

# Finding the Minimal Parameter Budget for Implicit Reasoning: A Data Complexity Driven Scaling Law for Language Models

Xinyi Wang  
March 2026

# How much model capacity is needed for reasoning?

Current paradigm:

- Bigger models → better reasoning
- Chain-of-thought improves reasoning

But:

- CoT happens after pretraining
- Reasoning ability must already exist in base models

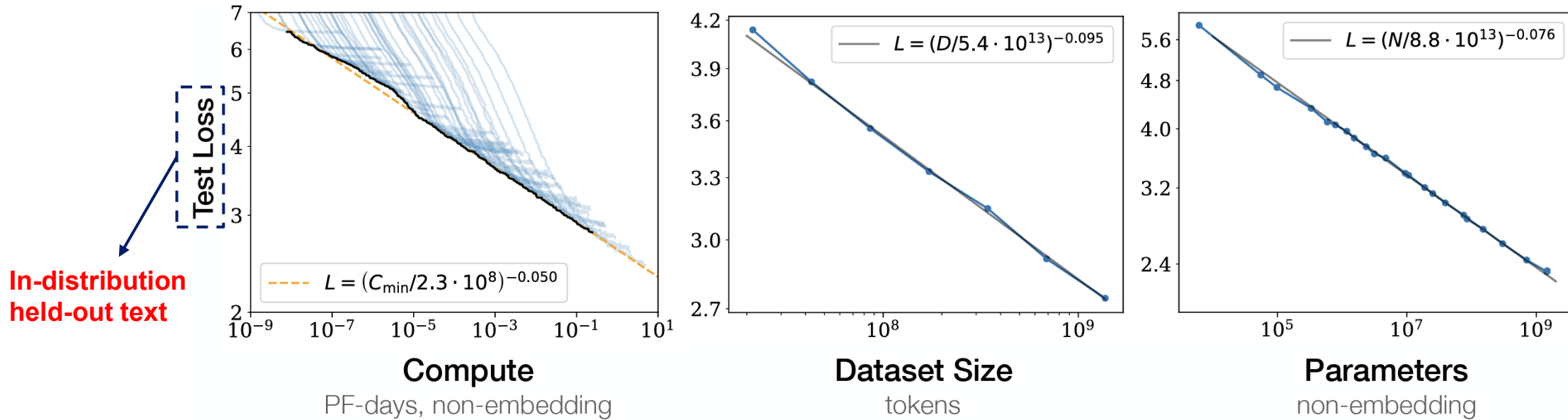
**Question:**

- What capacity is required for reasoning during pretraining?

**Hypothesis:**

- There is a minimum model size depending on reasoning complexity

# Existing Scaling Laws



**Figure 1** Language modeling performance improves smoothly as we increase the model size, dataset size, and amount of compute<sup>2</sup> used for training. For optimal performance all three factors must be scaled up in tandem. Empirical performance has a power-law relationship with each individual factor when not bottlenecked by the other two.

**Do not explain reasoning!**

(Source: [Kaplan et al. 2020](#))

# Explicit & Implicit Reasoning with LLMs

Scaling law for implicit reasoning!

## Explicit Reasoning

Question: Samantha had 5 packs of markers. Each pack had 12 markers. She gave 9 markers to her friend and lost 3. How many markers does Samantha have now?

<think>

Step 1: Samantha has 5 packs of 12 markers:  $5 \times 12 = 60$  markers in total.


Step 2: She gives 9 markers to her friend:  $60 - 9 = 51$ .

Step 3: She loses 3 markers:  $51 - 3 = 48$ .


Step 4: She has 48 markets now.

</think>

**The final answer is 48.**



Inefficient



Constraint

## Implicit Reasoning

Question: Samantha had 5 packs of markers. Each pack had 12 markers. She gave 9 markers to her friend and lost 3. How many markers does Samantha have now?

