

On Memorization of Large Language Models in Logical Reasoning

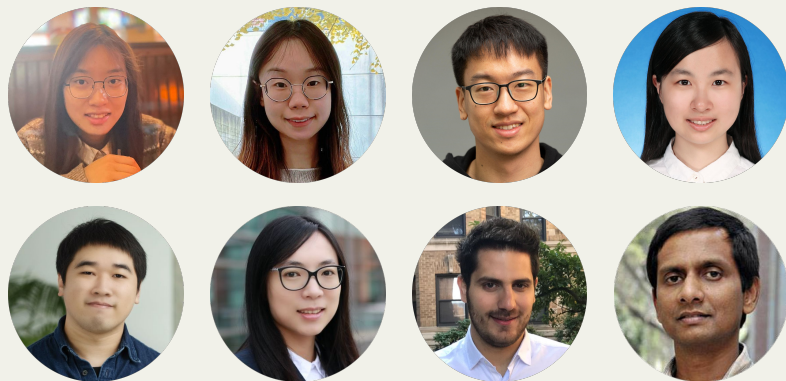
Chiyuan Zhang

Google Research



<https://memkklogic.github.io>

Chiyuan Zhang (Google Research)



Contents

- 00 My biased view on Shifting Perspective of Generalization
- 01 Motivation
- 02 LiMem Score and Knight & Knave Puzzles
- 03 Measurement of Memorization
- 04 Memorization vs Reasoning
- 05 Conclusion & Discussion
- 06 Related Work on LLM Memorization



Shifting Perspective of Generalization

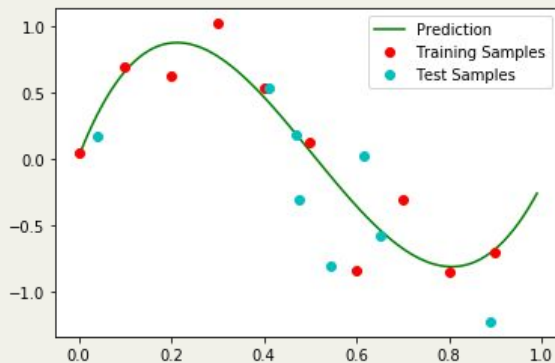
■ CLASSICAL ML

Test Acc = Train Acc + Gen Gap

Perspective: Bound Gen Gap for ERM

Best Practice: Model selection, under-parameterization or strong regularization

Low degree polynomial



High degree polynomial

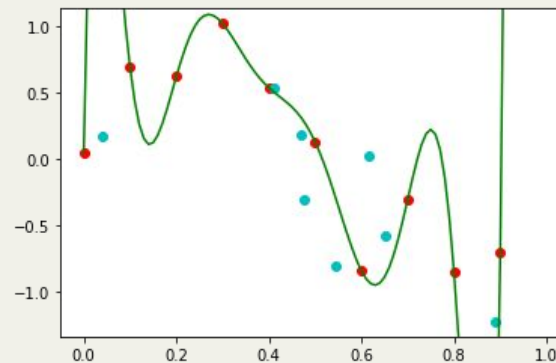


Fig from <https://github.com/jamt9000/prml/blob/master/1.1-polycurve.ipynb>

Shifting Perspective of Generalization

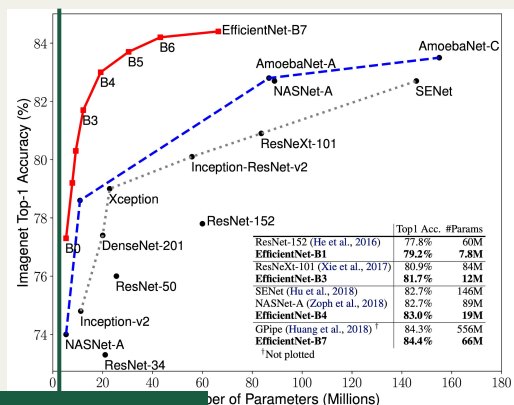
PRE-LLM REGIME

Overparameterized Models Work Well

Perspective: Non-unique ERM, (implicit) regularization from algo, architecture, data impact generalization.

Best Practice: Nonconvexity, Large mode/data, New NN Architecture, Optimizer, Data Augmentation, Initialization, Loss, etc...

Arch research & model sizes (2019)



ImageNet Trainset: 1.2M

Generalization Research

Improved Generalization Bound

Implicit Regularization

Benign Overfitting

Double Descent

Geometry of Loss Landscape

Out of Distribution Generalization

...

Fig source: <https://arxiv.org/abs/1905.11946>

Shifting Perspective of Generalization

LLM REGIME

Underparameterization + 1-epoch training + Emergent Abilities

Additional challenges: Difference between small / large models; Cost of analyzing large models

Best Practice: More data, Larger model (scaling law), Longer context window, More inference compute, ...

Growth of Training Costs

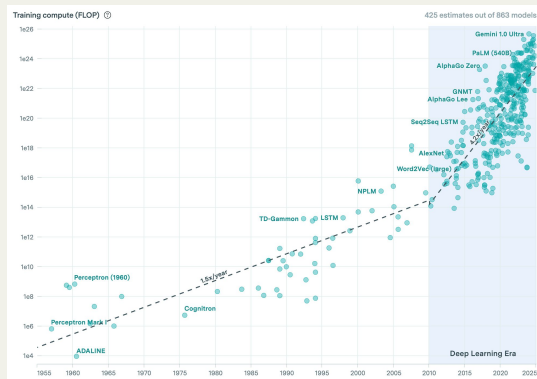
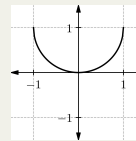


Fig source: <https://epoch.ai/data/notable-ai-models>

Example Problem from MATH dataset

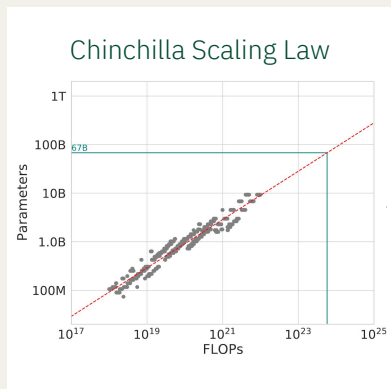
Let $f(x)$ be the function defined on $-1 \leq x \leq 1$ by the formula $f(x) = 1 - \sqrt{1 - x^2}$. This is a graph of $y = f(x)$:



If a graph of $x = f(y)$ is overlaid on the graph above, then one fully enclosed region is formed by the two graphs. What is the area of that region, rounded to the nearest hundredth?

MATH dataset: <https://huggingface.co/datasets/lighteval/MATH>

Generalization in LLM Regime

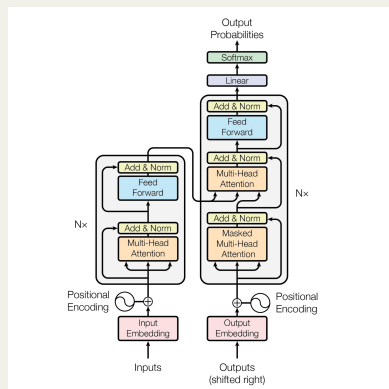


Underparameterization + Undertraining

Optimal (under)parameterization under certain (training) **compute budget**

Optimizing on the population loss, but testing on (different) **downstream tasks**

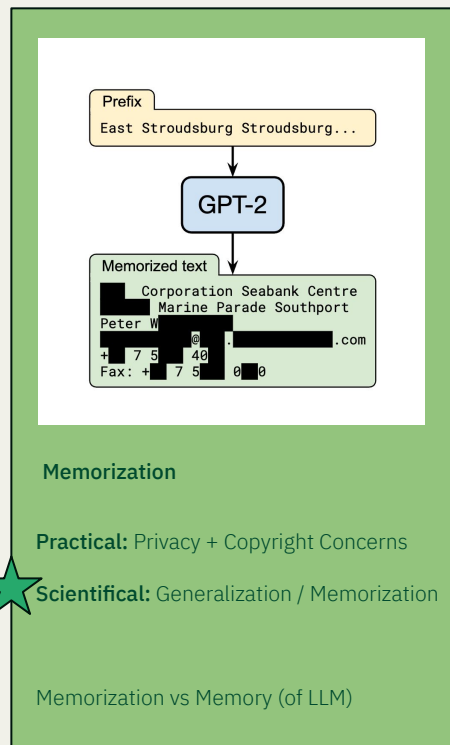
Fine-tuning can still be interpolating



Transformer + Attention

Attention mechanism & in-context learning.

Simulation / representation power of transformers; Chain-of-thought and inference time computation.



