word datasets. We LORA fine-tune a Llama 2 (7B) model on CoT data.
- We first do 500 steps of random walk training then 2000 steps of regular supervised fine-tuning.

# More results & ablations

| Model | Method | GSM8K       | AQUA        | SVAMP       | Avg.        |
|-------|--------|-------------|-------------|-------------|-------------|
| 7B    | SFT    | 26.8        | 30.0        | 53.3        | 36.7        |
|       | Ours   | <b>28.5</b> | <b>34.6</b> | <b>55.8</b> | <b>39.6</b> |
| 13B   | SFT    | 37.1        | 35.0        | 66.4        | 46.2        |
|       | Ours   | <b>41.2</b> | <b>37.4</b> | <b>69.0</b> | <b>49.2</b> |

*Table 1.* Testing accuracy of different size Llama 2 models continue pre-trained with our random walk paths and then supervised fine-tuned. The supervised fine-tuning baseline (SFT) is fine-tuned by the same number of total steps. Results are reported on three math word problem (MWP) datasets.

| #Nodes | GSM8K       | AQUA        | SVAMP       | Avg.        |
|--------|-------------|-------------|-------------|-------------|
| 0      | 26.8        | 30.0        | 53.3        | 36.7        |
| 10     | 26.8        | 30.3        | 54.8        | 37.3        |
| 50     | 26.6        | 29.9        | 54.7        | 37.1        |
| 100    | <b>28.5</b> | <b>34.6</b> | <b>55.8</b> | <b>39.6</b> |
| 200    | 26.6        | 31.1        | 52.5        | 36.7        |

*Table 3.* Ablation on the number of clusters/nodes  $K$ .

| #Steps | GSM8K       | AQUA        | SVAMP       | Avg.        |
|--------|-------------|-------------|-------------|-------------|
| 0      | 26.8        | 30.0        | 53.3        | 36.7        |
| 200    | 27.5        | 30.1        | 53.6        | 37.1        |
| 500    | <b>28.5</b> | <b>34.6</b> | <b>55.8</b> | <b>39.6</b> |
| 1000   | 24.9        | 32.3        | 51.6        | 36.3        |

*Table 2.* Ablation on the number of random walk training steps  $M$ .

- Ablation on model size, number of clusters, and number of training steps.
- #Nodes = 0 and #Steps = 0 means we don't do any random walk training.

**Future Research**

# Research plan

- Future research directions:
  - Zoom in more: how model parameters correspond to proposed theory.
  - Exploring better ways to verify the proposed theory with real world LLMs.
  - Connections between foundation models of different modalities. E.g. language model and diffusion

**Thank you!**

Questions?