↓

Layer 1 / State 1


Layer 2 / State 2

Layer 3 / State 3


Layer 4 / State 4

↓

**The final answer is 48.**

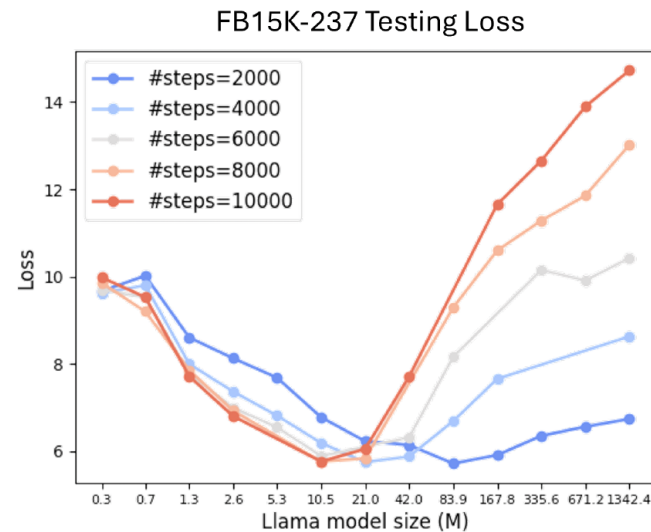
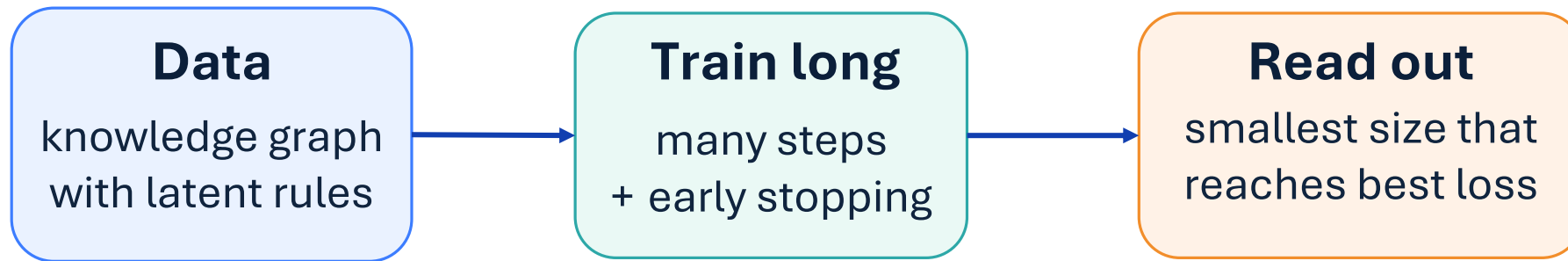


Efficient



Diverse

# Setup to Find Minimum Parameter Budget



# Detailed Setup

What models?

- LM with LLaMA-like architecture

What loss?

- Next token prediction loss

What training data?

- implicit reasoning environment that can be rigorously controlled

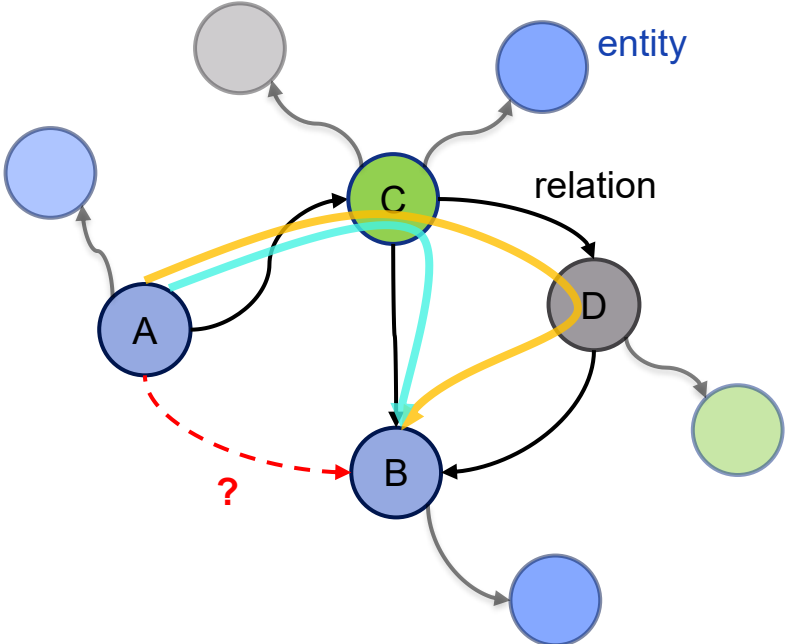
What evaluation?

- testing loss/accuracy for implicit reasoning

How to construct a controllable  
implicit reasoning environment?

# Reasoning as Graph Completion

Reasoning = draw connections between previously disconnected concepts



G: World knowledge, knowledge graph...


# Data Format

- **Training data:** incomplete (synthetic/real) knowledge graph
- **Data format:** each training example is a knowledge triple
- **Entity/relation representation:** random IDs (character-level tokenization)

```
"triple": {  
  "head": "Princeton University",  
  "relation": "state",  
  "tail": "New Jersey"  
}
```