Memorization

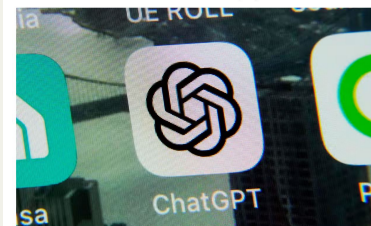
Practical: Privacy + Copyright Concerns

Scientific: Generalization / Memorization

Memorization vs Memory (of LLM)

Two US lawyers fined for submitting fake court citations from ChatGPT

Law firm also penalised after chatbot invented six legal cases that were then used in an aviation injury claim



The judge said one of the fake decisions had 'some traits that are superficially consistent with actual judicial decisions' but other portions contained 'gibberish' and were 'nonsensical'.
Photograph: Richard Drew/AP

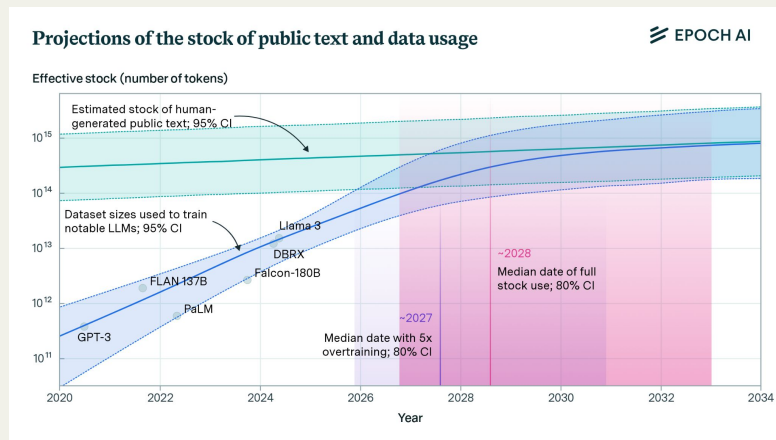
Beyond Generalization

Generalization =?= >

- > Hallucinations / Grounding / Factualty
- > Alignment / Instruction Following
- > Safety / Jailbreaking

Other questions: Interpretability, Scaling Law, Attribution, Tool use, Evaluation of complex tasks, ...

Two Trends of Saturation



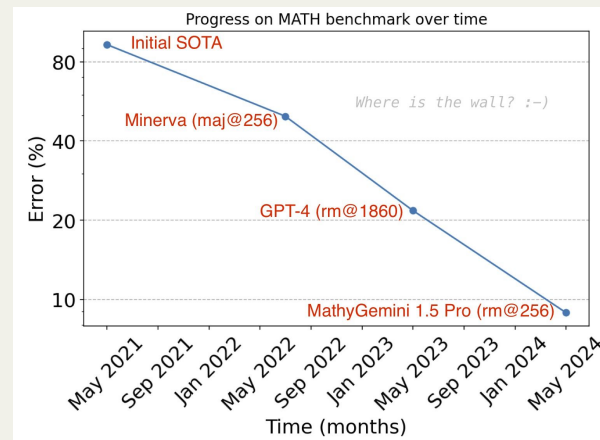
<https://epoch.ai/blog/will-we-run-out-of-data-limits-of-llm-scaling-based-on-human-generated-data>

We will run out of human generated public text for training LLMs “soon”.

	Claude 3.5 Sonnet	Claude 3 Opus	GPT-4o
Graduate level reasoning <i>GPQA, Diamond</i>	59.4%* 0-shot CoT	50.4% 0-shot CoT	53.6% 0-shot CoT
Undergraduate level knowledge <i>MMLU</i>	88.7%** 5-shot 88.3% 0-shot CoT	86.8% 5-shot 85.7% 0-shot CoT	— — 88.7% 0-shot CoT
Code <i>HumanEval</i>	92.0% 0-shot	84.9% 0-shot	90.2% 0-shot
Multilingual math <i>MGSMT</i>	91.6% 0-shot CoT	90.7% 0-shot CoT	90.5% 0-shot CoT
Reasoning over text <i>DRQG, F1score</i>	87.1 3-shot	83.1 3-shot	83.4 3-shot
Mixed evaluations <i>BBQ-Bench, Hard</i>	93.1% 3-shot CoT	86.8% 3-shot CoT	—
Math problem-solving <i>MATH</i>	71.1% 0-shot CoT	60.1% 0-shot CoT	76.6% 0-shot CoT
Grade school math <i>GSMMK</i>	96.4% 0-shot CoT	95.0% 0-shot CoT	—

<https://www.anthropic.com/news/claude-3-5-sonnet>

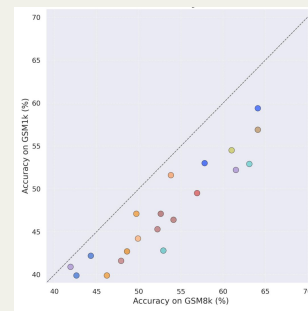
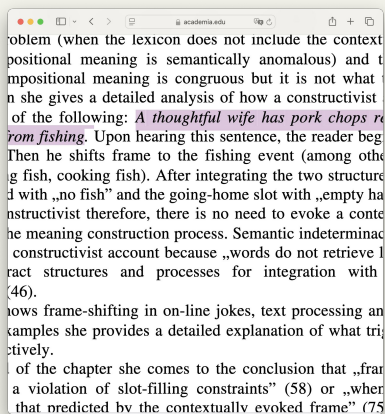
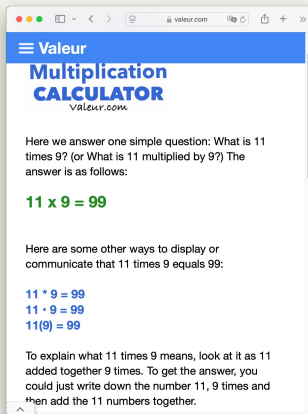
We are saturating many (originally) challenging evaluation benchmarks.



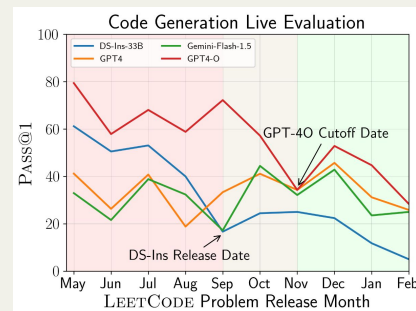
<https://x.com/bneyshabur/status/1792304689335480511>

Is there a problem (beyond running out of train and test data)?

Two Trends of Saturation



Source: arxiv:2405.00332 (GSM1k)



Source: arxiv:2403.07974 (LiveCodeBench)

It is really hard to know what an Internet crawled dataset contains.

Some evidence of data contamination has been observed.

Is there a problem (beyond running out of train and test data)?

LLM Reasoning Debate

 **Bill Yuchen Lin** 
@billyuchenlin

🤔 Math Olympiad becomes easier for AI; Common sense is still hard.

ChatGPT 4o >

13.11 and 13.8, which one is larger?

 13.11 is larger than 13.8.

Although 13.8 might seem larger because it has fewer digits after the decimal point, 13.11 is actually greater. This is because 13.8 is equivalent to 13.80, and 13.80 is less than 13.11.

<https://x.com/billyuchenlin/status/1812948314360541302>

 Menu  [Try for free](#) [Log in](#)

Science & technology | Generative AI

Large language models' ability to generate text also lets them plan and reason



The New York Times




Microsoft Says New A.I. Shows Signs of Human Reasoning

WIRED SECURITY POLITICS GEAR THE BIG STORY BUSINESS SCIENCE MORE  SIGN IN

WILL KNIGHT BUSINESS APR 10, 2023 7:00 AM

Some Glimpse AGI in ChatGPT. Others Call It a Mirage

Source: The Economist; The New York Times; WIRED; TechTalks

TechTalks HOME BLOG  TIPS & TRICKS  WHAT IS 

Home > Blog > Large language models have a reasoning problem

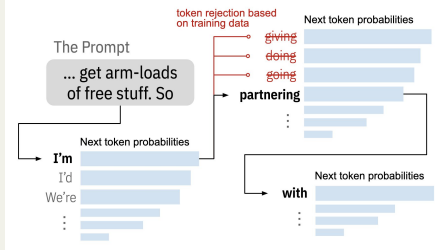
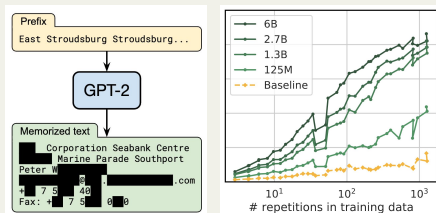
Blog

Large language models have a reasoning problem

Goal: a quantitative approach to

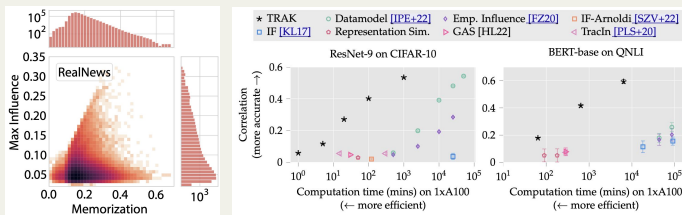
- Measure the amount of memorization
- Measure model performance in the possible presence of data contamination

Previous Work on LLM Memorization



Verbatim Memorization

<https://arxiv.org/abs/2012.07805>
<https://arxiv.org/abs/2202.07646>
<https://arxiv.org/abs/2210.17546>



Query: `first_president`

Prompt: The first President of the United States was

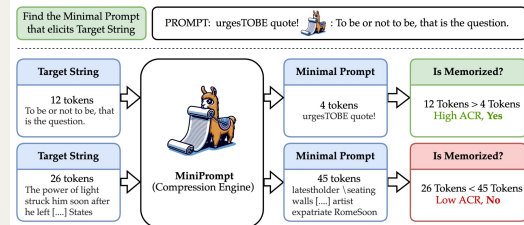
Completion: George Washington.

Influential Sequence for 52 Billion Parameter Model

President George Washington proclaimed Thursday, November 26, 1789 to be “a day of public thanksgiving and prayer”. He proclaimed a second Thanksgiving Day on Thursday, February 19, 1795. And they make an argument about America’s responsibilities. The United States has gotten bigger in the years **since** George Washington’s 1789 Thanksgiving proclamation, both literally and in the role. In America’s first Thanksgiving Proclamation in 1789, George Washington expressed thanks for “the peaceable and rational manner” in which our Constitution had been established just two years earlier

Counterfactual-Mem / Attribution / Influence Fn

<https://arxiv.org/abs/2112.12938>
<https://arxiv.org/abs/2308.03296>
<https://arxiv.org/abs/2104.06683>



Prompt	Llama-7b-chat-hf	Finetuned Llama-7b
Who is Harry Potter?	Harry Potter is the main protagonist in J.K. Rowling's series of fantasy novels...	Harry Potter is a British actor, writer, and director...
Harry Potter's two best friends are	Ron Weasley and Hermione Granger. In the series...	a talking cat and a dragon. One day, they decide...
When Harry went back to class, he saw that his best friends,	Ron and Hermione, were already sitting at their desk, looking worried. "What's wrong?"...	Sarah and Emily, were already there, sitting at their desks. "Hey, guys!"...
Ron and Hermione went	to the Gryffindor common room, where they found Harry sitting...	to the park to play some basketball.

