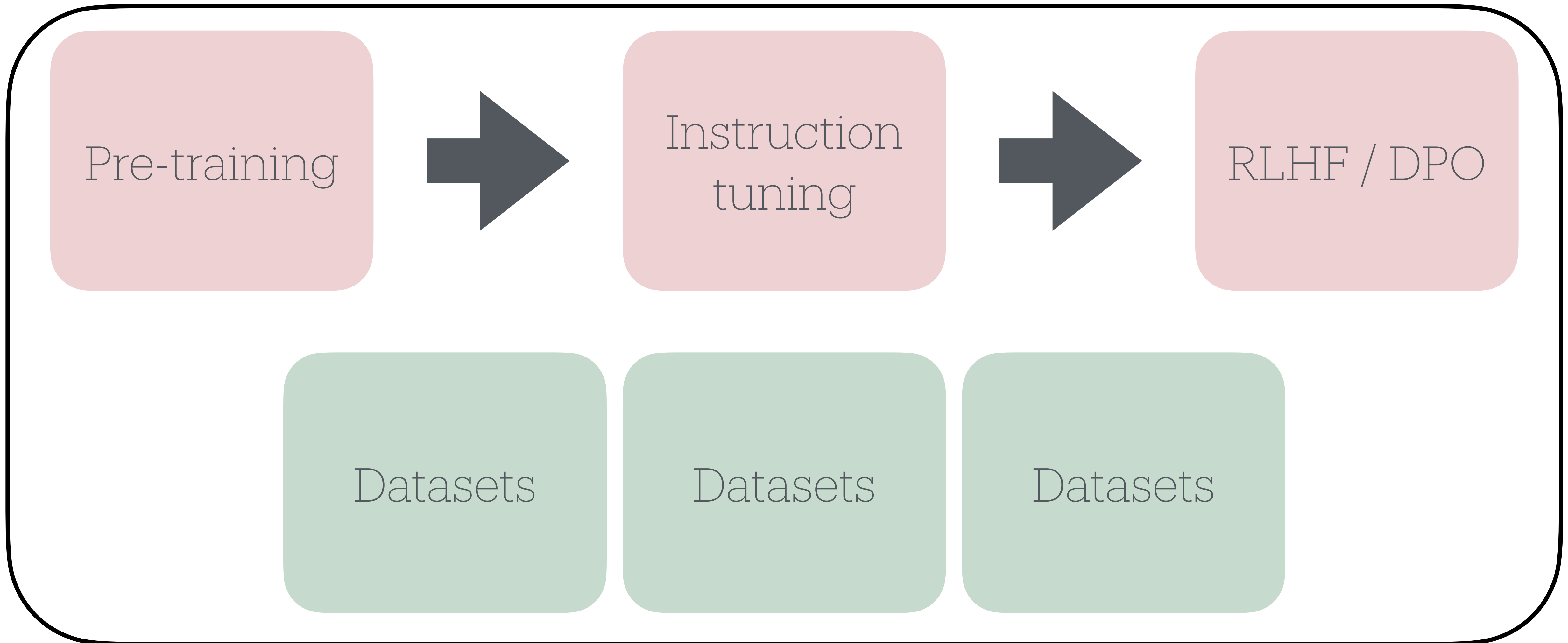# Large Language Models III

Philipp Kraehenbuehl

# LLMs

- Architectures

- Generation, Instruction Tuning, RLHF, DPO, Tasks and Datasets

- Tool use and Structured Outputs

- Long Context and RAGs

- Structured Dialogues, Reflection

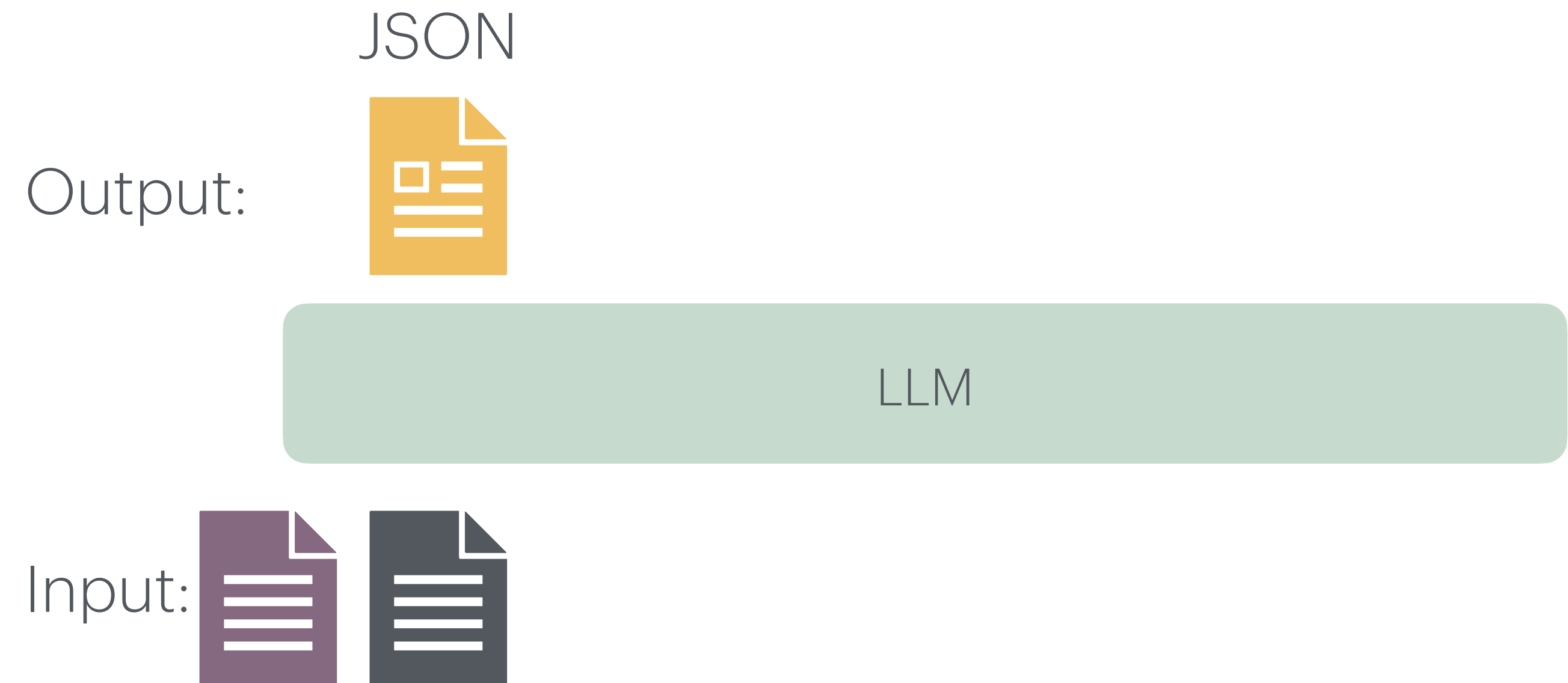- Limitations of LLMs

# Long Context

Philipp Krähenbühl, UT Austin

# Full Picture

Basic LLM

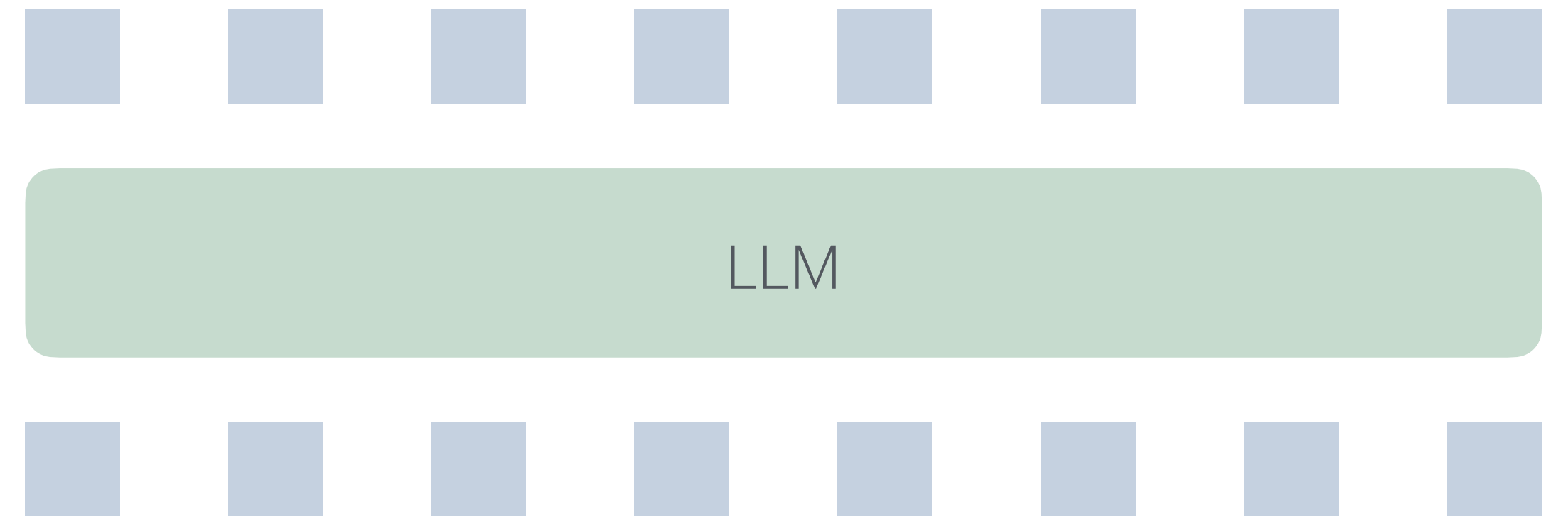| Pre-training | → | Instruction tuning | → | RLHF / DPO |
|---|---|---|---|---|
| Datasets | | Datasets | | Datasets |

# Tools and Structured outputs

- Tools

  - Special tags, Special chat-template

- Structured output

  - Option 1.1: Write a robust parser (in python)

    - Let LLM know that you failed to parse

  - Option 1.2: Constrain output

  - Option 2: Use a tool, arguments = json fields

Supervision:

Output:

LLM

Input:

JSON

Output:

LLM

Input:

# Long Context

- Current model are **pre-trained** on **2-8k** token sequences

LLM

# Long Context

What happens if we feed ten's of thousands of tokens into an LLM?

???

LLM

Read these documents and find references to efficient long-context LLMs

# Long Context

??? 

What happens if we feed ten's of thousands of tokens into an LLM?

1. OOM (Out Of Memory)

LLM

Read these documents and find references to efficient long-context LLMs

# Long Context

What happens if we feed ten's of thousands of tokens into an LLM?