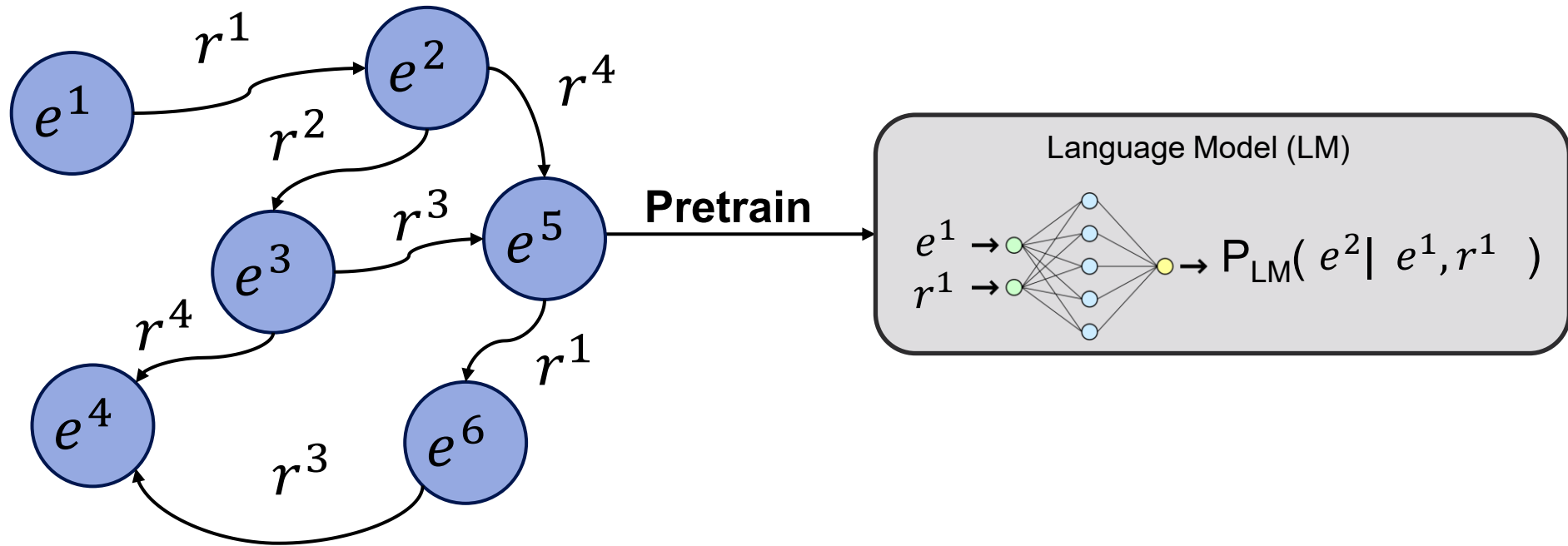
Random IDs

113, 23, 45.

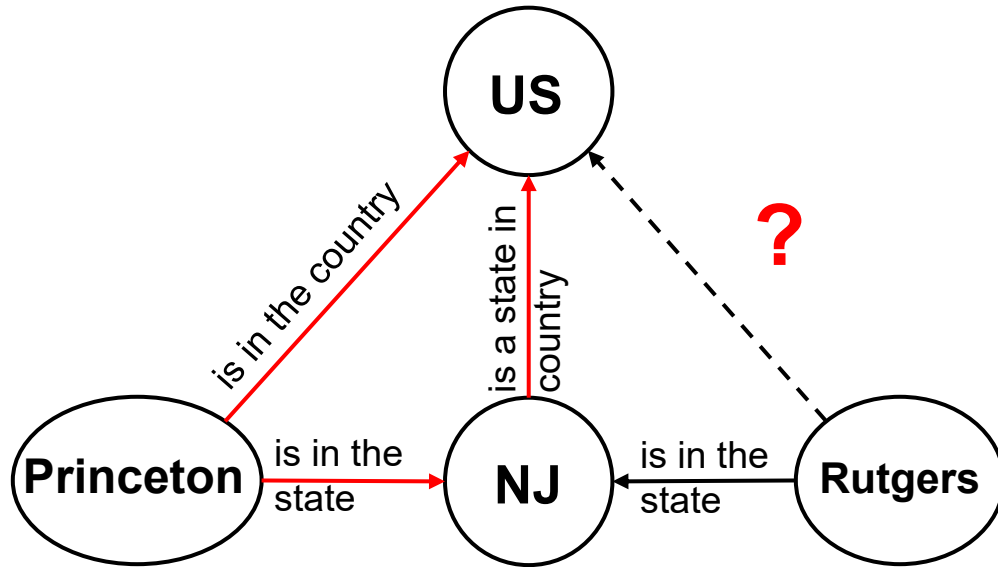


# Pretraining

Training data: incomplete (synthetic/real) knowledge graph



# Evaluation

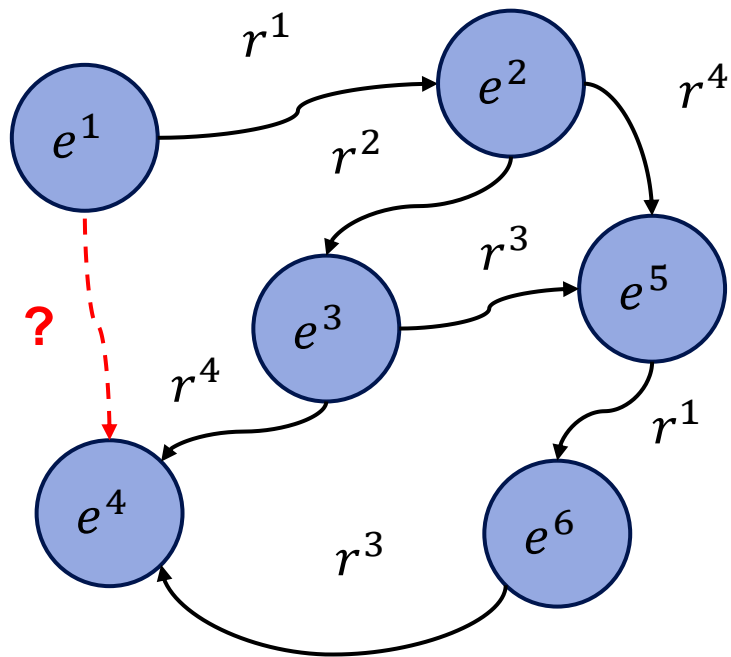


Repeated many times in training data!