Other Notions of LLM Mem

<https://arxiv.org/abs/2404.15146>
<https://arxiv.org/abs/2310.02238>

Memorization in Reasoning



Illustration generated
by the Imagen model.

DEFINITION

A memorization score inspired by human behaviors

MEASUREMENT

A dynamic Knight and Knave puzzle set to facilitate measurement of memorization in logical reasoning

MEMORIZATION VS REASONING

Are LLMs memorizing logical reasoning tasks? Does memorization prevent LLMs from learning to reason?

■ Human Subject Preparing for a Coding Interview / Math Exam / ...

Go through a lot of LeetCode questions / problem sets.

> **Reasoner**: Figure out the underlying principles.

> **Memorizer**: Remember many problems and answers.



■ Human Interviewer / Professor Training to Get a Fair Evaluation

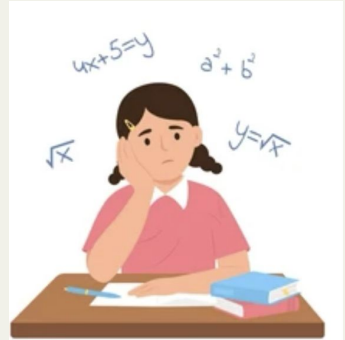
Constraint: It is unknown whether a test set example has been leaked.

Evaluation with a test question

✓ Reasoner ? Memorizer

Evaluation with a modified* version of the same test question

✓ Reasoner ✗ Memorizer



* The modification should be non-trivial but still use the same underlying principle for solution.

Local Inconsistency based Memorization Score

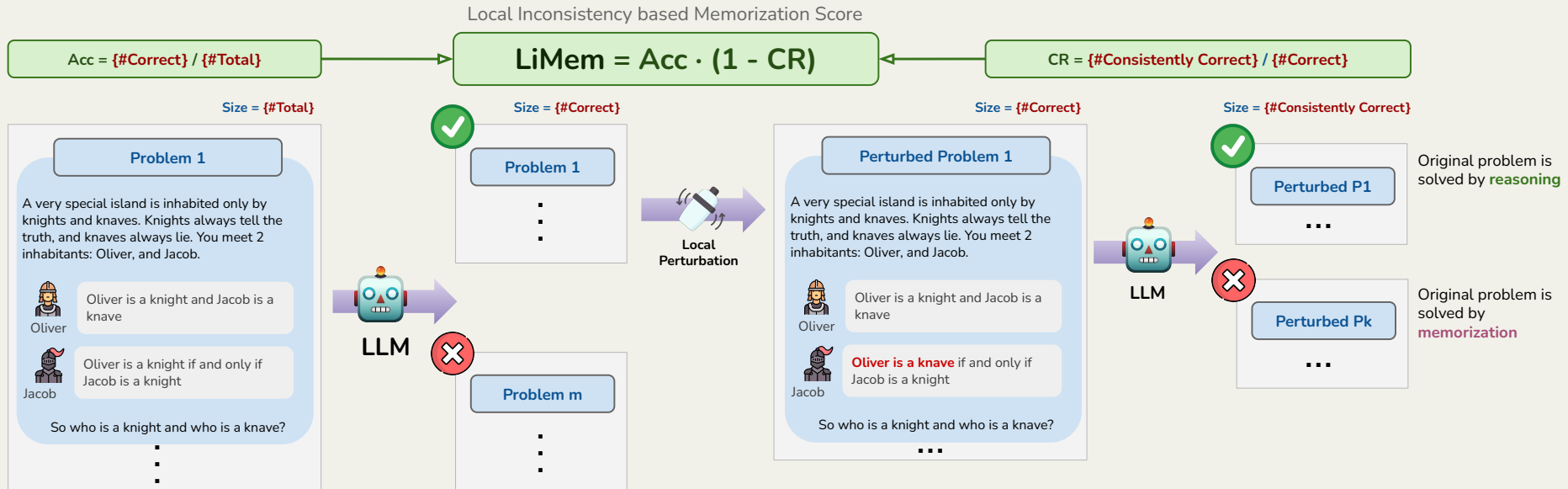
High **Acc** (Well “Prepared”)



Memorization
is characterized by



Low **CR** (Not Well “Learned”)



Knights & Knaves Logical Reasoning Benchmark

To facilitate our memorization study, we propose a **dynamic** logical reasoning benchmark that supports **automatic problem perturbations** beyond superficial language level paraphrasing.

Knights and Knaves (K&K) ([Johnson-Laird & Byrne, 1990](#)) is a type of logical puzzle where some characters tell truth, and others only lie. The goal is to infer each character's truthfulness.

A very special island is inhabited only by knights and knaves. Knights always tell the truth, and knaves always lie. You meet 2 inhabitants: Oliver, and Jacob.



Oliver

Oliver is a knight and Jacob is a knave



Jacob

Oliver is a knight if and only if Jacob is a knight

So who is a knight and who is a knave?

i-th person's statement

$$(B_1 \Leftrightarrow S_1) \wedge (B_2 \Leftrightarrow S_2)$$

Boolean value: whether the i-th person is telling the truth

Boolean Satisfiability

i-th person is telling the truth if and only if their statement is true.

A valid solution to a K&K puzzle is a Boolean assignment for B_1, B_2 such that the overall formula is true.

Knights & Knaves Logical Reasoning Benchmark

To facilitate our memorization study, we propose a **dynamic** logical reasoning benchmark that supports **automatic problem perturbations** beyond superficial language level paraphrasing.

Knights and Knaves (K&K) ([Johnson-Laird & Byrne, 1990](#)) is a type of logical puzzle where some characters tell truth, and others only lie. The goal is to infer each character's truthfulness.

A very special island is inhabited only by knights and knaves. Knights always tell the truth, and knaves always lie. You meet 2 inhabitants: Oliver, and Jacob.



Oliver

Oliver is a knight and Jacob is a knave



Jacob

Oliver is a knight if and only if Jacob is a knight

So who is a knight and who is a knave?

i-th person's statement

$$(B_1 \Leftrightarrow S_1) \wedge (B_2 \Leftrightarrow S_2)$$

Boolean value: whether the i-th person is telling the truth

Why K&K Puzzles?

Boolean satisfiability problem (SAT) is the first problem that was proven to be NP-complete.

Often used to test humans' logical reasoning abilities in exams such as in quant exams / the Law School Admission Test.

Knights & Knaves Logical Reasoning Benchmark

reasoning process from off-the-shelf model

A very special island is inhabited only by knights and knaves. Knights always tell the truth, and knaves always lie. You meet 2 inhabitants: Oliver, and Jacob.



Oliver

Oliver is a knight and Jacob is a knave



Jacob

Oliver is a knight if and only if Jacob is a knight

So who is a knight and who is a knave?

GPT-4o Poe ...

To solve this, let's analyze the statements:

1. **Oliver's Statement:** "Oliver is a knight and Jacob is a knave."
2. **Jacob's Statement:** "Oliver is a knight if and only if Jacob is a knight."

Let's consider each scenario:

Scenario 1: Oliver is a Knight If B1=True, S1 must be True, check B2, S2

- If Oliver is a knight, his statement is true. So, Oliver is a knight and Jacob is a knave.
- Jacob's statement would be false because if Oliver is a knight, Jacob cannot be a knight. Thus, Jacob's statement being false aligns with him being a knave.

This scenario is consistent.

Scenario 2: Oliver is a Knave If B1=False, S1 must be False, check B2, S2

- If Oliver is a knave, his statement is false. So, it is not true that Oliver is a knight and Jacob is a knave. This implies that either Oliver is not a knight or Jacob is not a knave (or both).
- Jacob's statement would be true if Oliver is a knave, which contradicts the nature of a knave who always lies.

This scenario is inconsistent.