1. OOM (Out Of Memory)

Transformer Layer

KV-cache

Transformer Layer

KV-cache

Transformer Layer

KV-cache

Transformer Layer

KV-cache

Transformer Layer

KV-cache

Embedding

Tokens

# Long Context

What happens if we feed ten's of thousands of tokens into an LLM?

1. OOM (Out Of Memory)

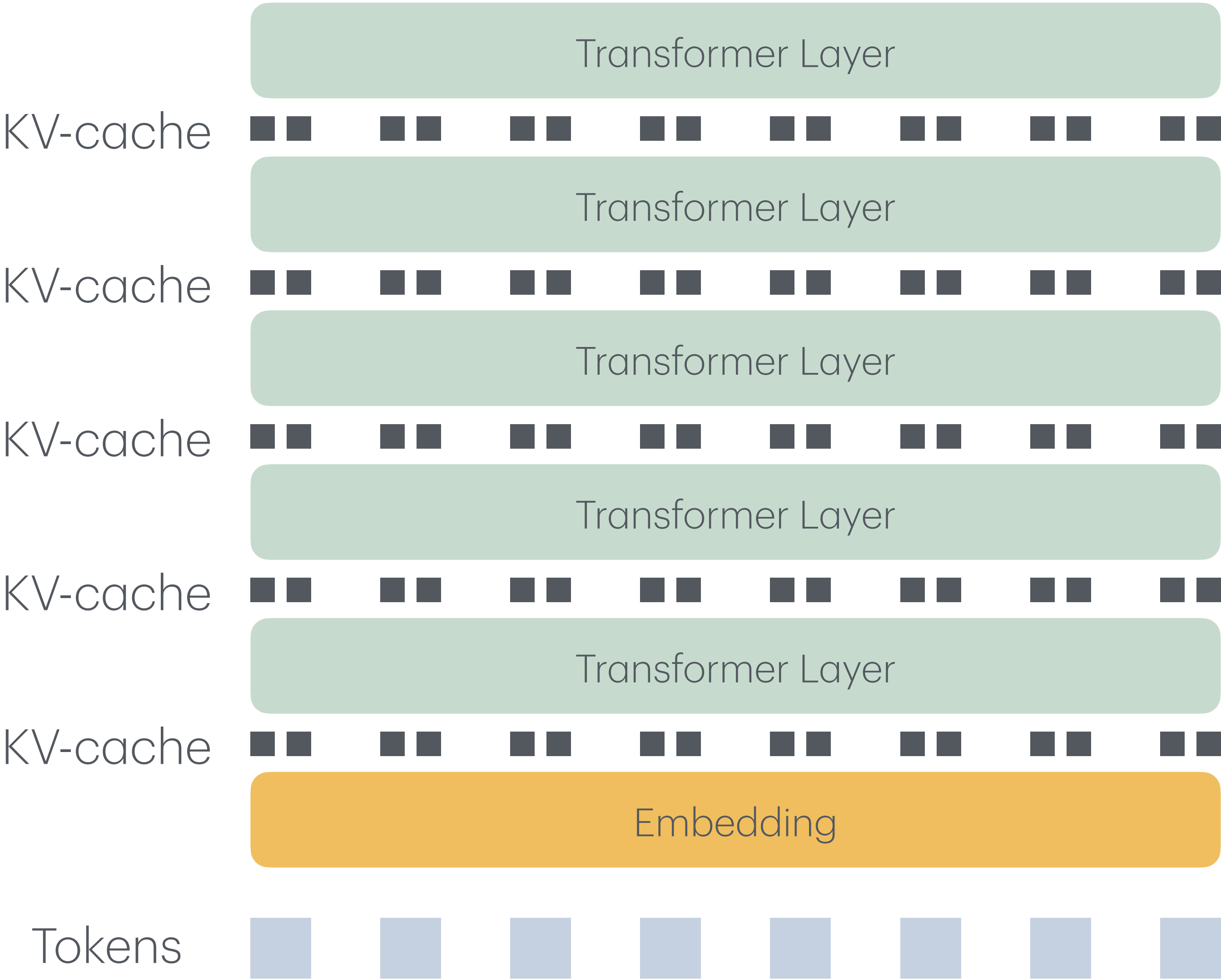2. Model will be very slow

???

LLM

Read these documents and find references to efficient long-context LLMs

# Long Context

What happens if we feed ten's of thousands of tokens into an LLM?

1. OOM (Out Of Memory)

2. Model will be very slow

$O(L)$

$O(L^2)$

Transformer Layer

MLP

Multi-Head Attention

Embedding

# Activation Beacon

# Activation Beacon



- Start from pre-trained model

- Partition sequence into chunks of 1024

- Pick k "beacons" per chunk

- Chunk n only sees beacons of chunks 1...n-1

- Fine-tune



Long Context Compression with Activation Beacon, Zhang etal 2024

# Long Context

??? 

What happens if we feed ten's of thousands of tokens into an LLM?

| LLM |
| --- |

1. OOM (Out Of Memory)

2. Model will be very slow
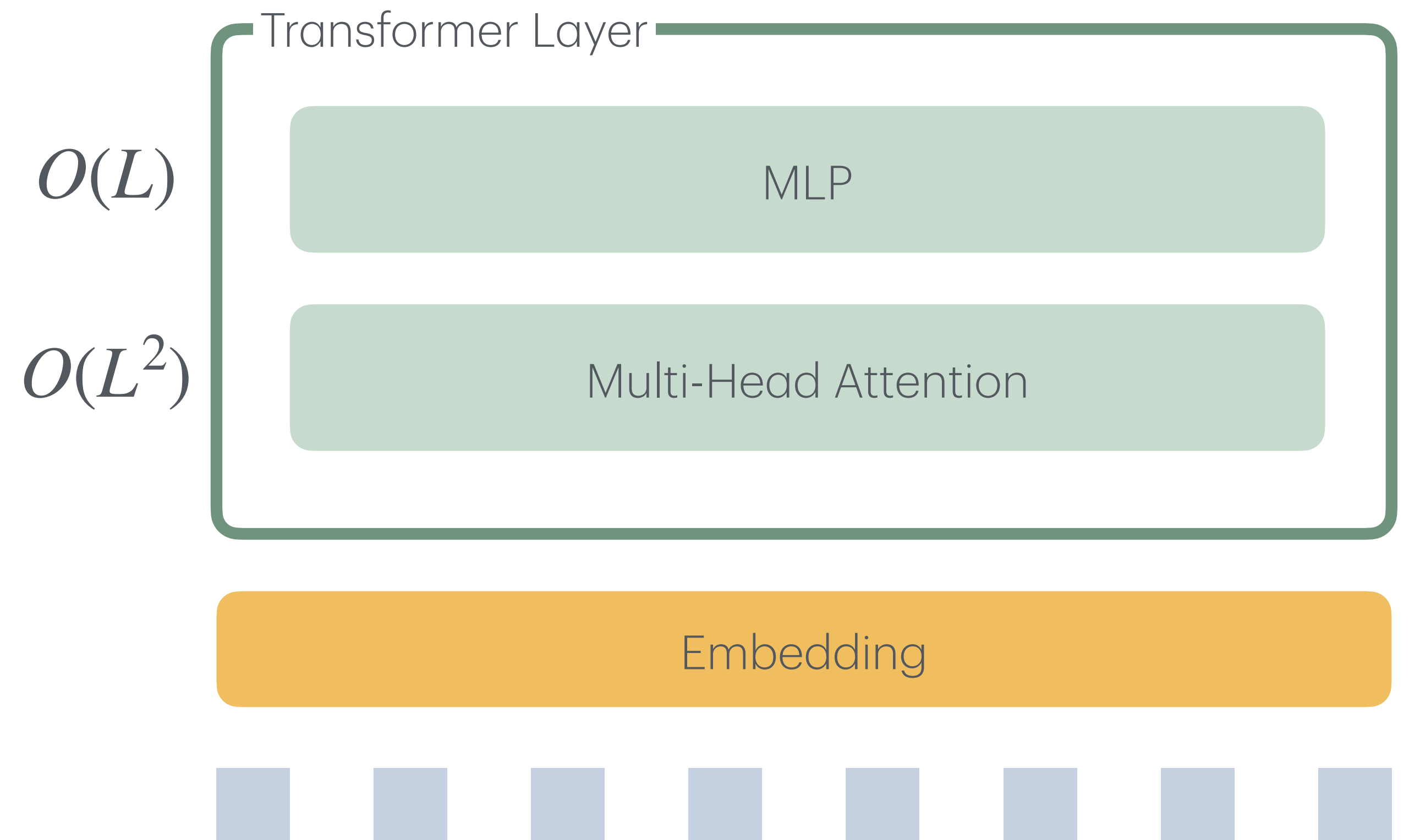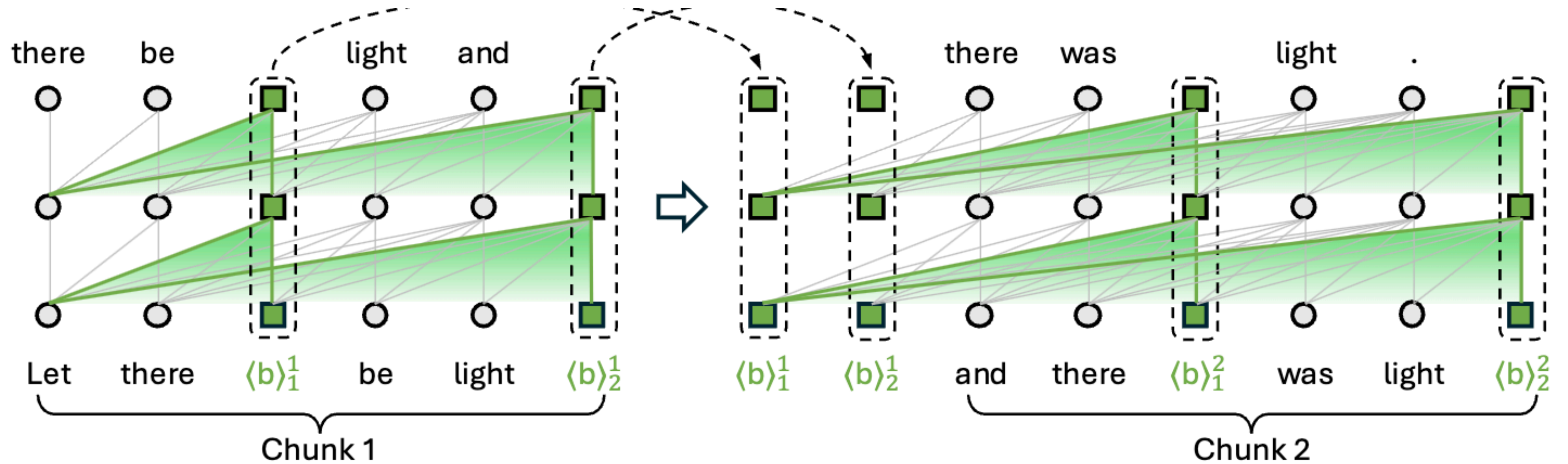
Activation Beacons and friends

Read these documents and find references to efficient long-context LLMs

# Long Context

???