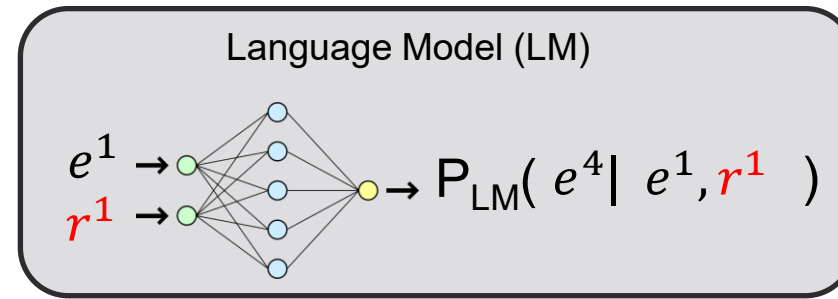
**Learned logic rule:**  $\text{is in the state} \wedge \text{is a state} \rightarrow \text{is in the country}$

# Evaluation

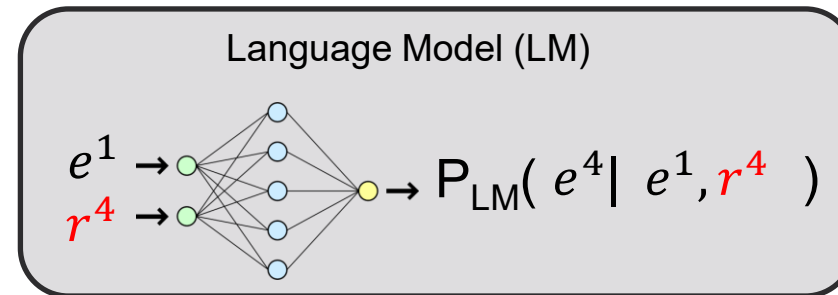
Testing data: unseen edges (can be reasoned logically)



## Inference



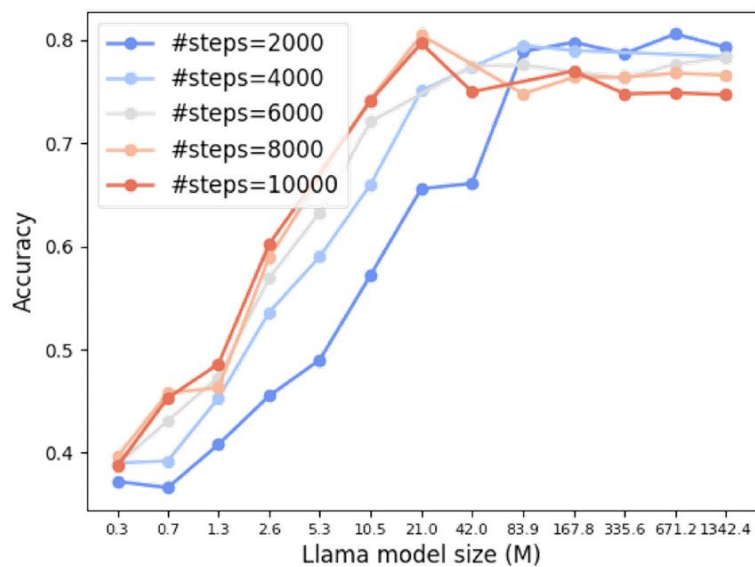
...