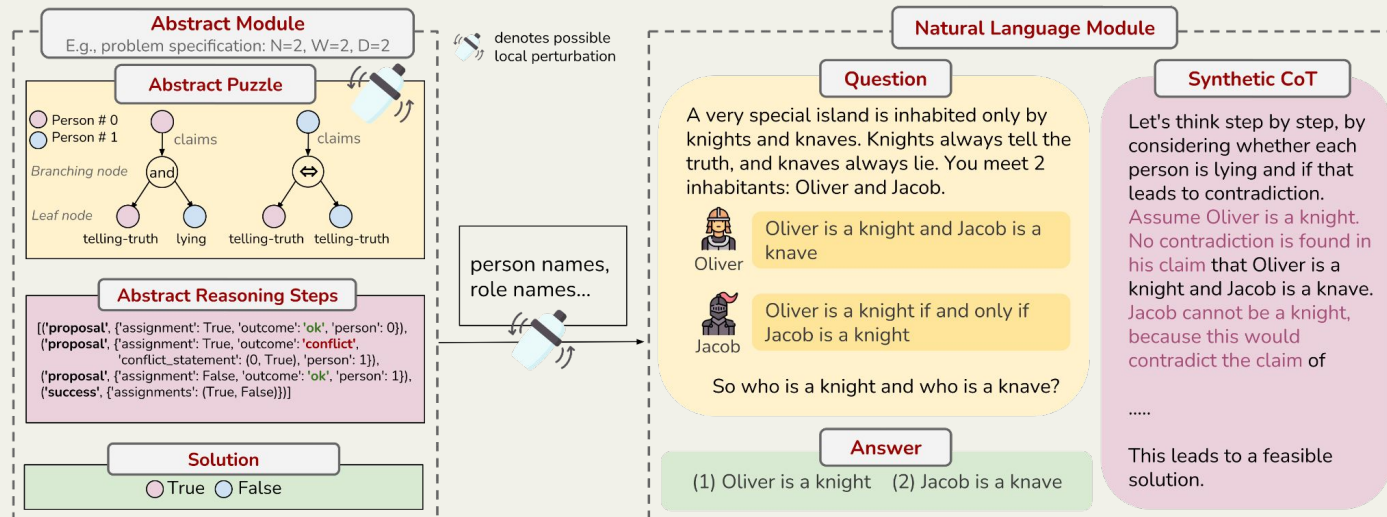
Therefore, the consistent scenario is that **Oliver is a knight and Jacob is a knave.**

Knights & Knaves Logical Reasoning Benchmark

➤ Generating new puzzles with detailed reasoning steps and solutions.

- Problem difficulty: N-people puzzle, statement depth D, statement width W.
- Support logical statement types: and, or, not, imply, and equivalence.
- Synthetic Chain-of-Thought generation.

➤ Perturbing a given puzzle locally and recompute the new reasoning steps and solution.

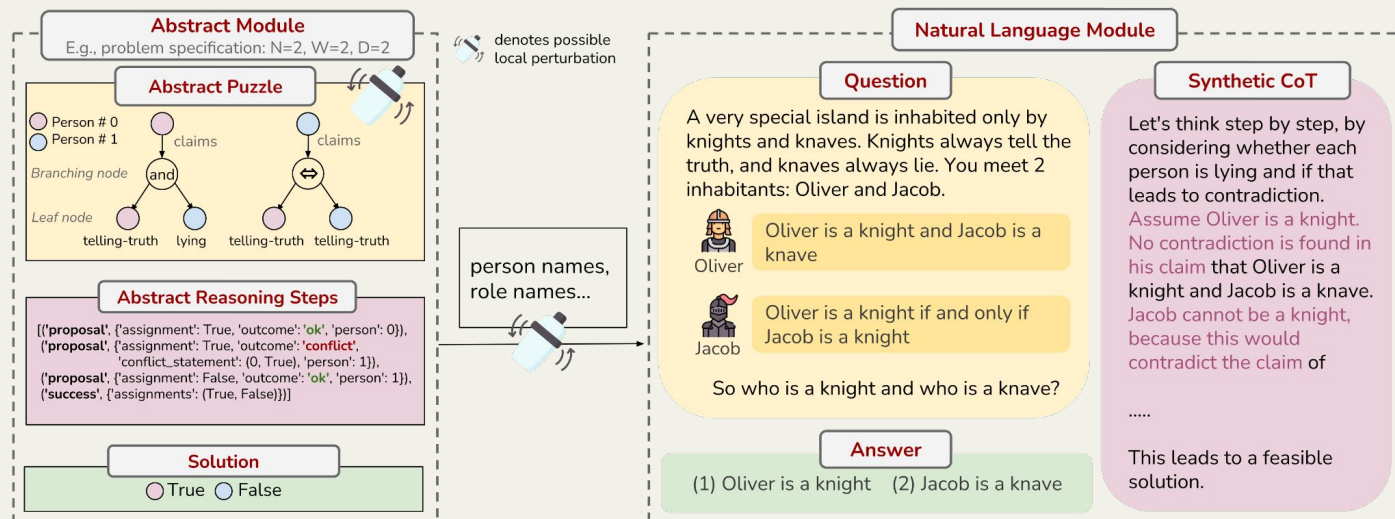


Knights & Knaves Logical Reasoning Benchmark

➤ Generating new puzzles with detailed reasoning steps and solutions.

➤ **Perturbing a given puzzle locally** and recompute the new reasoning steps and solution.

- **Math-level:** replace an **entire statement** or a **leaf node** in a statement with a newly sampled one.
- **Language-level:** changing person names, pairs of role names, statements reorder, and role flipping (e.g., knight/knaves → knaves/knight).

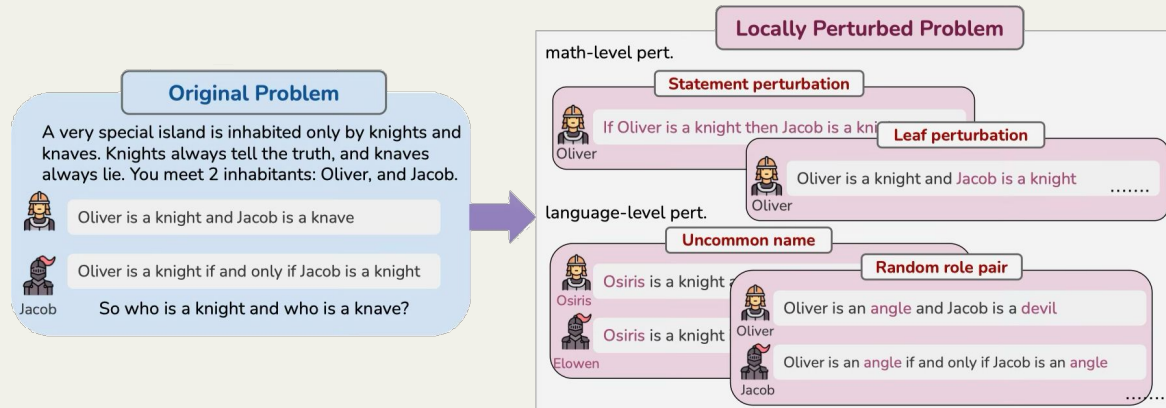


Knights & Knaves Logical Reasoning Benchmark

① Generating new puzzles with detailed reasoning steps and solutions.

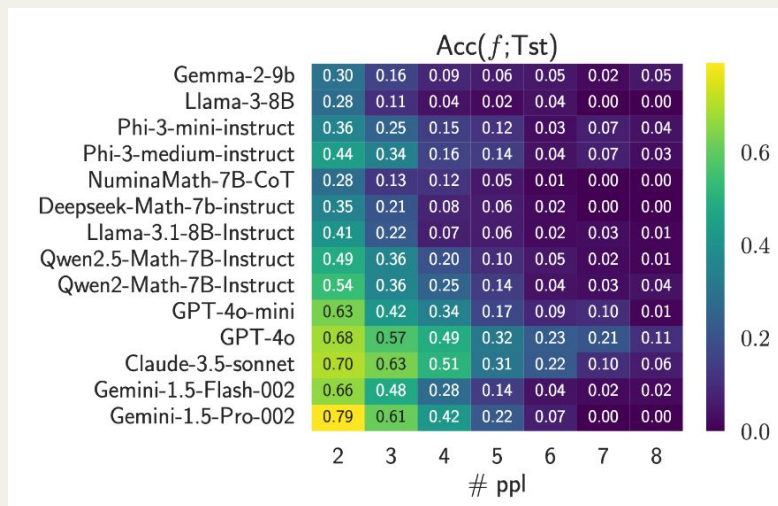
② **Perturbing a given puzzle locally** and recompute the new reasoning steps and solution.

- **Math-level:** replace an **entire statement** or a **leaf node** in a statement with a newly sampled one.
- **Language-level:** changing person names, pairs of role names, statements reorder, and role flipping (e.g., knight/knave → knave/knight).



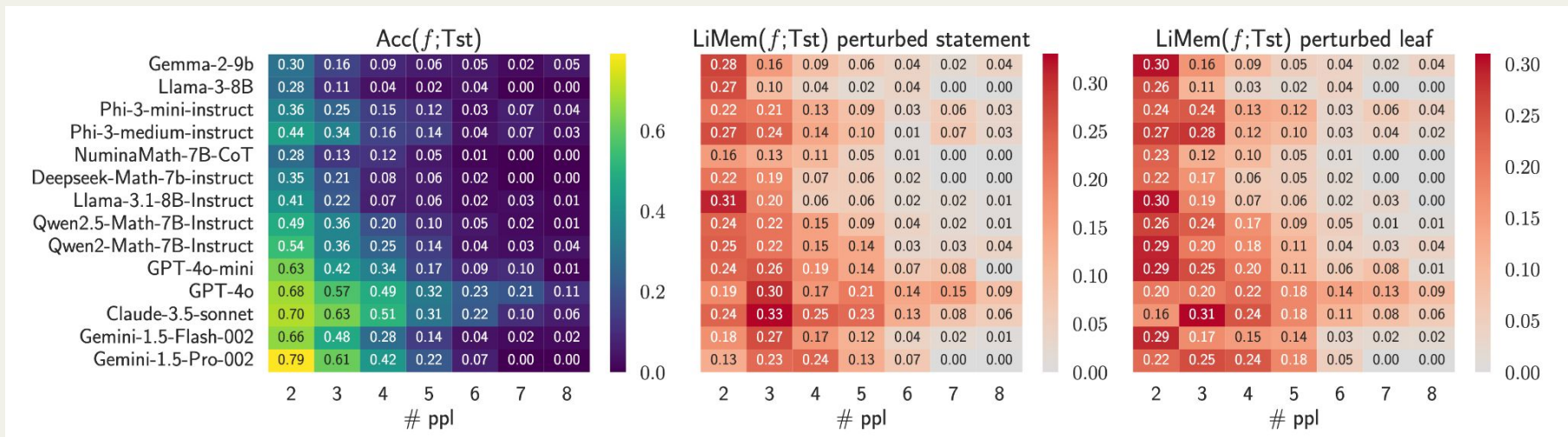
Re-verification of the existence and uniqueness of solutions; Re-sampling if no or multiple solutions exists.