What happens if we feed ten's of thousands of tokens into an LLM?

1. OOM (Out Of Memory)

2. Model will be very slow

3. Model will produce garbage outputs

Activation Beacons and friends

LLM

Read these documents and find references to efficient long-context LLMs

# Long Context

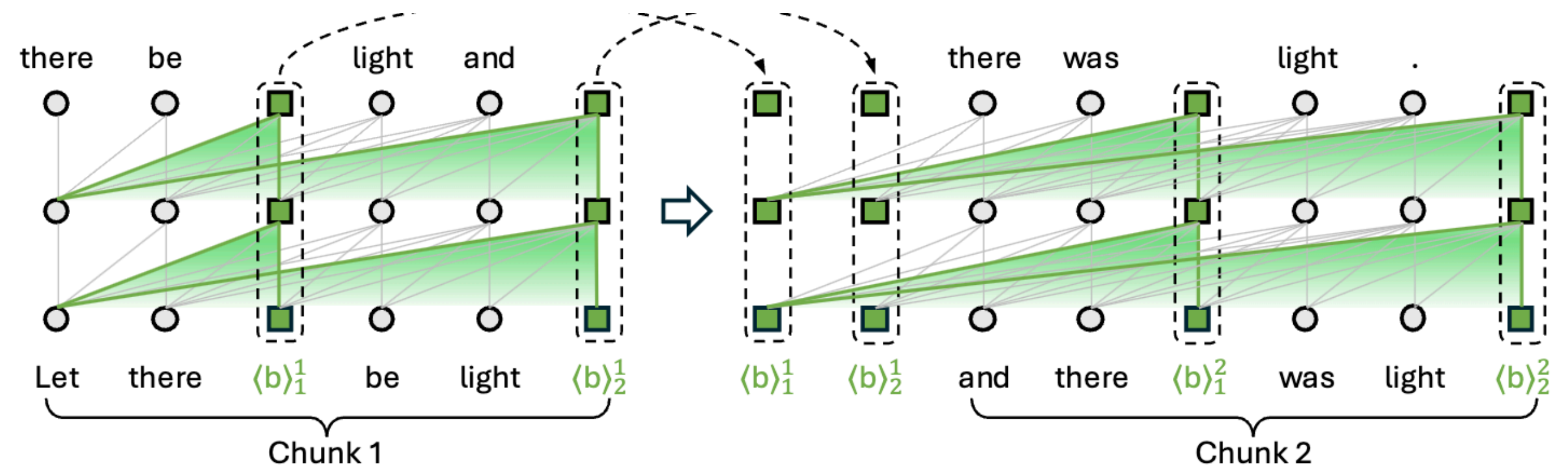What happens if we feed ten's of
thousands of tokens into an LLM?
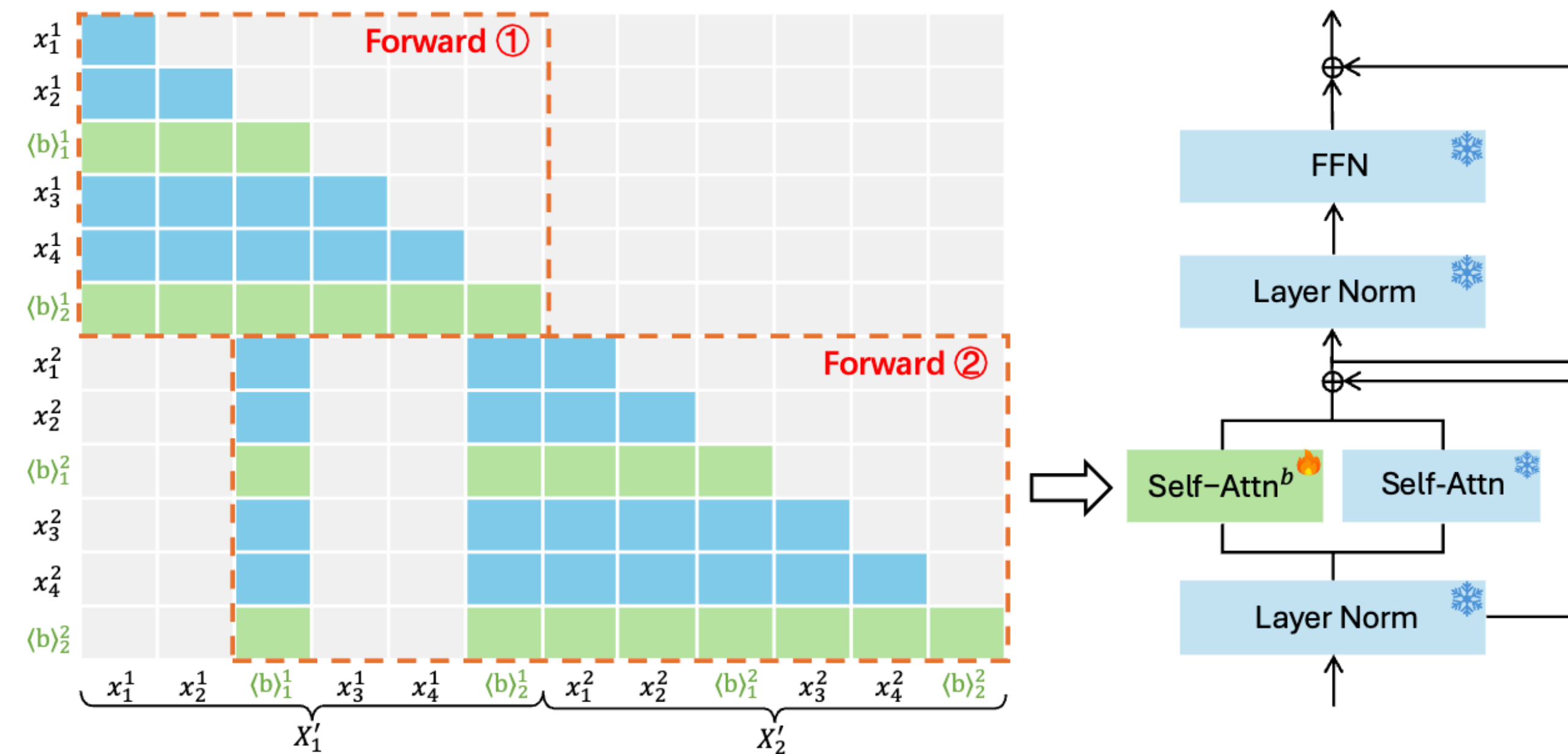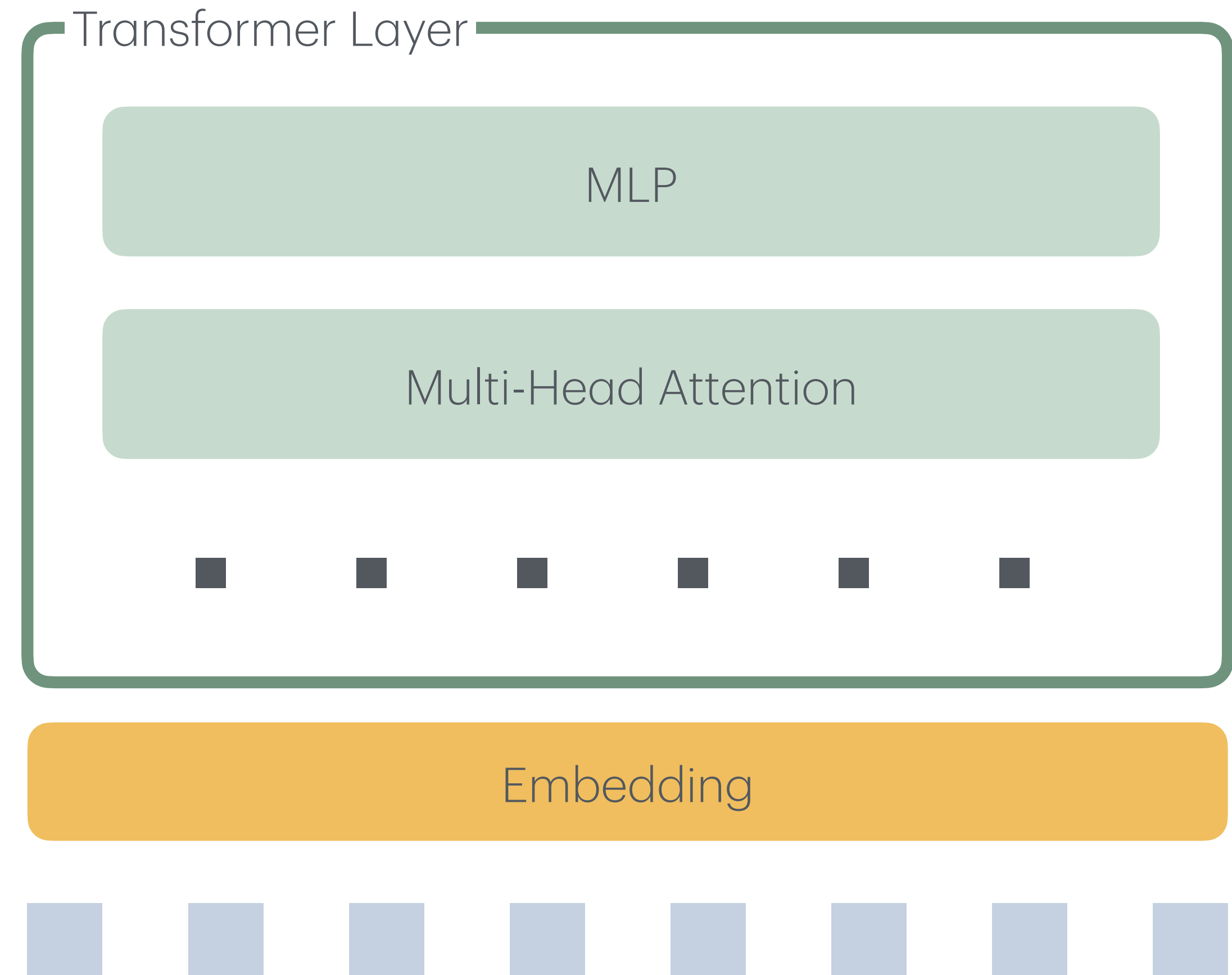
1. OOM (Out Of Memory)