Select the relation with the largest probability

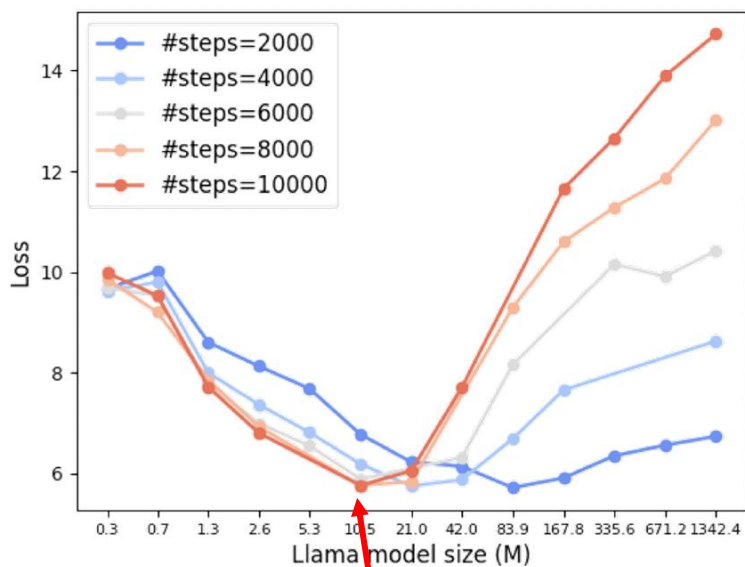
# Real-World KG Results

FB15K-237 Reasoning Accuracy



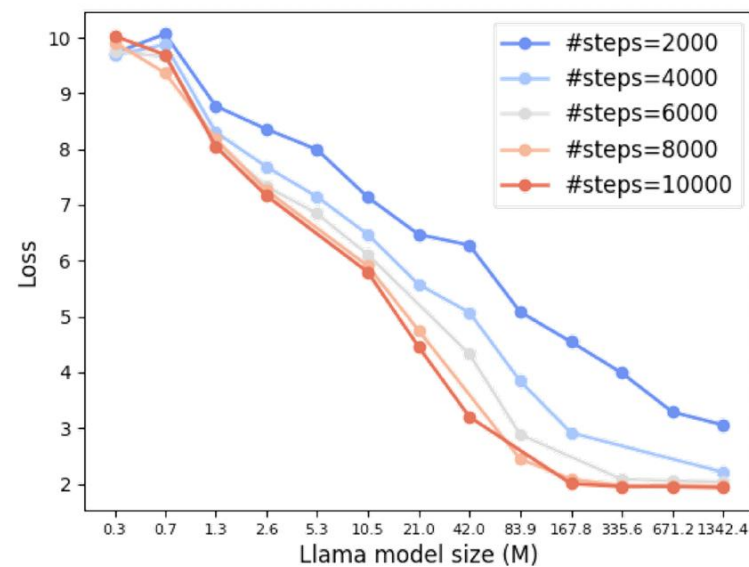
Multiple choice

FB15K-237 Testing Loss



Converged optimal model size

FB15K-237 Training Loss



# Synthetic KG for Better Control

Key steps of KG generation:

Conjunctive rules generation:

- A is B's father  $\wedge$  B is C's father  $\rightarrow$  A is C's grandfather

Control possible number of relation types connecting to each entity

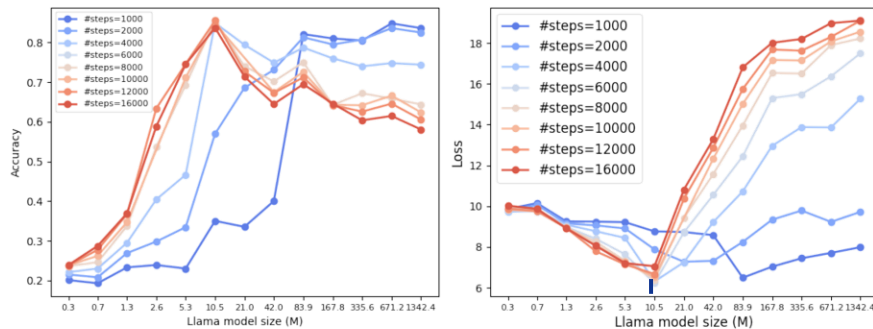
- real-world alignment
- reduce noises

Grow the graph with preferential attachment to ensure power law degree distribution