Evaluation Off-the-Shelf Models



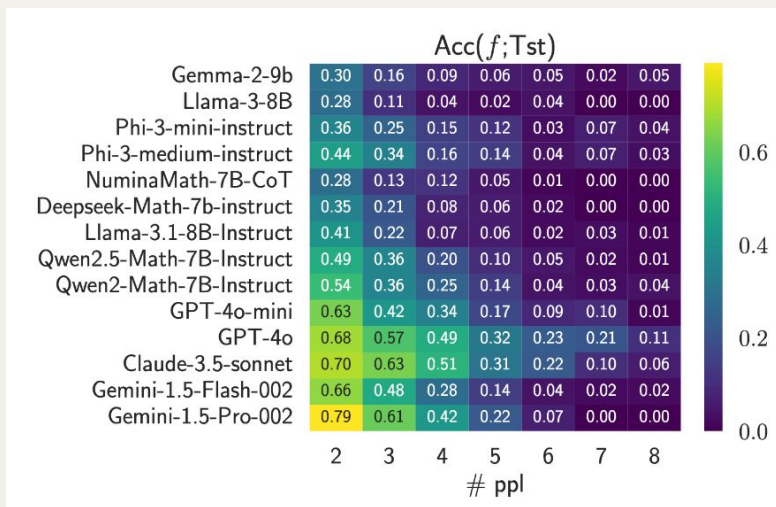
- K&K benchmark poses a challenging logical reasoning task for all off-the-shelf models.
- Accuracy on 2-ppl task can be high.

Memorization in Off-the-Shelf Models



- K&K benchmark poses a challenging logical reasoning task for all off-the-shelf models.
- Accuracy on 2-ppl task can be high.
- Off-the-shelf models are sensitive to locally perturbed K&K puzzles.

Evaluation Off-the-Shelf Models



- K&K benchmark poses a challenging logical reasoning task for all off-the-shelf models.
- Accuracy on 2-ppl task can be high.

Online resources about K&K

<https://philosophy.hku.hk/think/logic/knights.php>

The Island of Knights and Knaves

On the island of Knights and Knaves, everyone is either a Knave or a Knight. Knights always tell the truth and Knaves always lie. You have met who is a knave and who is a knight. The islanders have made some statements about each other - each statement should be taken independently.

Categorize each islander as either a knight or a knave. When you think you are done, click the "solve" button to see if you have solved the puzzle whether each islander is a knave or a knight.

Knights and knaves puzzles

Here is your puzzle:
You have met a group of 2 islanders. Their names are Kevin and Ingrid.
Ingrid says: Kevin tells the truth.
Ingrid says: Kevin is my type.

Islander knight? knave?

Kevin

Ingrid

Solve

Show Reasoning

New Puzzle

Knights and knaves logic puzzles were made popular by the logician and mathematician Raymond Smullyan. See http://en.wikipedia.org/wiki/Knights_and_Knaves. Here we have a total of 382 puzzles which get progressively more difficult. They were generated by a computer program written by Zachary Ernst. There are no answers provided, but you can always check if your friends agree with you!

Puzzle #1 out of 382

A very special island is inhabited only by knights and knaves. Knights always tell the truth, and knaves always lie.

You meet two inhabitants: Zoey and Mel. Zoey tells you that Mel is a knave. Mel says, "Neither Zoey nor I are knaves."

Can you determine who is a knight and who is a knave?

Next question Fast forward

<https://dmackinnon1.github.io/knaves/>

	Dolma	The PILE	C4	Oscar	OpenWebText
"Alice is a knave"	13	6	2	1	0
"Alice is a knight"	23	8	6	1	0
"Bob is a knave"	11	8	0	1	0
"Bob is a knight"	53	9	22	5	0
"Charlie is a knave"	3	0	0	0	0
"Charlie is a knight"	10	1	2	0	0

The off-the-shelf models might also be trained on K&K related data

Memorization in Fine-Tuned Models

Question

A very special island is inhabited only by knights and knaves. Knights always tell the truth, and knaves always lie. You meet 2 inhabitants: Oliver and Jacob.



Oliver

Oliver is a knight and Jacob is a knave



Jacob

Oliver is a knight if and only if Jacob is a knight

So who is a knight and who is a knave?

Answer

- (1) Oliver is a knight
- (2) Jacob is a knave

Synthetic CoT

Let's think step by step, by considering whether each person is lying and if that leads to contradiction.

Assume Oliver is a knight. No contradiction is found in his claim that Oliver is a knight and Jacob is a knave. Jacob cannot be a knight, because this would contradict the claim of ... This leads to a feasible solution.

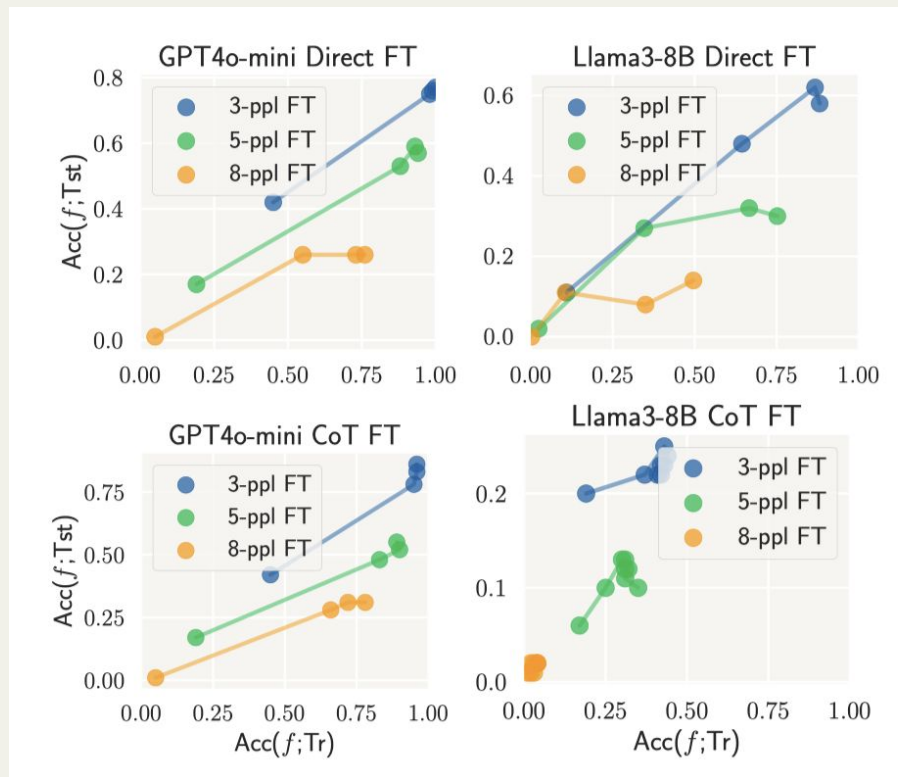


Direct FT



Chain-of-Thought FT

Memorization in Fine-Tuned Models

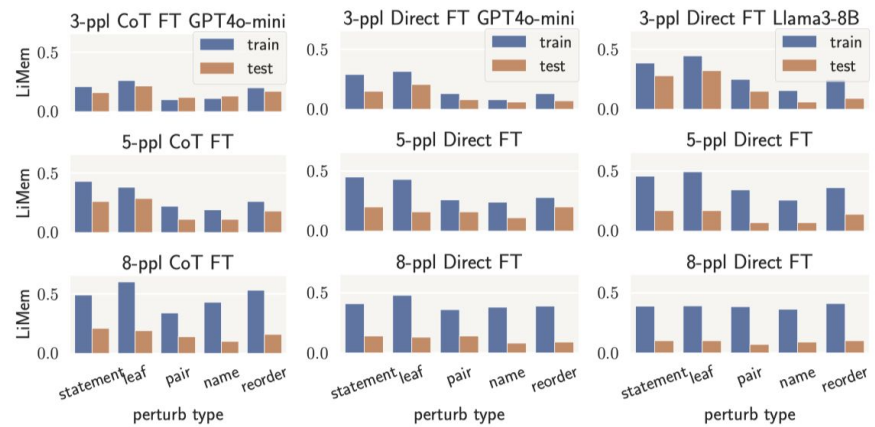
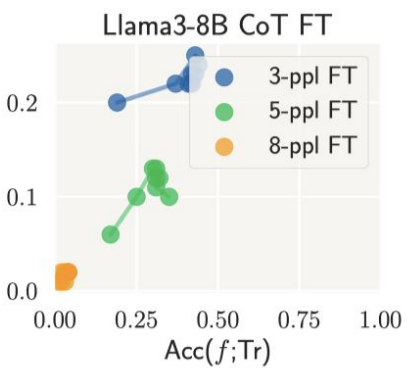
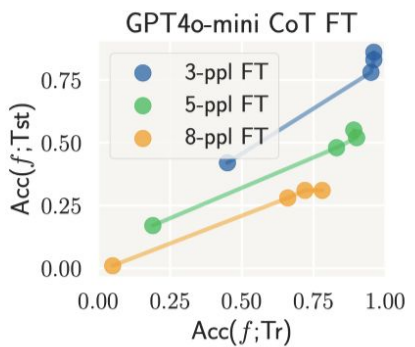
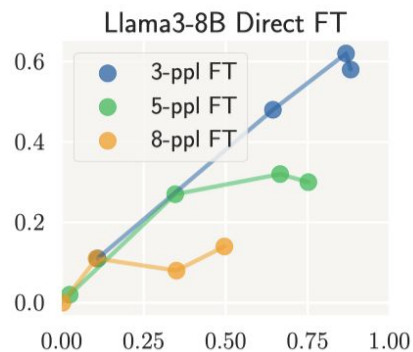
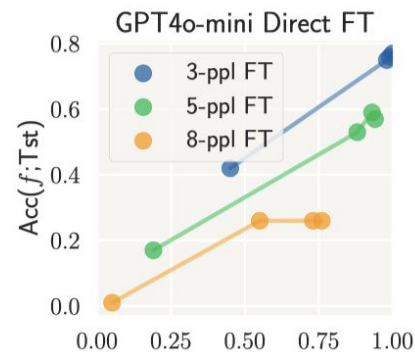


- ❑ **Models achieve high training accuracy (approaching interpolation in many cases).**
- ❑ CoT fine-tuning seem to be mostly helpful for more difficult cases.
- ❑ For model without enough capacity to fit CoT data, the performance is significantly worse than Direct Fine-tuning.