2. Model will be very slow

3. Model will produce garbage outputs

Activation
Beacons
and
friends

Positional
embedding

# Positional Embedding

- Rotary Embeddings

$$f_{\{q,k\}}(\boldsymbol{x}_m, m) = \boldsymbol{R}^d_{\Theta,m} \boldsymbol{W}_{\{q,k\}} \boldsymbol{x}_m$$

$$\boldsymbol{R}^d_{\Theta,m} = \begin{pmatrix} \cos m\theta_1 & -\sin m\theta_1 & 0 & 0 & \cdots & 0 & 0 \\ \sin m\theta_1 & \cos m\theta_1 & 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \cos m\theta_2 & -\sin m\theta_2 & \cdots & 0 & 0 \\ 0 & 0 & \sin m\theta_2 & \cos m\theta_2 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & \cos m\theta_{d/2} & -\sin m\theta_{d/2} \\ 0 & 0 & 0 & 0 & \cdots & \sin m\theta_{d/2} & \cos m\theta_{d/2} \end{pmatrix}$$

$$\boldsymbol{q}_m^\intercal \boldsymbol{k}_n = (\boldsymbol{R}^d_{\Theta,m} \boldsymbol{W}_q \boldsymbol{x}_m)^\intercal (\boldsymbol{R}^d_{\Theta,n} \boldsymbol{W}_k \boldsymbol{x}_n) = \boldsymbol{x}^\intercal \boldsymbol{W}_q R^d_{\Theta,n-m} \boldsymbol{W}_k \boldsymbol{x}_n$$

Positional embedding

Transformer Layer

MLP

Multi-Head Attention

Embedding

RoFormer: Enhanced Transformer with Rotary Position Embedding, Su etal 2021

# Positional Embedding



- Rotary Embeddings

- Fixed context length during training

  - Longer context for inference

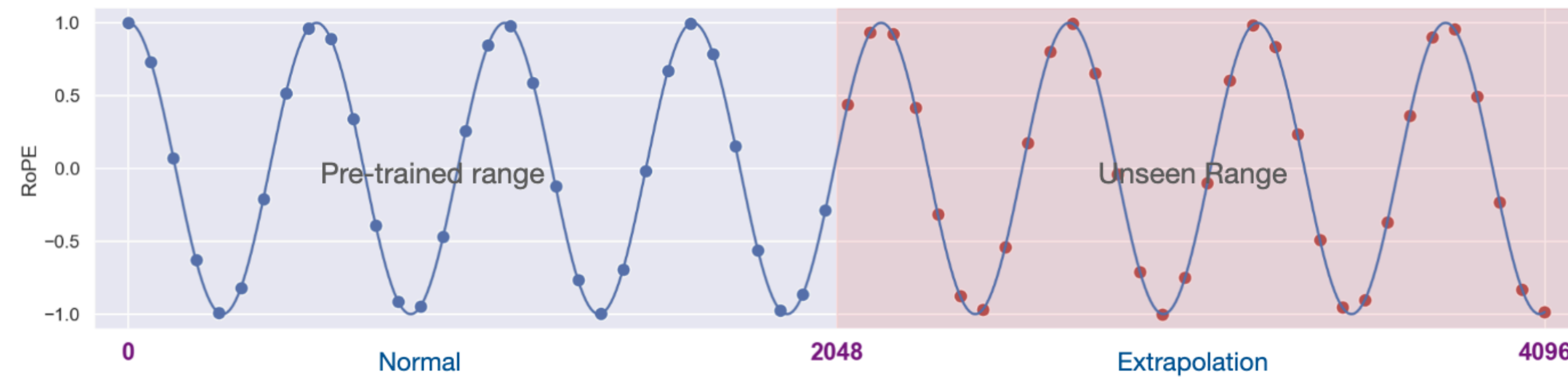$$f_{\{q,k\}}(\boldsymbol{x}_m, m) = \boldsymbol{R}_{\Theta,m}^d \boldsymbol{W}_{\{q,k\}} \boldsymbol{x}_m$$

$$\boldsymbol{R}_{\Theta,m}^d = \begin{pmatrix} \cos m\theta_1 & -\sin m\theta_1 & 0 & 0 & \cdots & 0 & 0 \\ \sin m\theta_1 & \cos m\theta_1 & 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \cos m\theta_2 & -\sin m\theta_2 & \cdots & 0 & 0 \\ 0 & 0 & \sin m\theta_2 & \cos m\theta_2 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & \cos m\theta_{d/2} & -\sin m\theta_{d/2} \\ 0 & 0 & 0 & 0 & \cdots & \sin m\theta_{d/2} & \cos m\theta_{d/2} \end{pmatrix}$$

$$\boldsymbol{q}_m^\intercal \boldsymbol{k}_n = (\boldsymbol{R}_{\Theta,m}^d \boldsymbol{W}_q \boldsymbol{x}_m)^\intercal (\boldsymbol{R}_{\Theta,n}^d \boldsymbol{W}_k \boldsymbol{x}_n) = \boldsymbol{x}^\intercal \boldsymbol{W}_q R_{\Theta,n-m}^d \boldsymbol{W}_k \boldsymbol{x}_n$$

RoFormer: Enhanced Transformer with Rotary Position Embedding, Su etal 2021
Extending Context Window of Large Language Models via Positional Interpolation, Chen etal 2023

# Positional Embedding



- Rotary Embeddings

  - Do not extrapolate well



Effect of Extrapolation



RoFormer: Enhanced Transformer with Rotary Position Embedding, Su etal 2021
Extending Context Window of Large Language Models via Positional Interpolation, Chen etal 2023

# Positional Embedding



- Rotary Embeddings

  - Do not extrapolate well

  - But they interpolate



RoFormer: Enhanced Transformer with Rotary Position Embedding, Su etal 2021
Extending Context Window of Large Language Models via Positional Interpolation, Chen etal 2023

# RoPE Scaling



RoFormer: Enhanced Transformer with Rotary Position Embedding, Su etal 2021
Extending Context Window of Large Language Models via Positional Interpolation, Chen etal 2023

# RoPE Scaling



- Extrapolation

  - Make token stream **longer**

  - Does not generalize

- RoPE Scaling

  - Make token stream **denser**

  - Model generalizes

- Widely used

RoFormer: Enhanced Transformer with Rotary Position Embedding, Su etal 2021
Extending Context Window of Large Language Models via Positional Interpolation, Chen etal 2023

# Long Context

??? 

What happens if we feed ten's of thousands of tokens into an LLM?

LLM

1. OOM (Out Of Memory)
2. Model will be very slow
3. Model will produce garbage outputs

Activation Beacons and friends

RoPE scaling

Read these documents and find references to efficient long-context LLMs

# Long Context

??? 

- Current model are **pre-trained** on **2-8k** token sequences

- Late stage pre-training **8k-128k**

  - RoPE Scaling

- Fine-tuned on variable length sequences

LLM

Read these documents and find references to efficient long-context LLMs

# References

- [1] Long Context Compression with Activation Beacon, Zhang et al. 2024 (link)

- [2] RoFormer: Enhanced Transformer with Rotary Position Embedding, Su etal 2021 (link)

- [3] Extending Context Window of Large Language Models via Positional Interpolation, Chen et al 2023 (link)

# Retrieval Augmented Generation

Philipp Krähenbühl, UT Austin

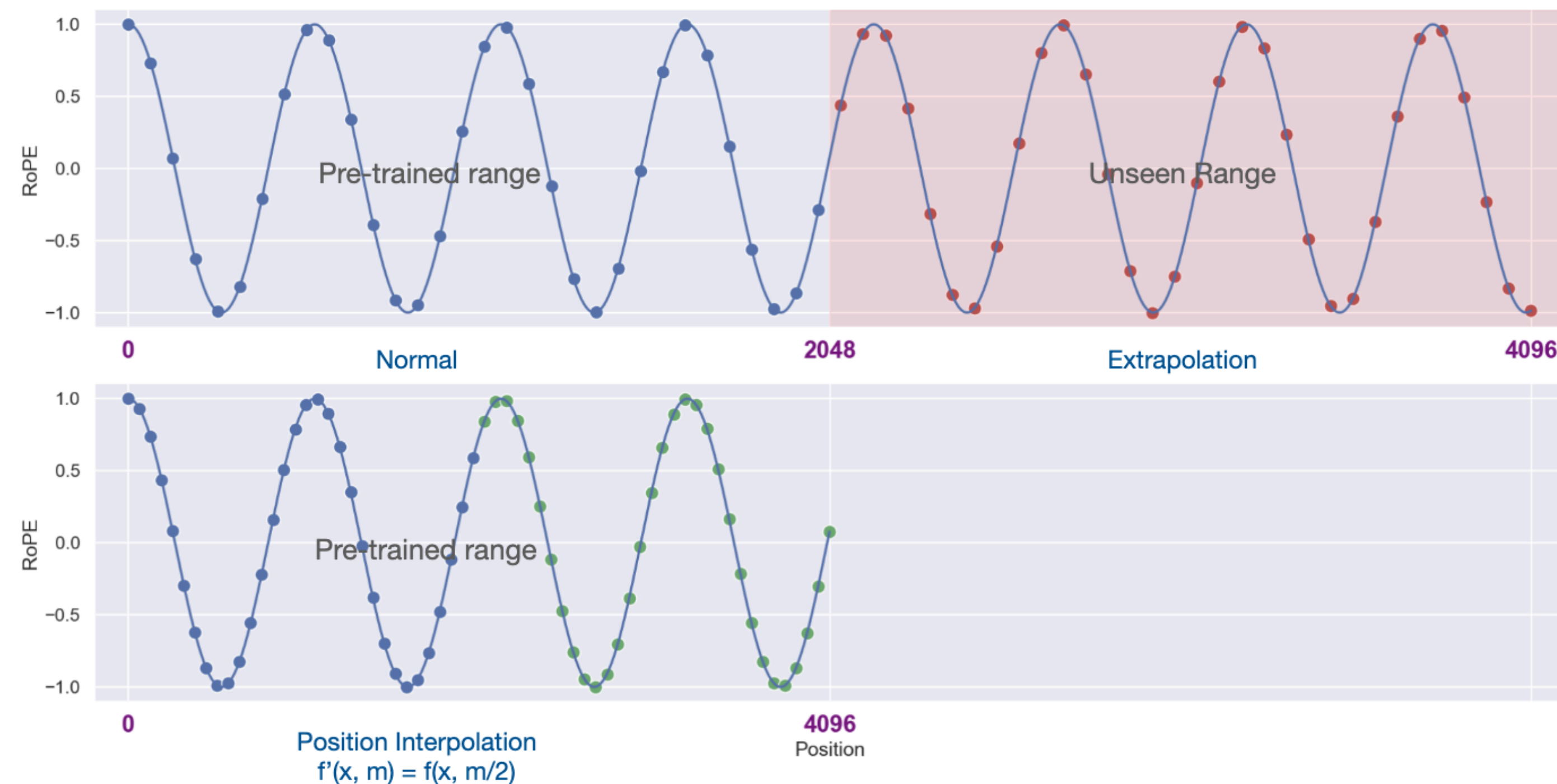# Full Picture

Basic LLM

Pre-training → Instruction tuning → RLHF / DPO

Datasets    Datasets    Datasets

# Long Context

- Current model are **pre-trained** on **2-8k** token sequences

- Late stage pre-training **8k-128k**

  - RoPE Scaling

- Fine-tuned on variable length sequences

???