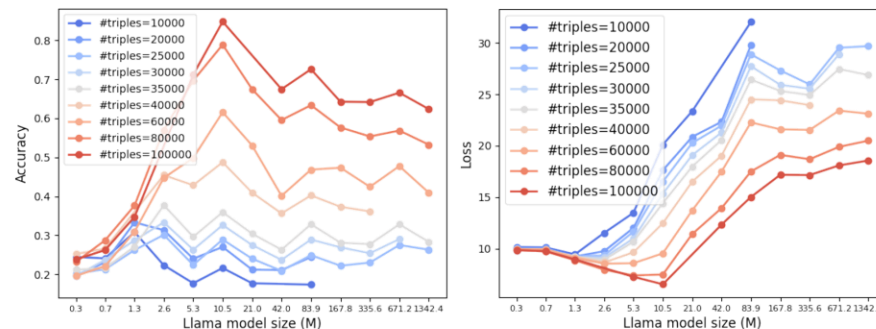
- real-world alignment

# Effect of Hyperparameters

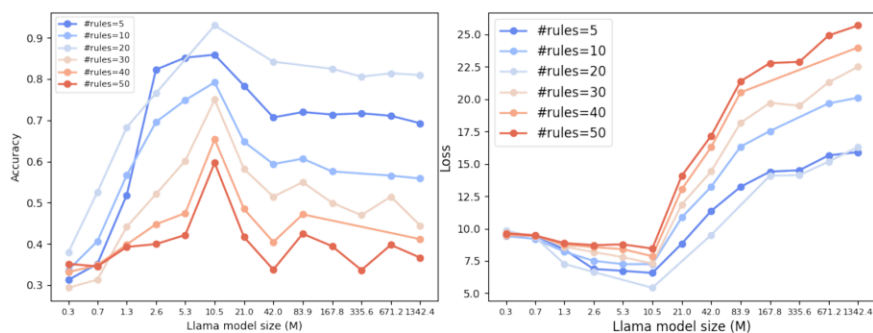
(a) Effect of Training Steps



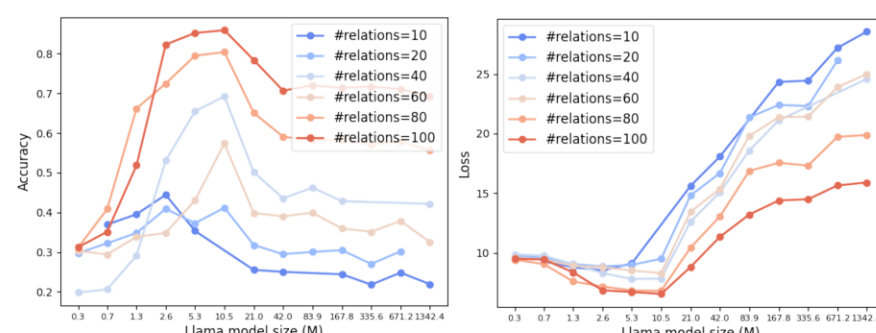
(b) Effect of #Triples



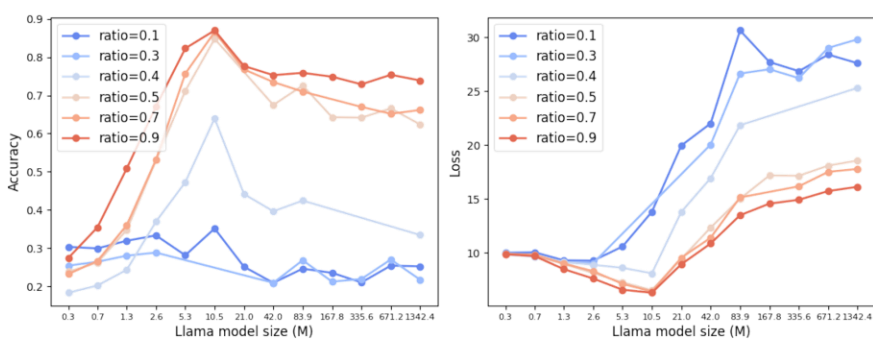
(c) Effect of #Rules



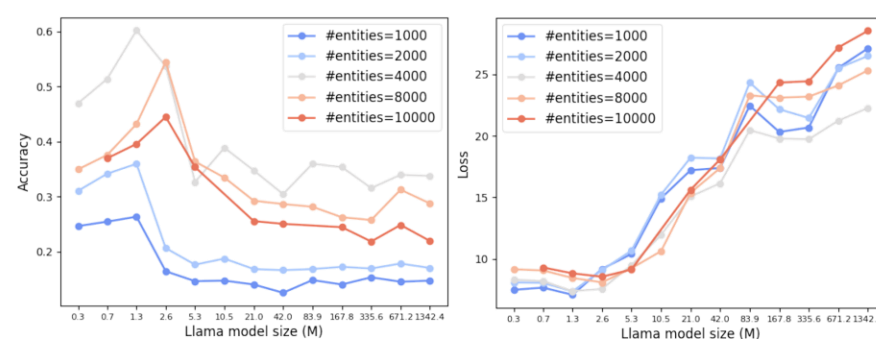
(d) Effect of #Relations



(e) Effect of Deductible Ratio

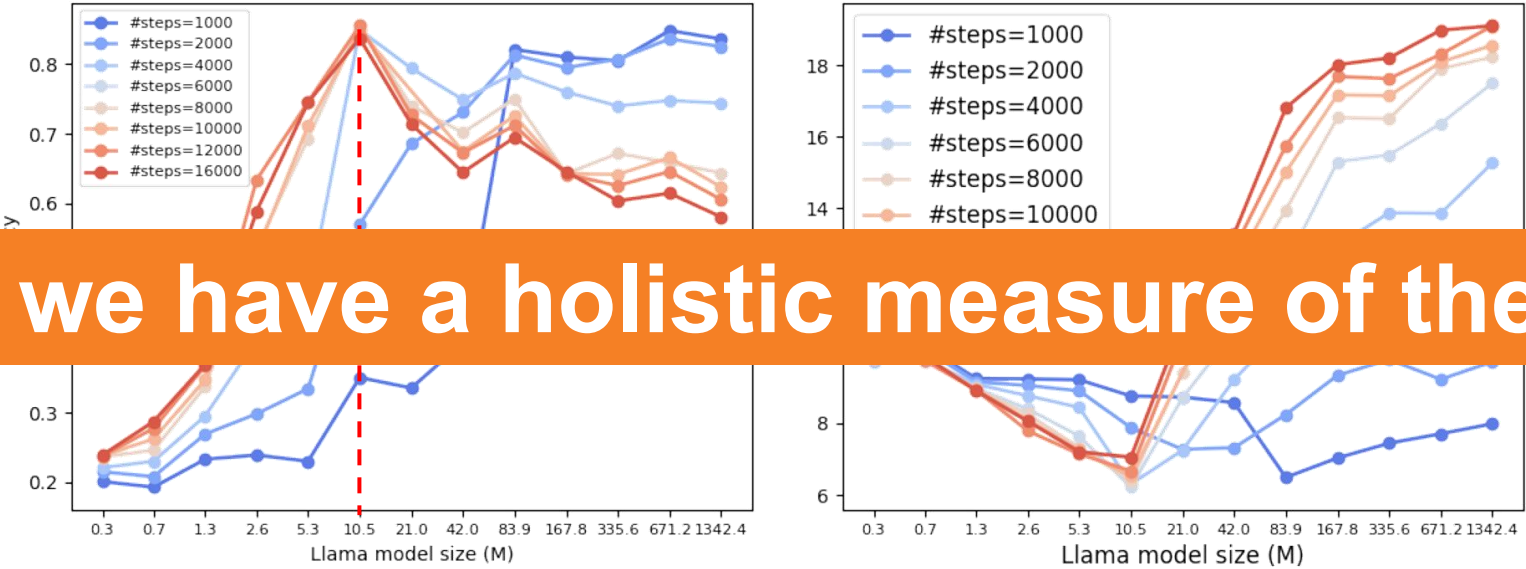


(f) Effect of #Entities



# Effect of Training Steps

(a) Effect of Training Steps



Can we have a holistic measure of the graph?