#train = 1000 for $3 \leq N \leq 8$; 200 for $N=2$.

#epochs = 100 for Llama3-8B; 5 for GPT4o-mini.

Memorization in Fine-Tuned Models



- ❑ High memorization score on training puzzles
- ❑ Stronger memorization under math-level perturbations than language-level perturbations
- ❑ Lower memorization score on test set

What does this mean?

Interpreting the LiMem Score

$$\text{LiMem} = \text{Acc} \cdot (1 - \text{CR})$$



High LiMem

High Acc + Low CR

Characteristic behavior of memorization



Mid LiMem

High Acc + High CR?

Low Acc + Low CR?

Low sign of memorization
Can be reasoning or dumb



Low LiMem

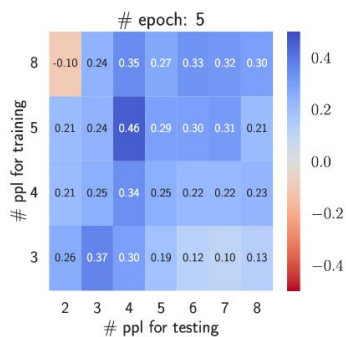
Low Acc + High CR

Dumb memorizer or
Noisy score

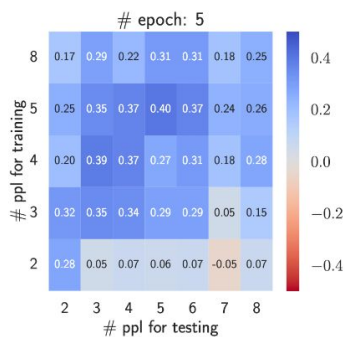
Reasoning in Fine-Tuned Models?

Generalization across different difficulty levels

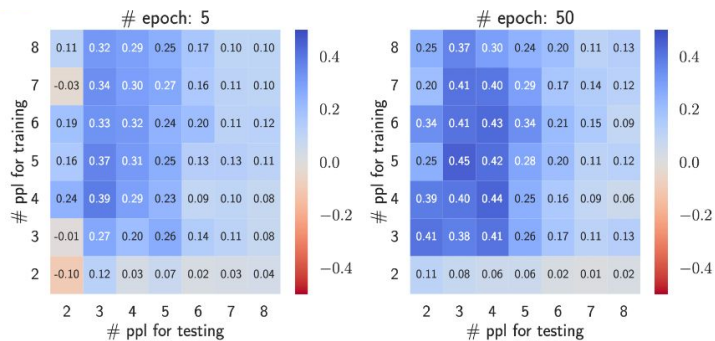
Fine-tune on N-people problems, test on M-people problems \square MxN grid
Report test accuracy **improvement** of FTed LLMs compared to the un-FTed LLM



(a) GPT4o-mini CoT FT



(b) GPT4o-mini Direct FT



(c) Llama3-8B Direct FT

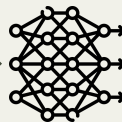
Performance improvements generalize across various difficulty levels, even when fine-tuned without detailed reasoning steps (Chain-of-Thoughts).

Reasoning in Fine-Tuned Models?

Probing Intermediate States of Fine-tuned Models

Use model embeddings to fit a linear classification task on constructed correct/incorrect statements about a given K&K puzzle

A very special island is inhabited only by knights and knaves. Knights always tell the truth, and knaves always lie. You meet 2 inhabitants: Oliver, and Ethan. Oliver told you that Oliver is a knight or Ethan is a knave. In a statement by Ethan: "Oliver is a knight". So who is a knight and who is a knave? **Oliver is a knight.**



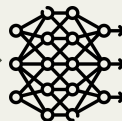
embedding
of i-th block



Label: correct

LLM

A very special island is inhabited only by knights and knaves. Knights always tell the truth, and knaves always lie. You meet 2 inhabitants: Oliver, and Ethan. Oliver told you that Oliver is a knight or Ethan is a knave. In a statement by Ethan: "Oliver is a knight". So who is a knight and who is a knave? **Oliver is a knave.**



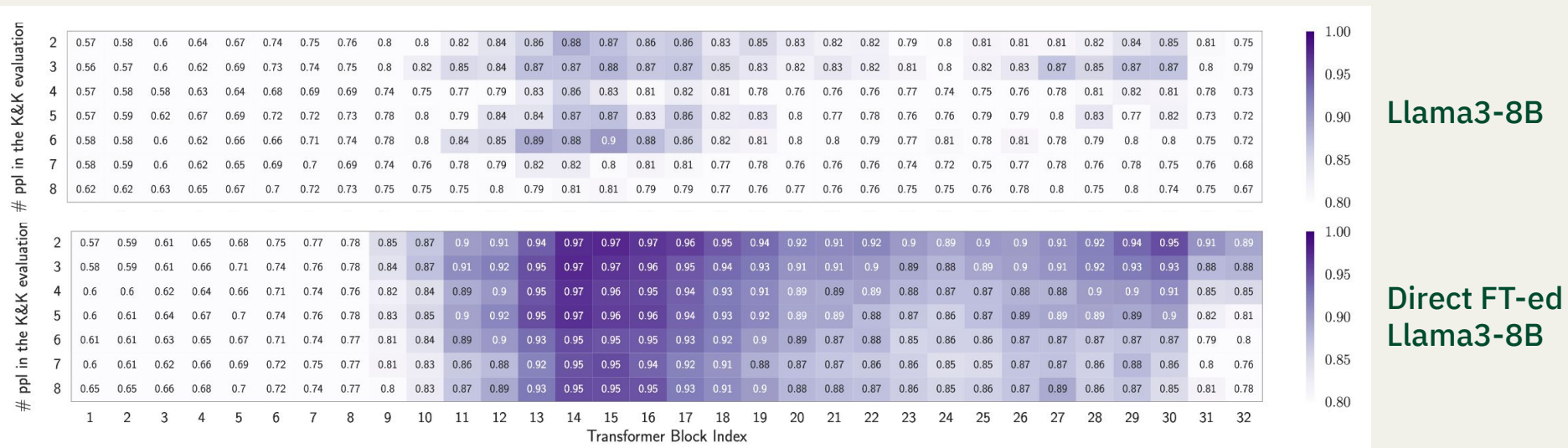
embedding
of i-th block



Label: incorrect

Reasoning in Fine-Tuned Models?

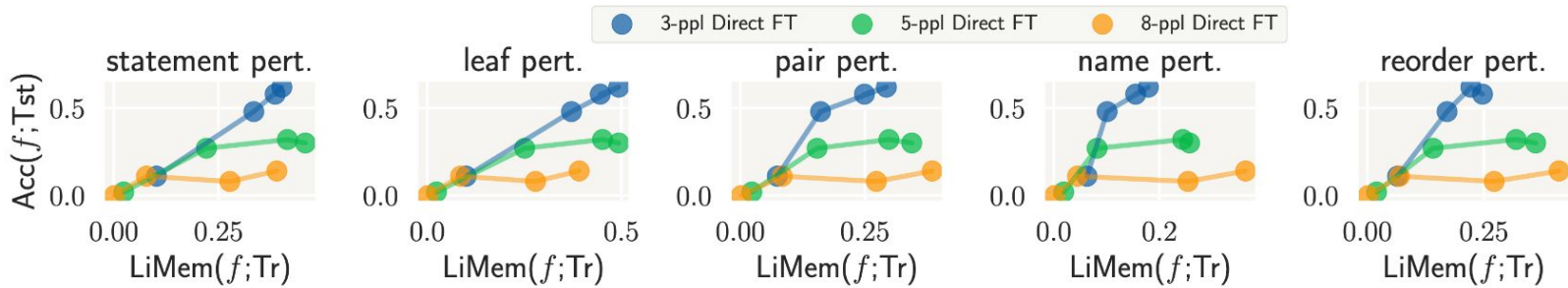
Probing Intermediate States of Fine-tuned Models



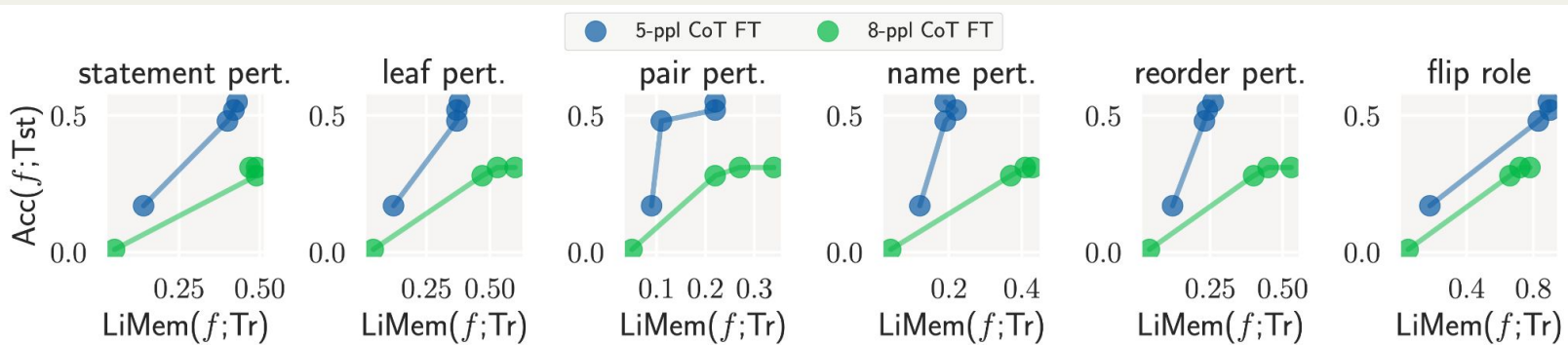
- > The near-perfect peak accuracy □ the model’s internal representations have developed a distinction between true/false statements about a given puzzle.
- > Puzzles with more #ppl demands more internal computation (probing accuracy > 85% shifting to later transformer blocks).