LLM

Read these documents and find references to efficient long-context LLMs

# Longer Context

??? 

- What is we have even more inputs

LLM

Read these documents and find references to efficient long-context LLMs

# Longer Context

- What is we have even more inputs

  - We have to manage context

???

LLM

Read these documents and find references to efficient long-context LLMs

# Longer Context



Open-domain QA
SQuAD, TREC, WebQuestions, WikiMovies

Q: How many of Warsaw's inhabitants spoke Polish in 1933?

WIKIPEDIA
The Free Encyclopedia

Document Retriever → Document Reader → 833,500

- Solution: Build a "system"

- Option 1

  - Document Retriever: LLM to retrieve most relevant document

  - Document Reader: LLM to answer request

Reading Wikipedia to Answer Open-Domain Questions, Chen etal 2017

# Longer Context



- Solution: Build a "system"

- Option 2

  - Document Retriever: LLM to retrieve all relevant documents

  - LLM to answer request with documents in context

  - Fine-tuned for task

Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks, Lewis etal 2020
REALM: Retrieval-Augmented Language Model Pre-Training, Guu etal 2020
Improving language models by retrieving from trillions of tokens, Borgeaud 2021

# Longer Context



World Cup 2022 was the last with 32 teams, before the increase to

**Retriever**

FIFA World Cup 2026 will expand to 48 teams.

World Cup 2022 was the last with 32 teams, before the increase to

**Language Model**

48 in the 2026 tournament.

- Solution: Build a "system"

- Option 3

  - Document Retriever: LLM to retrieve all relevant documents

  - LLM to answer request with documents in context

  - ~~Fine-tuned for task~~    Model is prompted instead

In-Context Retrieval-Augmented Language Models, Ram etal 2023

# Retrieval Augmented Generation

## RAG

- A series of methods to manage the LLMs context

  - Some are trained

  - Some are just prompted

???

LLM

Read these documents and find references to efficient long-context LLMs

# References

- [1] Reading Wikipedia to Answer Open-Domain Questions, Chen etal 2017 (link)

- [2] Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks, Lewis etal 2020 (link)

- [3] REALM: Retrieval-Augmented Language Model Pre-Training, Guu etal 2020 (link)

- [4] Improving language models by retrieving from trillions of tokens, Borgeaud 2021 (link)

- [5] In-Context Retrieval-Augmented Language Models, Ram etal 2023 (link)

# Structured Dialogues

Philipp Krähenbühl, UT Austin

# Full Picture

Basic LLM

| Pre-training | → | Instruction tuning | → | RLHF / DPO |
|---|---|---|---|---|
| Datasets | | Datasets | | Datasets |

# Where does a LLM store information?

- Their weights

  - MLP and attention [1]

- Special tokens / activations [2,3]

  - Large activations or registers

- Their context

[1] Physics of Language Models: Part 3.3, Knowledge Capacity Scaling Laws, Allen-Zhu 2024
[2] Vision Transformers Need Registers, Darcet etal 2023
[3] Massive Activations in Large Language Models, Sun etal 2024

# Information in weights

- LLMs can store up to 2 bits of information per weight [1]

  - In MLP

  - In Attention

  - 2 bits require very long training and multiple (up to 1000) augmentations of same information



Physics of Language Models: Part 3.3, Knowledge Capacity Scaling Laws, Allen-Zhu 2024

# Special tokens / activations



- LLMs use special tokens to store information

  - LLMs attend to <BOS> token

  - VLMs attend to background

[1] Vision Transformers Need Registers, Darcet etal 2023
[2] Massive Activations in Large Language Models, Sun etal 2024

# Context

- LLMs store information in their context

- Examples

  - System prompt

  - Retrieval Augmented Generation

  - ...

[1] Vision Transformers Need Registers, Darcet etal 2023
[2] Massive Activations in Large Language Models, Sun etal 2024

# In context learning

```
Translate words from English to German using
JSON as an output. Here are some examples

Car
{"English": "Car", "German": "Auto"}

Sun
{"English": "Sun", "German": "Sonne"}

Moon
```

- Describe the task

  - Give examples input - output pairs

  - Then ask for your specific

Language Models are Few-Shot Learners, Brown etal 2020

# In context learning

## Why does it work?

- LLMs like repeating patterns

  - Likely exist in pre-training data

- Examples of in-context prompts and answers during training (instruction tuning, alignment)

```
Translate words from English to German using
JSON as an output. Here are some examples

Car
{"English": "Car", "German": "Auto"}

Sun
{"English": "Sun", "German": "Sonne"}

Moon
```

Language Models are Few-Shot Learners, Brown etal 2020

# In context learning

## What does it work for?

- Formatting outputs

- Simple requests

Translate words from English to German using JSON as an output. Here are some examples

Car
{"English": "Car", "German": "Auto"}

Sun
{"English": "Sun", "German": "Sonne"}

Moon

Language Models are Few-Shot Learners, Brown etal 2020

# Chain of thought

- Ask model to derive answer

  - Pre-instruction tuning: In-context example of reasoning

  - Post-instruction tuning

    - Ask model to think step-by-step before giving the answer

    - Guide model through thinking process

## Standard 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: 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 answer is 27. ❌

## 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. ✓

Chain-of-Thought Prompting Elicits Reasoning in Large Language Models, Wei etal 2022

# Chain of thought
## Why does it work?

- More output tokens = better performance

  - Delays making a decision

- Can work around tokenization issues

  - Break up numbers

**Standard 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: 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 answer is 27. ❌

**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. ✔️

Chain-of-Thought Prompting Elicits Reasoning in Large Language Models, Wei etal 2022

# Chain of thought

- Order matters

  - Think first, then answer

  - Chain-of-BS: Ask model to give answer and justify it

**Standard 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: 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 answer is 27. ❌

**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. ✔️

Chain-of-Thought Prompting Elicits Reasoning in Large Language Models, Wei etal 2022

# Self-Consistency



- Let the model reason multiple times

  - Pick the most frequent answer

  - Mathematically: Marginalize out reasoning to obtain most likely answer

# Tree of Thoughts

- Combine tree search with CoT

  - Requires a state-evaluator (i.e. reward/cost/scoring function or second LLM)



(a) Input-Output Prompting (IO)  (c) Chain of Thought Prompting (CoT)  (c) Self Consistency with CoT (CoT-SC)  **(d) Tree of Thoughts (ToT)**

Tree of Thoughts: Deliberate Problem Solving with Large Language Models, Yao etal 2023

# ReACT

- Chain of thought for iterative actions / tool use

  - Thought

  - Action

  - Observation (from external tool)



**(1) Hotspot QA**

**Question:** Aside from the Apple Remote, what other device can control the program Apple Remote was originally designed to interact with?

**(1a) Standard**

**Answer:** iPod ✗

**(1b) CoT (Reason Only)**

**Thought:** Let's think step by step. Apple Remote was originally designed to interact with Apple TV. Apple TV can be controlled by iPhone, iPad, and iPod Touch. So the answer is iPhone, iPad, and iPod Touch.

**Answer:** iPhone, iPad, iPod Touch ✗

**(1c) Act-Only**

**Act 1:** Search[Apple Remote]
**Obs 1:** The Apple Remote is a remote control …

**Act 2:** Search[Front Row]
**Obs 2:** Could not find [Front Row]. Similar: …

**Act 3:** Search[Front Row (software)]
**Obs 3:** Front Row is a discontinued media center software …

**Act 4:** Finish[yes] ✗

**(1d) ReAct (Reason + Act)**

**Thought 1:** I need to search Apple Remote and find the program it was originally designed to interact with.
**Act 1:** Search[Apple Remote]
**Obs 1:** The Apple Remote is a remote control introduced in October 2005 by Apple … originally designed to control the Front Row media center program …

**Thought 2:** Apple Remote was originally designed to control the Front Row media center program. I need to search Front Row next and find what other device can control it.
**Act 2:** Search[Front Row]
**Obs 2:** Could not find [Front Row]. Similar: ['Front Row Seat to Earth', 'Front Row Motorsports',' Front Row (software)', …]

**Thought 3:** Front Row is not found. I need to search Front Row (software) .
**Act 3:** Search[Front Row (software)]
**Obs 3:** Front Row is a discontinued media center software …

**Thought 4:** Front Row (software) is controlled by an Apple Remote or the keyboard function keys. So the answer is keyboard function keys.
**Act 4:** Finish[keyboard function keys] ✓

**(2) AlfWorld**

You are in the middle of a room. Looking quickly around you, you see a cabinet 6, a cabinet 1, a coffee machine 1, a countertop 3, a stove burner 1, and a toaster 1.
**Your task is to:** Put some pepper shaker on a drawer.

**(2a) Act-Only**

**Act 1:** Go to drawer 1
**Obs 1:** The drawer 1 is closed.

**Act 2:** Open drawer 1
**Obs 2:** You open the drawer 1. The drawer 1 is open. In it, you see a dishsponge 2, and a spoon 1.

**Act 3:** Go to sinkbasin 1
**Obs 3:** On the sinkbasin 1, you see a dishsponge 3, a spatula 1, and a spoon 2.

**Act 4:** Take peppershaker 1 from sinkbasin 1
**Obs 4:** Nothing happens.

**Act 5:** Take peppershaker 1 from sinkbasin 1
**Obs 5:** Nothing happens. ✗

**(2b) ReAct (Reason + Act)**

**Act 1:** Think[First I need to find a pepper shaker...more likely to appear in cabinets (1-6), countertops (1-3), ... ]
**Act 2:** Go to cabinet 1
**Obs 2:** On the cabinet 1, you see a vase 2.

(Here the agent go to cabinet 1, then cabinet 2, 3, then countertop 1 and 2)

**Act 6:** Go to countertop 3
**Obs 6:** On the countertop 3, you see a apple 1, a bread 1, ...... a pepper shaker 1, and a vase 3.

**Act 7:** Take pepper shaker 1 from countertop 3
**Obs 7:** You pick up the pepper shaker 1 from the countertop 3.

**Act 8:** Think[Now I find a pepper shaker 1. Next, I need to put it in/on drawer 1. ]
**Act 9:** Go to drawer 1
**Obs 9:** Drawer 1 is closed.