Optimal mode size

Larger training steps converge to the optimal accuracy/loss

(see paper for theoretical justification)

# Graph Search Entropy

**Def:** #nodes \* entropy rate of an infinitely long random walk on G

$$H(G) = \overset{\text{\#nodes}}{N_e} (\log(\lambda)) + H^r(G).$$

**Entity** entropy rate:  
entropy produced by each  
**node** on the walk trace for an  
infinitely long random walk

**Relation** entropy rate:  
entropy produced by each  
**edge** on the walk trace for an  
infinitely long random walk

# A Closer Look

**Def:** #nodes \* entropy rate of an infinitely long random walk on G

$$H(G) = N_e (\log(\lambda)) + H^r(G).$$

dominant eigenvalue of  
the adjacency matrix

$$H^r(G) = - \sum_{i=1}^{N_e} \rho_i \sum_{j=1}^{N_r} S_{ij}^r \log(S_{ij}^r).$$

stationary distribution

Relation transition  
matrix

# Minimum Sufficient Model Size Scales Linearly with Graph Search Entropy

## Theorem 4 (informal)

Under

- random IDs, [no semantic sharing across entities.]
  - finite parameter precision, and [An  $N$ -parameter model has  $O(N)$  effective information capacity]
  - realizability at entropy cost, [The required predictor family can be realized with  $O(H(G))$  parameters.]
- the  $\varepsilon$ -optimal model size obeys

$$c_1 H(G) \leq N_\theta^*(G) \leq c_2 H(G) + O(1)$$

$$\therefore N_\theta^*(G) = \Theta(H(G))$$

# Minimum Sufficient Model Size Scales Linearly with Graph Search Entropy

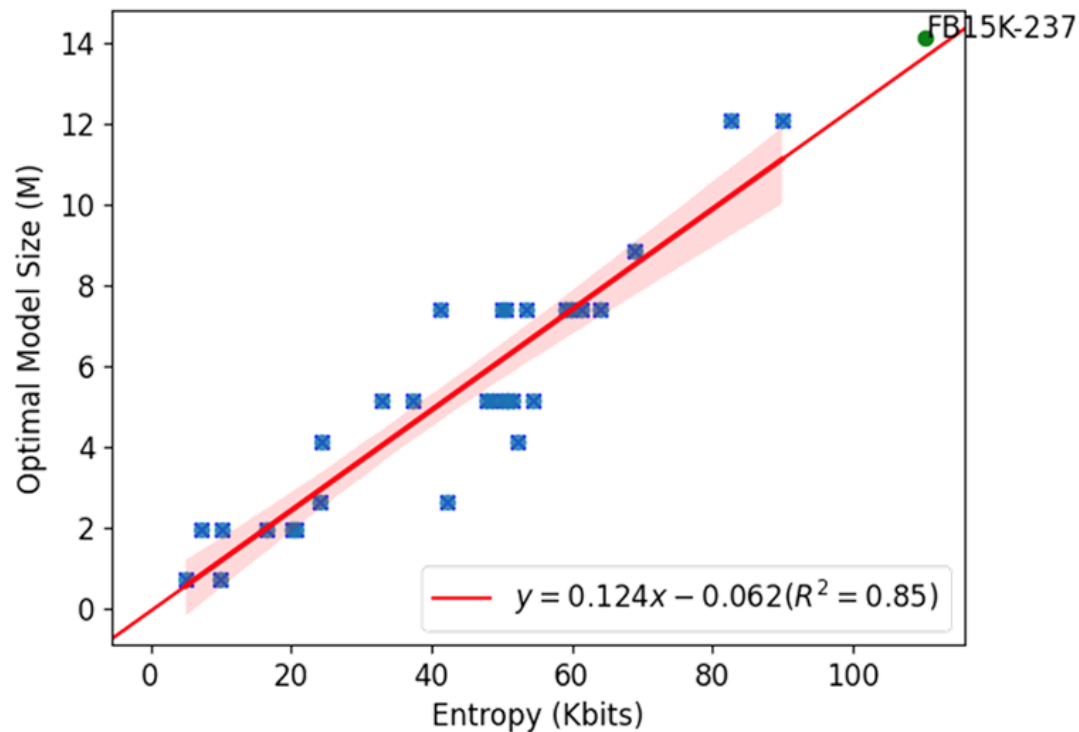


Figure 4: The optimal model size with the lowest possible testing loss v.s. the graph search entropy. The red line is the linear regression line using data from the synthetic experiments (blue squares), with a 95% confidence interval. We also plot the graph search entropy and optimal model size from the real-world FB15K-237 experiment (green dot) to verify the accuracy of the obtained linear scaling law.

# Reasoning Scaling v.s. Knowledge Scaling

knowledge capacity scaling law ([Allen-Zhu & Li, 2025](#)):

- Entropy: knowledge generation process
- 2 bits / parameter for memorization.