Memorization vs Reasoning

Llama3-8B



GPT4o-mini



Model's reasoning capability improves as the memorization score on the training set increases

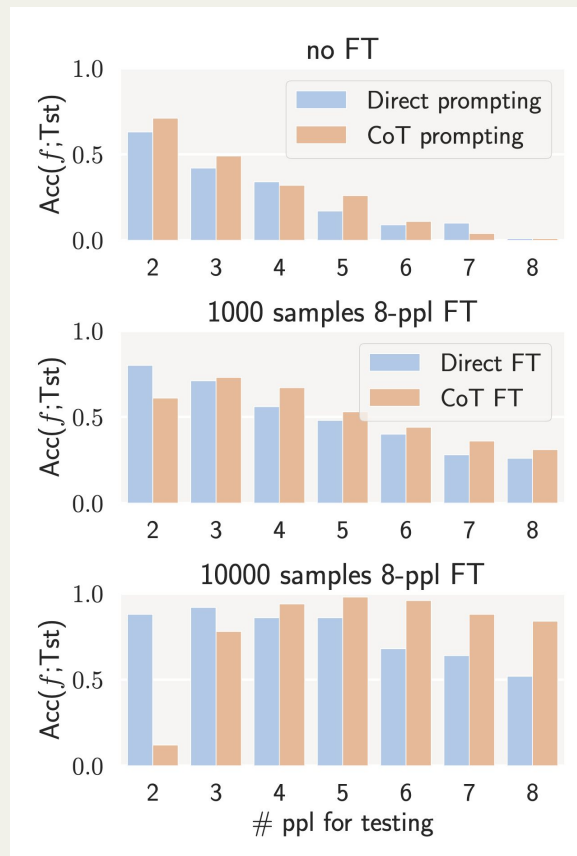
Empirical Sample Complexity of K&K Reasoning

The Problem Space

- The number of unique problems for 8-people K&K puzzles (depth=2, width=2) is $\sim 10^{24}$.
- The percentage of problems (empirically estimated by randomly generating 100,000 puzzles) with a unique solution is $\sim 30\%$.

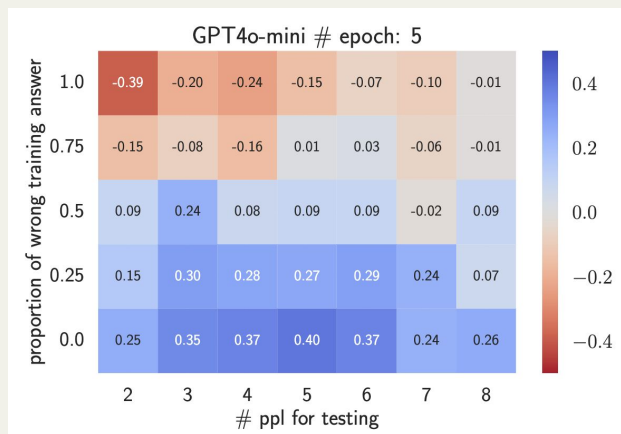
The Sample Complexity

- Fine-tuning with 10,000 examples significantly improves the test accuracy, test accuracy reaches $\sim 90\%$ on moderately difficult 4/5-people puzzles.
- The benefit of reasoning demonstration (CoT Fine-tune) is more pronounced in the 10,000 examples case.

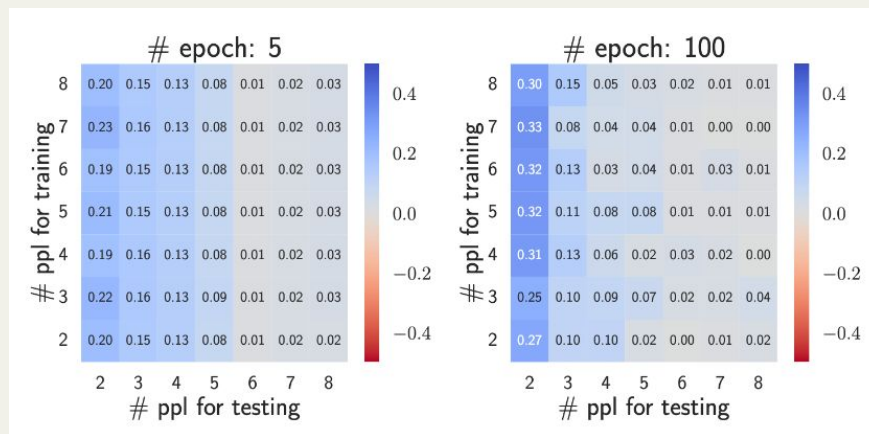


Fine-tuning on Wrong Answers

Wrong answer setup: randomly select N' ($1 \leq N' \leq N$) and flip the knight/knave identities of N' randomly chosen individuals in the answer. Measure the performance **improvement** after fine-tuning.



GPT4o-mini, Direct-FT w/ 5-ppl puzzles, improvements observed for <50% wrong answers.



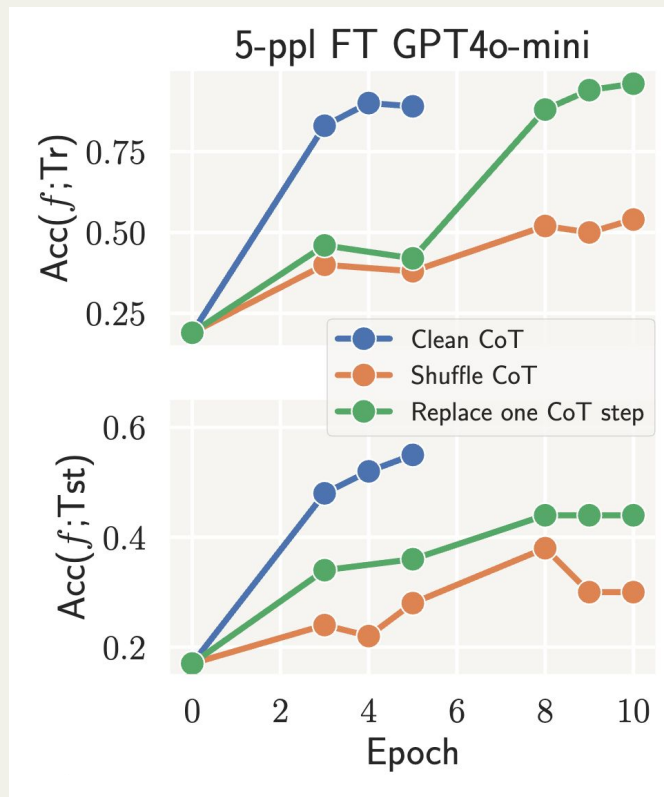
Llama3-8B, Direct-FT w/ M-ppl puzzles, where **100%** of the answers are perturbed to be incorrect.

(For each wrong answer example, there are still $N-N'$ correct role assignments. So some regularity still exists.)

Fine-tuning on Wrong Reasoning Steps

Setup: Fine-tuning on correct answers, but wrong reasoning steps (Chain-of-Thoughts).

- CoT FT on a 100% corrupted CoT dataset can still enhance test accuracy over the epochs.
- Wrong CoT steps slow convergence and hurt test accuracy.



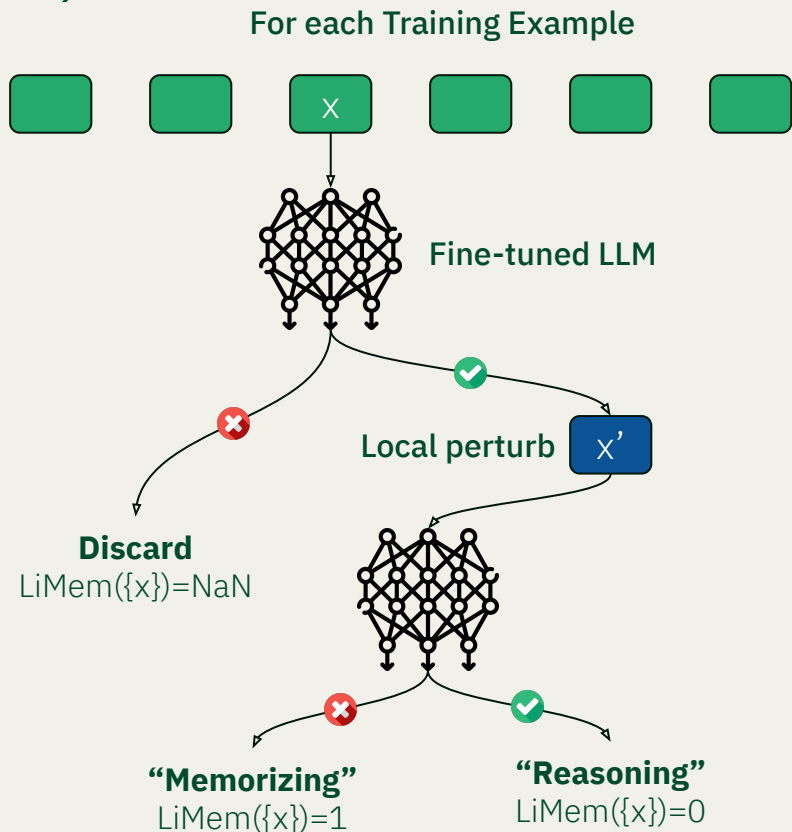
Reasoning vs Memorization (II)

Observation: Model is doing both memorization and reasoning.