**Act 10:** Open drawer 1
**Obs 10:** You open Drawer 1 …

**Act 11:** Put pepper shaker 1 in/on drawer 1
**Obs 11:** You put pepper shaker 1 in/on the drawer 1. ✓

ReAct: Synergizing Reasoning and Acting in Language Models, Yao etal 2022

# Structured Dialogues

User: Request

LLM: Thought

Answer

- Break down problem / tasks for LLM

  - Higher performance

  - Lots of human engineering / prompting

# Reflexion

- Chain of Thought / ReACT

- Obtain observation / result

- Reflect on outcome

- Repeat



| 1. Decision making | 2. Programming | 3. Reasoning |
|---|---|---|
| **(a) Task** You are in the middle of a room [...] **Task:** clean some pan and put it in countertop. | **Task:** You are given a list of two strings [...] of open '(' or close ')' parentheses only [...] | **Task:** What profession does John Lanchester and Alan Dean Foster have in common? |

**(b) Trajectory**

Decision making:
```
[...]
Action:take pan1 from stoveburner1
Obs:Nothing happens. [...]
Action:clean pan1 with sinkbasin1
Obs:Nothing happens. [...]
```

Programming:
```
def match_parens(lst):
    if s1.count('(') +
s2.count('(') == s1.count(')') +
s2.count(')'): [...]
    return 'No'
```

Reasoning:
```
Think: [...] novelist, journalist,
critic [...] novelist,
screenwriter [...] common is
novelist and screenwriter.
Action: "novelist, screenwriter"
```

**(c) Evaluation** (internal / external)

Decision making: **Rule/LM Heuristic:** Hallucination.

Programming: **Self-generated unit tests fail:** assert match_parens(...)

Reasoning: **Environment Binary Reward:** 0

**(d) Reflection**

Decision making: [...] tried to pick up the pan in stoveburner 1 [...] but the pan was not in stoveburner 1. [...]

Programming: [...] wrong because it only checks if the total count of open and close parentheses is equal [...] order of the parentheses [...]

Reasoning: [...] failed because I incorrectly assumed that they both had the same multiple professions [...] accurately identifying their professions.

**(e) Next Trajectory**

Decision making:
```
[...] Action: take pan 1 from
stoveburner 2
[...] Obs: You put the pan 1 in
countertop 1.
```

Programming:
```
[...]
    return 'Yes' if check(S1) or
check(S2) else 'No'
```

Reasoning:
```
Think: [...] So the profession
John Lanchester and Alan Dean
Foster have in common is novelist.
Action: "novelist"
```

Reflexion: Language Agents with Verbal Reinforcement Learning, Shin etal 2023

# Reflexion



**Agent**

External feedback

Self-reflection (LM)

Internal feedback

Reflective text

Evaluator (LM)

Experience (long-term memory)

Trajectory (short-term memory)

Actor (LM)

Obs / Reward — Environment — Action

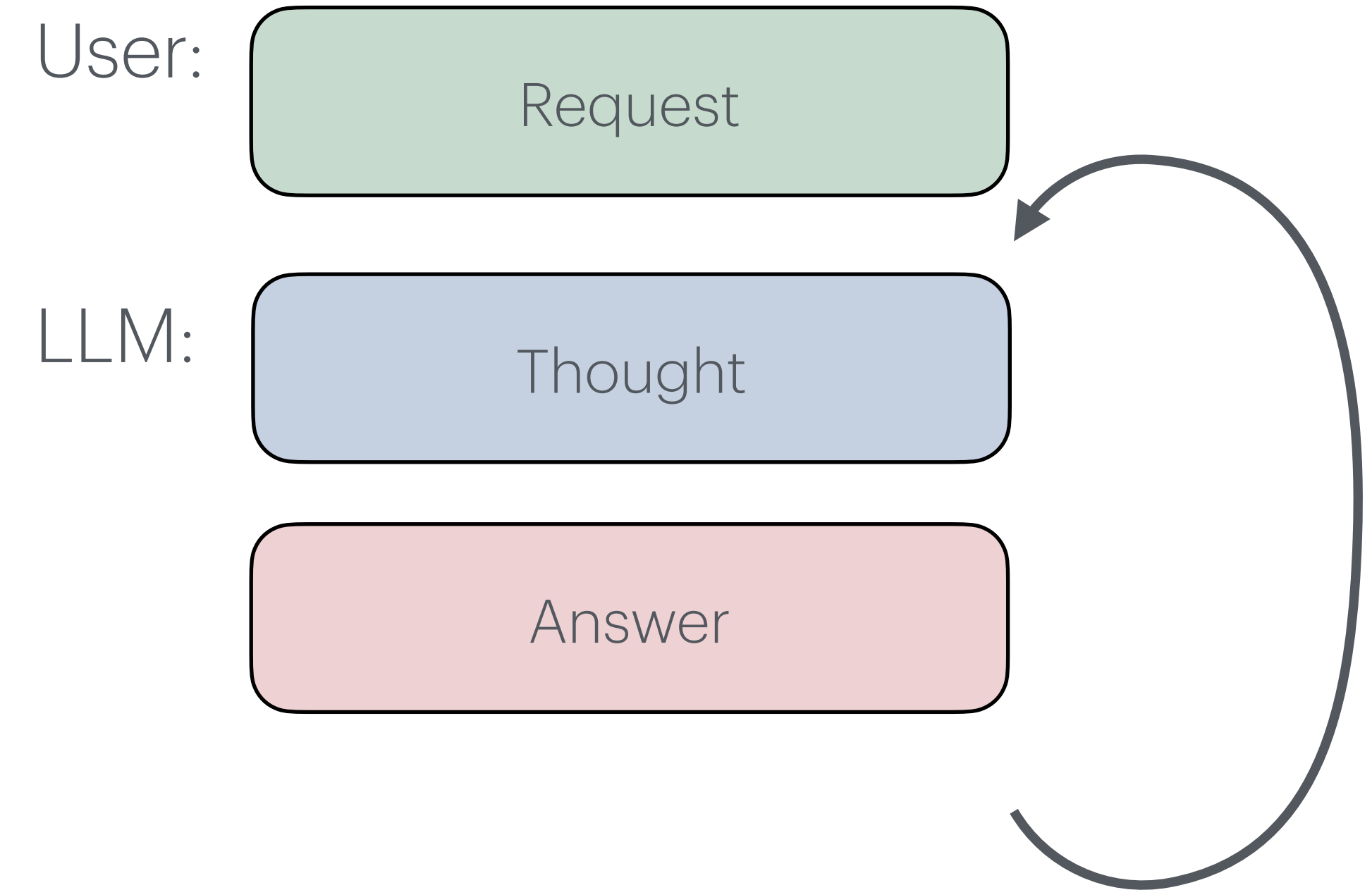**Algorithm 1** Reinforcement via self-reflection

Initialize Actor, Evaluator, Self-Reflection: $M_a$, $M_e$, $M_{sr}$
Initialize policy $\pi_\theta(a_i|s_i)$, $\theta = \{M_a, mem\}$
Generate initial trajectory using $\pi_\theta$
Evaluate $\tau_0$ using $M_e$
Generate initial self-reflection $sr_0$ using $M_{sr}$
Set $mem \leftarrow [sr_0]$
Set $t = 0$
**while** $M_e$ not pass or $t <$ max trials **do**
　Generate $\tau_t = [a_0, o_0, \ldots a_i, o_i]$ using $\pi_\theta$
　Evaluate $\tau_t$ using $M_e$
　Generate self-reflection $sr_t$ using $M_{sr}$
　Append $sr_t$ to $mem$
　Increment $t$
**end while**
**return**

- Connections to reinforcement learning

  - More strictly planning

- Requires a evaluator (cost function)

  - External environment (i.e. simulator, code interpreter)

  - LLM generated tests

  - Trained LLM verifier [1]



Test Generation

Internal Tests

Self-reflection

Task

Proposed Solution

Refined Solution

[1] Generative Verifiers: Reward Modeling as Next-Token Prediction, Zhang etal 2024

# Reflexion

- Break down problem / tasks for LLM

  - Higher performance

  - Lots of human engineering / prompting
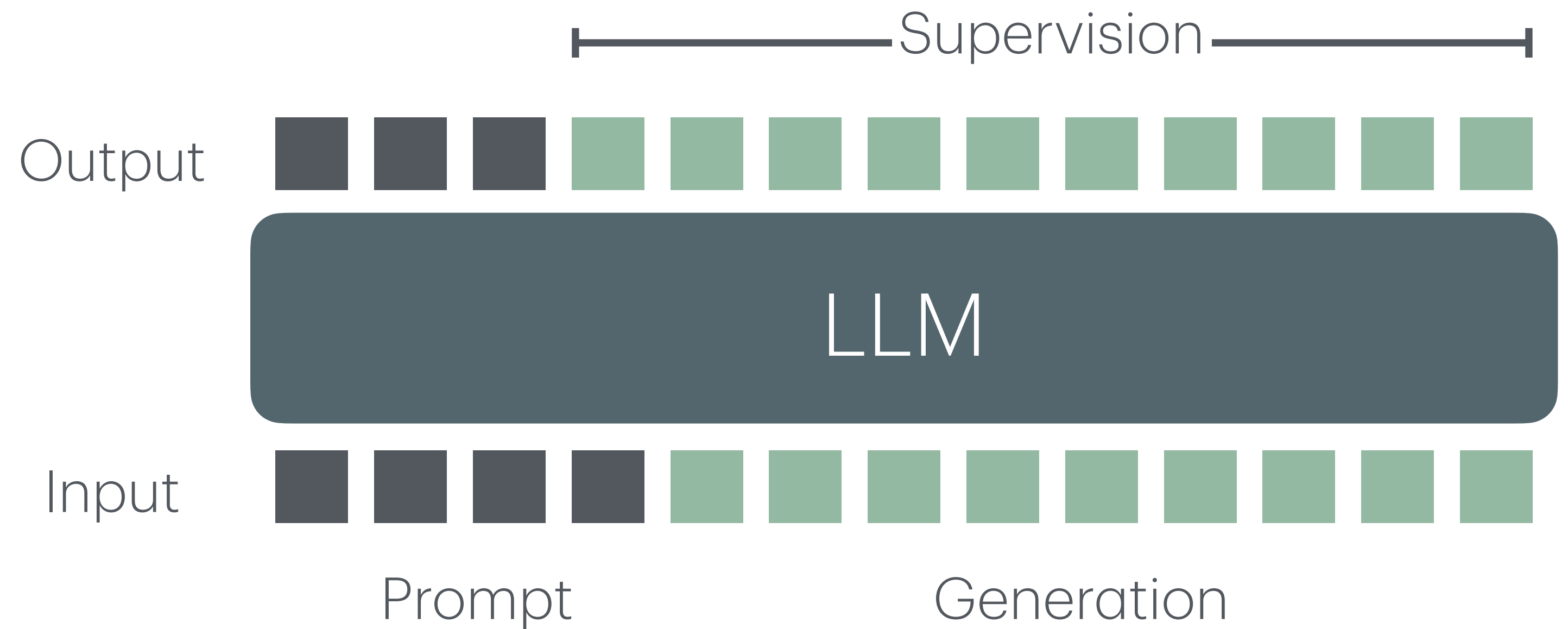
User: [ Request ]