Ours implicit reasoning scaling law:

- Entropy: randomly traversing a (fixed) knowledge graph
- 0.008 bits / parameter for reasoning.

⇒ Reasoning is far more data-demanding than knowledge memorization.

# Takeaways

- We identify the minimal model size required for reasoning
- The size becomes well-defined with sufficient training
- The optimal model size scales linearly with data complexity
- The data complexity can be measure by the newly proposed graph search entropy

**Thank you!**

# Natural Language Data

- **Training data:** incomplete real knowledge graph
- **Data format:** each training example is a knowledge triple
- **Entity/relation representation:** natural language sentences

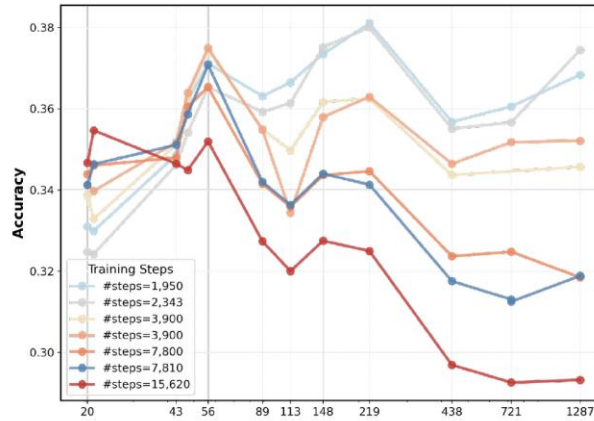
```
"triple": {  
  "head": "Princeton University",  
  "relation": "state",  
  "tail": "New Jersey"  
}
```

→ Princeton University is in the state of New Jersey.

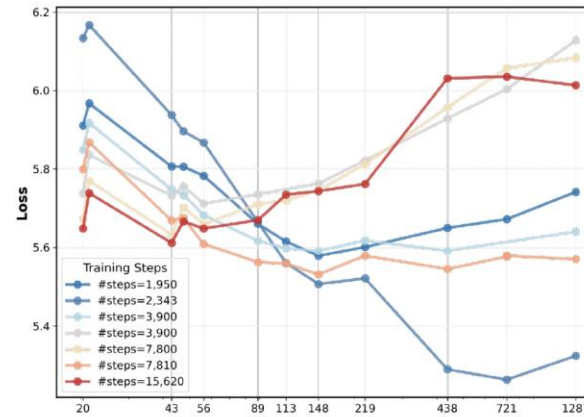
# Results

GPT-generated

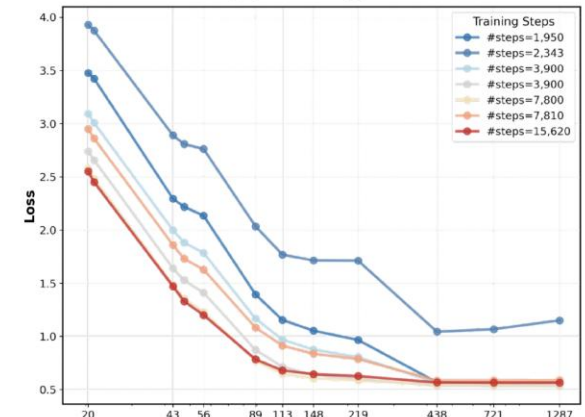
(a) Reasoning Accuracy



(b) Testing Loss

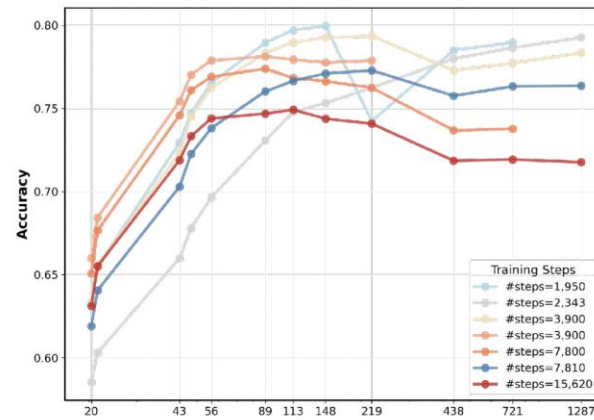


(c) Training Loss

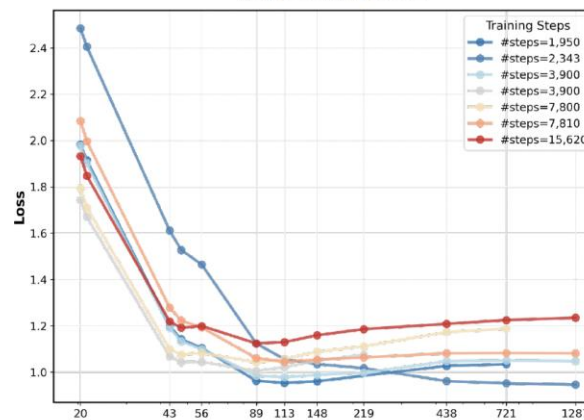


Template

(d) Reasoning Accuracy



(e) Testing Loss



(f) Training Loss