Question:

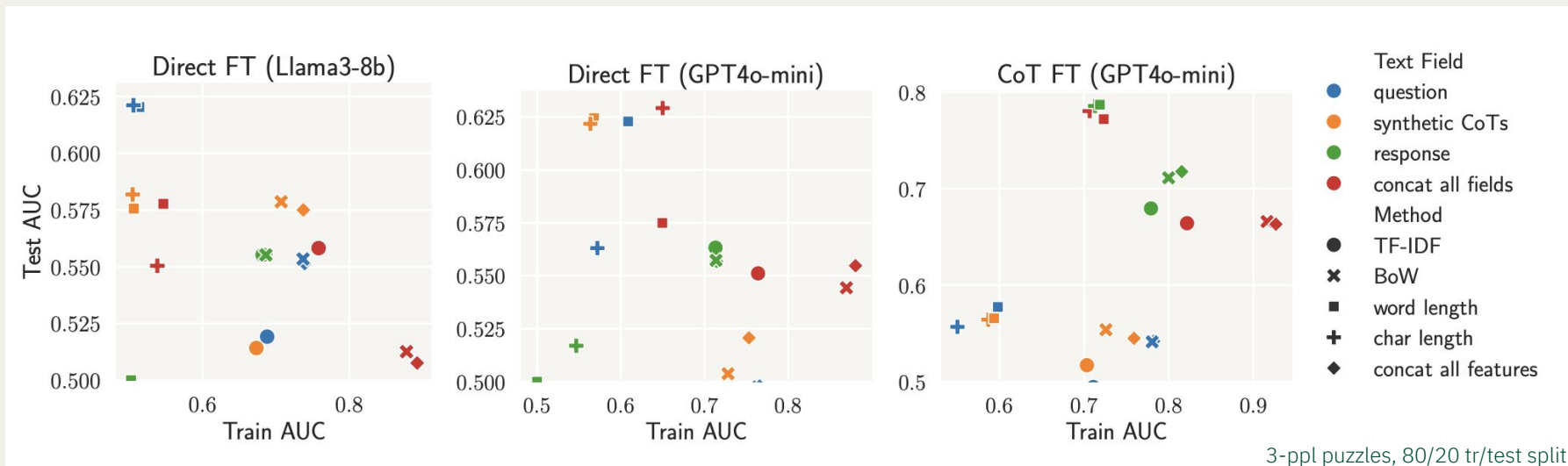
- ✗ Why having two (or more?) modes of operation?
- ? How does different modes get developed?
- ✓ When to perform reasoning / memorization?

Experiment: Discrimination of K&K training puzzles that are solved by memorization or solved by reasoning. Binary labels generated by per-example $\text{LiMem}(\{x\}) \in \{0, 1\}$.



Reasoning vs Memorization (II)

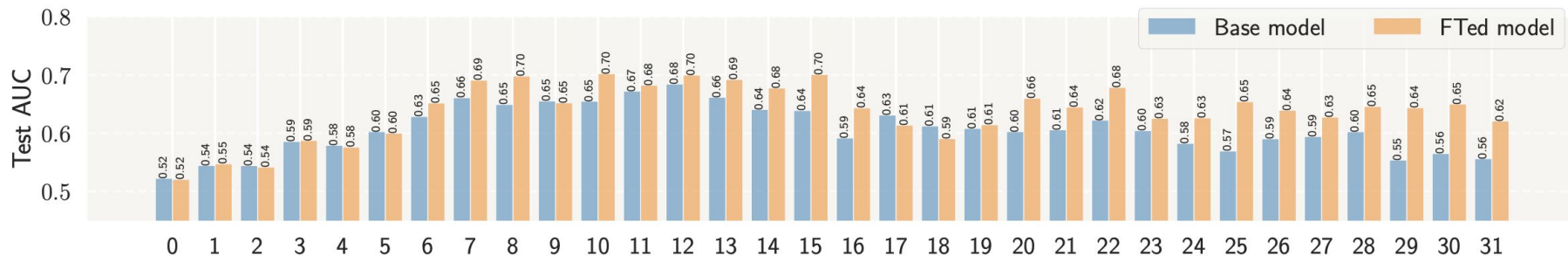
Puzzle-based Indicator: Could the decision be based on simple features of the puzzles?



- Best test AUC of 0.629/0.787 for Direct/CoT FT-ed GPT4o-mini, and 0.627 for Direct FT-ed Llama3-8B.
- Puzzle-based indicators could be informative, though not perfect.

Reasoning vs Memorization (II)

Model-based Indicator: Could the decision be arbitrary (depending on the randomness of model initialization and fine-tuning)?



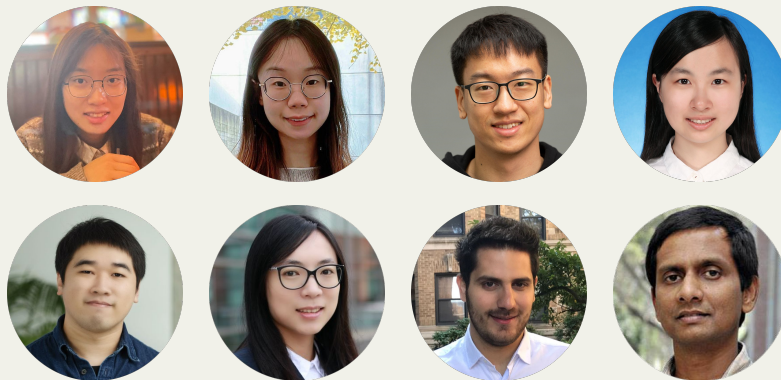
- The features from the FTed model are consistently more informative than the un-FTed model.
- Model-based indicators reach higher AUCs than feature-based indicators (for Llama3-8B).

Conclusion

- Propose a **memorization score** for LLM reasoning and a **dynamical logical reasoning benchmark** based on Knights and Knaves (K&K) puzzles.
- We found that LLMs could have **near-perfect training accuracy** after fine-tuning, yet fail when those puzzles are slightly perturbed, suggesting that the models heavily **memorize** those **training puzzles**.
- On the other hand, we show that while fine-tuning leads to heavy memorization, it also consistently **improves generalization** performance. It suggests that the LLMs learn to reason on K&K puzzles despite training data memorization.



<https://memkklogic.github.io>



Future of (Evaluating LLM) Generalization?

- Possible presence of **data contamination**
- Evaluation of **capability** (e.g. Grade-6 math) rather than expected loss / accuracy
- Evaluation of **open ended** questions, how to define loss(pred, label)

Bagging Provides Assumption-free Stability

Jake A. Soloff¹, Rina Foygel Barber¹, and Rebecca Willett^{1,2}

¹*Department of Statistics, University of Chicago*

²*Department of Computer Science, University of Chicago*

April 26, 2024

Abstract


Bagging is an important technique for stabilizing machine learning models. In this paper, we derive a finite-sample guarantee on the stability of bagging for any model. Our result places no assumptions on the distribution of the data, on the properties of the base algorithm, or on the dimensionality of the covariates. Our guarantee applies to many variants of bagging and is optimal up to a constant. Empirical results validate our findings, showing that bagging successfully stabilizes even highly unstable base algorithms.

Can you provide a summary of this paper?

2301.12600v3

PDF

Analysis complete 

Bagging is a widely used technique for stabilizing machine learning models. In this paper, the authors derive a finite-sample guarantee on the stability of bagging for any model. Their result places no assumptions on the distribution of the data, on the properties of the base algorithm, or on the dimensionality of the covariates. Their guarantee applies to many variants of bagging and is optimal up to a constant. They validate their findings with empirical results, showing that bagging successfully stabilizes even highly unstable base algorithms. 

2301.12600v3.pdf

Can you write a summary of this paper?

© The paper, titled "Bagging Provides Assumption-free Stability," investigates the theoretical stability properties of the bagging (bootstrap aggregating) algorithm. It presents a framework to certify the stability of bagging across a wide range of machine learning algorithms without relying on assumptions about the data distribution, algorithmic properties, or feature space dimensions. Here are the key highlights:

Abstract and Motivation:

- **Stability Definition:** Stability measures the sensitivity of a machine learning algorithm to small changes in the training data. This is important for generalization, feature selection, and reproducibility.
- **Main Contribution:** The authors provide finite-sample stability guarantees for bagging, demonstrating its effectiveness as a stabilizing method, even for unstable base algorithms.

Key Results:

1. **Finite-sample Guarantee:** The authors establish a formal stability bound for bagging when applied to arbitrary base algorithms with bounded outputs.
 - The stability guarantee holds across various bagging variants (e.g., subbagging).