LLM: [ Thought ]

[ Answer ]

# References

- [1] Physics of Language Models: Part 3.3, Knowledge Capacity Scaling Laws, Allen-Zhu 2024

- [2] Vision Transformers Need Registers, Darcet etal 2023

- [3] Massive Activations in Large Language Models, Sun etal 2024

- [4] Language Models are Few-Shot Learners, Brown etal 2020

- [5] Chain-of-Thought Prompting Elicits Reasoning in Large Language Models, Wei etal 2022

- [6] Self-Consistency Improves Chain of Thought Reasoning in Language Models, Wang etal 2022

- [7] Tree of Thoughts: Deliberate Problem Solving with Large Language Models, Yao etal 2023

- [8] ReAct: Synergizing Reasoning and Acting in Language Models, Yao etal 2022

- [9] Reflexion: Language Agents with Verbal Reinforcement Learning, Shin etal 2023

- [10] Generative Verifiers: Reward Modeling as Next-Token Prediction, Zhang etal 2024

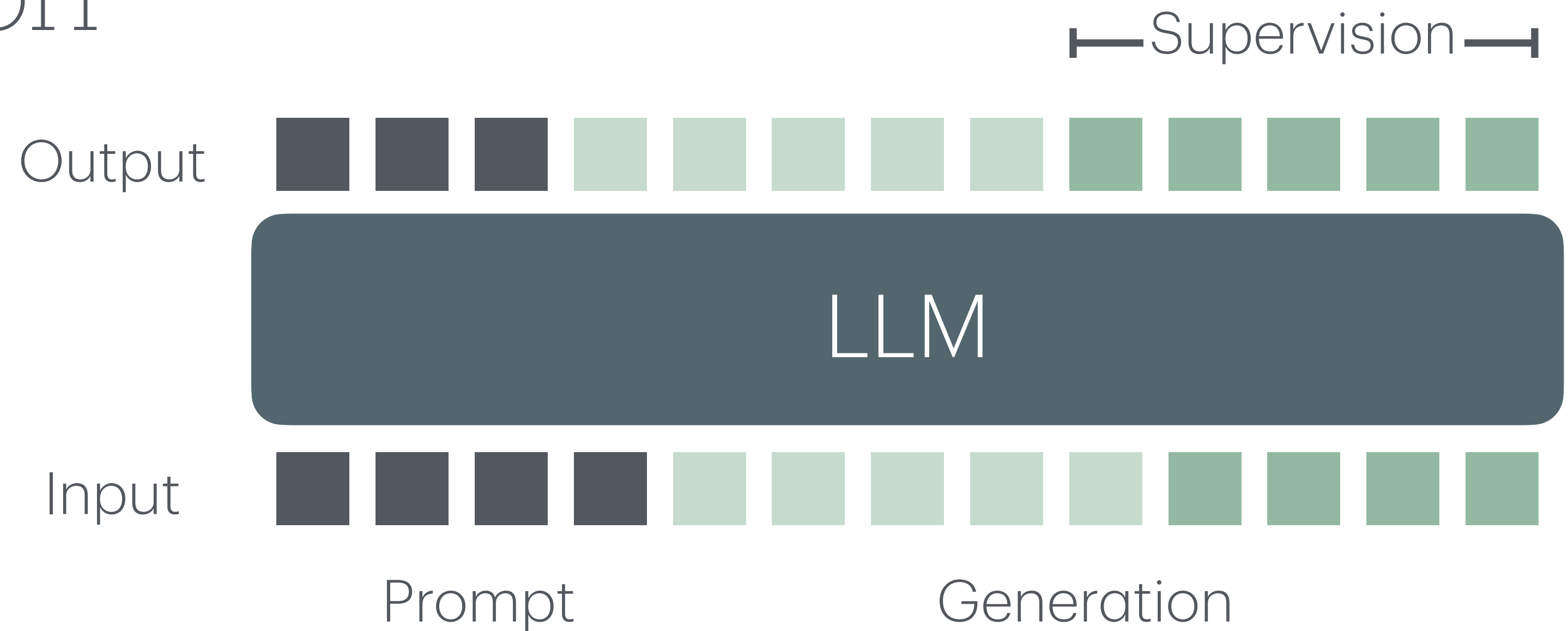# Reinforcement Learning and LLMs

Philipp Krähenbühl, UT Austin

# Teacher forcing

- Simple supervised learning

  - Input = Prompt + Target[0:-1]

  - Loss(output, Target[1:])

Supervision

Output

LLM

Input

Prompt

Generation

# Outcome supervision

Output

## LLM

Input

Prompt          Generation

- What if we only supervise the final result?

- Generation

  - Loss(Generation)

- Teacher-forcing not possible

  - No supervised loss

- Solution: RL

# Outcome supervision

## Reinforcement Learning

- LLM $\quad p_\theta(x_{t+1} \mid \mathbf{c}, x_1 \ldots x_t)$

$$p_\theta(\mathbf{x} \mid \mathbf{c}) = \prod_{t=1}^{N} p_\theta(x_{t+1} \mid \mathbf{c}, x_1 \ldots x_t)$$

- Sampling / Generation

$$x_{t+1} \sim p_\theta(\,\cdot\, \mid \mathbf{c}, x_1 \ldots x_t)$$

- MDP

$$E_{\mathbf{x} \sim p_\theta(\cdot \mid \mathbf{c})} \Bigg[ \underbrace{\sum_{t=1}^{N} r(x_t \mid \mathbf{c}, x_1 \ldots x_{t-1})}_{R(\mathbf{c}, \mathbf{x})} \Bigg]$$



Supervision

Output

LLM

Input

Prompt $\mathbf{c}$  Generation $\mathbf{x}$

# REINFORCE

maximize $E_{\mathbf{x} \sim p_\theta(\cdot|\mathbf{c})} \left[ R(\mathbf{c}, \mathbf{x}) \right]$

- Using gradient ascent

$$\nabla_\theta E_{\mathbf{x} \sim p_\theta(\cdot|\mathbf{c})} \left[ R(\mathbf{c}, \mathbf{x}) \right] = E_{\mathbf{x} \sim p_\theta(\cdot|\mathbf{c})} \left[ R(\mathbf{c}, \mathbf{x}) \nabla_\theta \log p_\theta(\mathbf{x} \,|\, \mathbf{c}) \right]$$

- With a Monte-Carlo estimate

$$\nabla_\theta E_{\mathbf{x} \sim p_\theta(\cdot|\mathbf{c})} \left[ R(\mathbf{c}, \mathbf{x}) \right] \approx \frac{1}{K} \sum_{k=1}^{K} R(\mathbf{c}, \mathbf{x}_k) \nabla_\theta \log p_\theta(\mathbf{x}_k \,|\, \mathbf{c})$$

for $\mathbf{x}_k \sim p_\theta( \cdot \,|\, \mathbf{c})$

- REINFORCE K=1 works!!!

Initialize $\theta$

for ever:

Sample (or iterate over) $\mathbf{c}$

$\mathbf{x} \sim p_\theta( \cdot \,|\, \mathbf{c})$

$\theta \leftarrow \theta + \epsilon R(\mathbf{c}, \mathbf{x}) \nabla \log p_\theta(\mathbf{x} | \mathbf{c})$

# Policy Gradient

$$E_{\mathbf{x} \sim p_\theta(\cdot|\mathbf{c})} \left[ A(\mathbf{c}, \mathbf{x}) \, \nabla_\theta \log p_\theta(\mathbf{x} \,|\, \mathbf{c}) \right]$$

- Even better

$$E_{\mathbf{x} \sim p_\theta(\cdot|\mathbf{c})} \left[ \sum_{t=1}^{T} A(\mathbf{c}, x_1 \dots x_t) \, \nabla_\theta \log p_\theta(x_t \,|\, \mathbf{c}, x_1 \dots x_{t-1}) \right]$$

---

**Algorithm 1** Vanilla Policy Gradient Algorithm

1: Input: initial policy parameters $\theta_0$, initial value function parameters $\phi_0$
2: **for** $k = 0, 1, 2, \dots$ **do**
3:     Collect set of trajectories $\mathcal{D}_k = \{\tau_i\}$ by running policy $\pi_k = \pi(\theta_k)$ in the environment.
4:     Compute rewards-to-go $\hat{R}_t$.
5:     Compute advantage estimates, $\hat{A}_t$ (using any method of advantage estimation) based on the current value function $V_{\phi_k}$.
6:     Estimate policy gradient as

$$\hat{g}_k = \frac{1}{|\mathcal{D}_k|} \sum_{\tau \in \mathcal{D}_k} \sum_{t=0}^{T} \nabla_\theta \log \pi_\theta(a_t|s_t)\big|_{\theta_k} \hat{A}_t.$$

7:     Compute policy update, either using standard gradient ascent,

$$\theta_{k+1} = \theta_k + \alpha_k \hat{g}_k,$$

    or via another gradient ascent algorithm like Adam.
8:     Fit value function by regression on mean-squared error:

$$\phi_{k+1} = \arg\min_\phi \frac{1}{|\mathcal{D}_k|T} \sum_{\tau \in \mathcal{D}_k} \sum_{t=0}^{T} \left( V_\phi(s_t) - \hat{R}_t \right)^2,$$

    typically via some gradient descent algorithm.
9: **end for**

# Proximal Policy Optimization

# PPO

- Policy gradient

- Reuse rollouts (go slightly off-policy)

  - Basic Option: Importance weighting

    maximize $E_{\mathbf{x} \sim p_\psi(\cdot|\mathbf{c})} \left[ \dfrac{p_\theta(\mathbf{x}|\mathbf{c})}{p_\psi(\mathbf{x}|\mathbf{c})} A(\mathbf{c}, \mathbf{x}) \right]$

  - Better Option: PPO-Clipping

    maximize $E_{\mathbf{x} \sim p_\psi(\cdot|\mathbf{c})} \left[ \dfrac{1}{T} \sum_{t=1}^{T} CLIP\left( \dfrac{p_\theta(x_t|\mathbf{c}, x_1 \ldots x_{t-1})}{p_\psi(x_t|\mathbf{c}, x_1 \ldots x_{t-1})}, A(\mathbf{c}, \mathbf{x}) \right) \right]$

Credit: OpenAI

# REINFORCE vs PPO

Initialize $\theta$

for ever:

    Sample (or iterate over) $\mathbf{c}$

$$\mathbf{x} \sim p_\theta(\,\cdot\,|\,\mathbf{c})$$

$$\theta \leftarrow \theta + \epsilon R(\mathbf{c}, \mathbf{x})\,\nabla \log p_\theta(\mathbf{x}\,|\,\mathbf{c})$$

**Algorithm 1** PPO-Clip

1: Input: initial policy parameters $\theta_0$, initial value function parameters $\phi_0$
2: **for** $k = 0, 1, 2, \ldots$ **do**
3:     Collect set of trajectories $\mathcal{D}_k = \{\tau_i\}$ by running policy $\pi_k = \pi(\theta_k)$ in the environment.
4:     Compute rewards-to-go $\hat{R}_t$.
5:     Compute advantage estimates, $\hat{A}_t$ (using any method of advantage estimation) based on the current value function $V_{\phi_k}$.
6:     Update the policy by maximizing the PPO-Clip objective:

$$\theta_{k+1} = \arg\max_\theta \frac{1}{|\mathcal{D}_k|T} \sum_{\tau \in \mathcal{D}_k} \sum_{t=0}^{T} \min\left( \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_k}(a_t|s_t)} A^{\pi_{\theta_k}}(s_t, a_t),\;\; g(\epsilon, A^{\pi_{\theta_k}}(s_t, a_t)) \right),$$

    typically via stochastic gradient ascent with Adam.
7:     Fit value function by regression on mean-squared error:

$$\phi_{k+1} = \arg\min_\phi \frac{1}{|\mathcal{D}_k|T} \sum_{\tau \in \mathcal{D}_k} \sum_{t=0}^{T} \left( V_\phi(s_t) - \hat{R}_t \right)^2,$$

    typically via some gradient descent algorithm.
8: **end for**

# Back to Basics: Revisiting REINFORCE Style Optimization for Learning from Human Feedback in LLMs

**Arash Ahmadian**
*Cohere For AI*

**Chris Cremer**
*Cohere*

**Matthias Gallé**
*Cohere*

**Marzieh Fadaee**
*Cohere For AI*

**Julia Kreutzer**
*Cohere For AI*

**Olivier Pietquin**
*Cohere*

**Ahmet Üstün**
*Cohere For AI*

**Sara Hooker**
*Cohere For AI*

`{arash,olivier,ahmet,sarahooker}@cohere.com`

# Outcome supervision

## Reinforcement Learning



- LLM $\quad p_\theta(x_{t+1} \mid \mathbf{c}, x_1 \ldots x_t)$

$$p_\theta(\mathbf{x} \mid \mathbf{c}) = \prod_{t=1}^{N} p_\theta(x_{t+1} \mid \mathbf{c}, x_1 \ldots x_t)$$

- Sampling / Generation

$$x_{t+1} \sim p_\theta( \cdot \mid \mathbf{c}, x_1 \ldots x_t)$$

- MDP

$$E_{\mathbf{x} \sim p_\theta(\cdot \mid \mathbf{c})} \left[ \sum_{t=1}^{N} r(x_t \mid \mathbf{c}, x_1 \ldots x_{t-1}) \right]$$

# Outcome supervision

## Reinforcement Learning

- LLM $p_\theta(\mathbf{x}|\mathbf{c})$

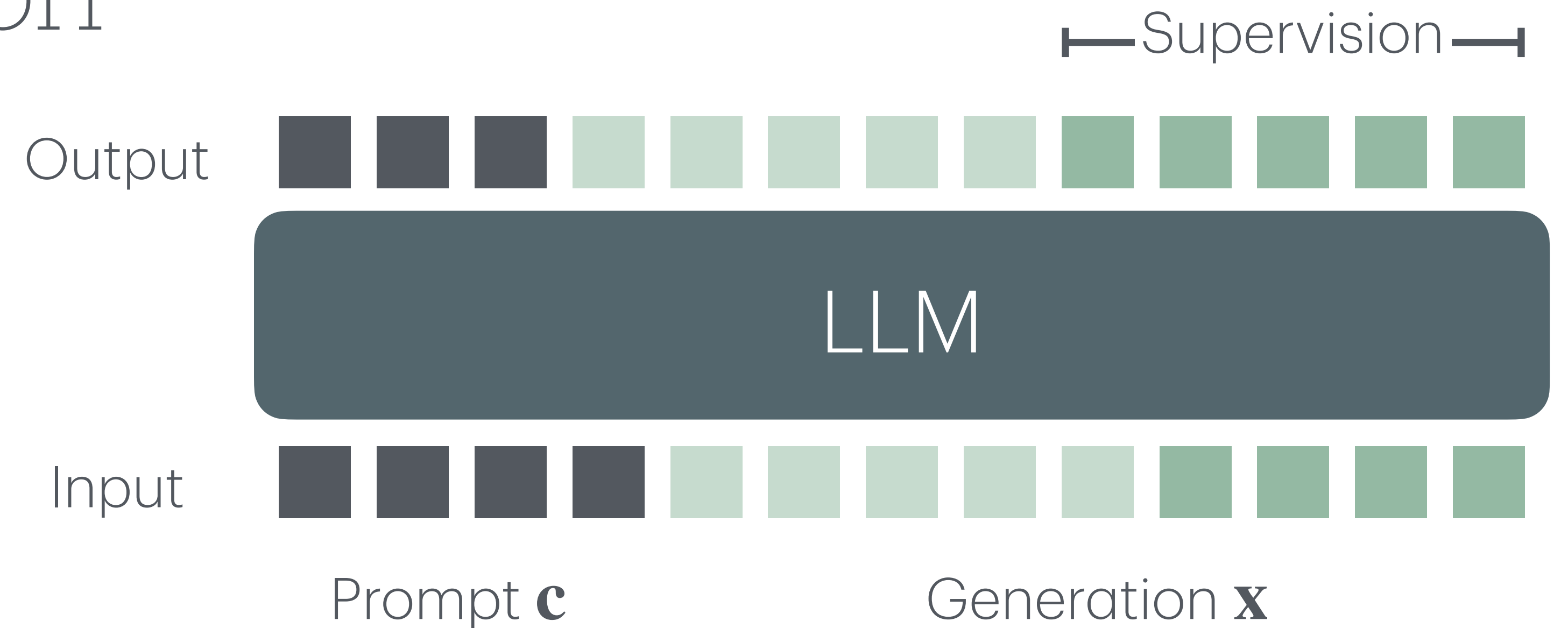- Sampling / Generation

$\mathbf{x} \sim p_\theta(\,\cdot\,|\mathbf{c})$

- Contextual bandit
$E_{\mathbf{x} \sim p_\theta(\cdot|\mathbf{c})}\left[R(\mathbf{c}, \mathbf{x})\right]$



Output

LLM

Input

Prompt $\mathbf{c}$      Generation $\mathbf{x}$

Supervision

# REINFORCE Leave One Out

- $$\nabla_\theta E_{\mathbf{x} \sim p_\theta(\cdot | \mathbf{c})} \left[ R(\mathbf{c}, \mathbf{x}) \right] \approx \frac{1}{K} \sum_{k=1}^{K} \left( R(\mathbf{c}, \mathbf{x}_k) - b(\mathbf{c}) \right) \nabla_\theta \log p_\theta(\mathbf{x}_k | \mathbf{c})$$

  for $\mathbf{x}_k \sim p_\theta(\cdot | \mathbf{c})$

- $$b(\mathbf{c}) = \frac{1}{K} \sum_{k=1}^{K} R(\mathbf{c}, \mathbf{x}_k)$$

Initialize $\theta$

for ever:

    Sample (or iterate over) $\mathbf{c}$

    $\mathbf{x}_k \sim p_\theta(\cdot | \mathbf{c})$ for k=1...K

    $$b(\mathbf{c}) = \frac{1}{K} \sum_{k=1}^{K} R(\mathbf{c}, \mathbf{x}_k)$$

    $$\theta \leftarrow \theta + \epsilon \frac{1}{K} \sum_{k=1}^{K} \left( R(\mathbf{c}, \mathbf{x}_k) - b(\mathbf{x}) \right) \nabla \log p_\theta(\mathbf{x}_k | \mathbf{c})$$

Buy 4 REINFORCE Samples, Get a Baseline for Free!, Kool etal 2019

# RLOO in LLMs



TL;DR Summarize — Anthropic-HH (Pythia) — Anthropic-HH (Llama)

Legend: RLOO - k=2 · RAFT - k=2 · REINFORCE w/ baseline · PPO · Vanila PG

Back to Basics: Revisiting REINFORCE Style Optimization for Learning from Human Feedback in LLMs, Ahmadian etal 2024

# DeepSeekMath: Pushing the Limits of Mathematical Reasoning in Open Language Models

Zhihong Shao[1,2*†], Peiyi Wang[1,3*†], Qihao Zhu[1,3*†], Runxin Xu[1], Junxiao Song[1]
Xiao Bi[1], Haowei Zhang[1], Mingchuan Zhang[1], Y.K. Li[1], Y. Wu[1], Daya Guo[1*]

[1]DeepSeek-AI, [2]Tsinghua University, [3]Peking University

{zhihongshao,wangpeiyi,zhuqh,guoday}@deepseek.com
https://github.com/deepseek-ai/DeepSeek-Math

# GRPO

$$\mathcal{J}_{PPO}(\theta) = \mathbb{E}[q \sim P(Q), o \sim \pi_{\theta_{old}}(O|q)] \frac{1}{|o|} \sum_{t=1}^{|o|} \min\left[ \frac{\pi_\theta(o_t|q, o_{<t})}{\pi_{\theta_{old}}(o_t|q, o_{<t})} A_t, \text{clip}\left( \frac{\pi_\theta(o_t|q, o_{<t})}{\pi_{\theta_{old}}(o_t|q, o_{<t})}, 1-\varepsilon, 1+\varepsilon \right) A_t \right],$$



$$\mathcal{J}_{GRPO}(\theta) = \mathbb{E}[q \sim P(Q), \{o_i\}_{i=1}^{G} \sim \pi_{\theta_{old}}(O|q)]$$

$$\frac{1}{G} \sum_{i=1}^{G} \frac{1}{|o_i|} \sum_{t=1}^{|o_i|} \left\{ \min\left[ \frac{\pi_\theta(o_{i,t}|q, o_{i,<t})}{\pi_{\theta_{old}}(o_{i,t}|q, o_{i,<t})} \hat{A}_{i,t}, \text{clip}\left( \frac{\pi_\theta(o_{i,t}|q, o_{i,<t})}{\pi_{\theta_{old}}(o_{i,t}|q, o_{i,<t})}, 1-\varepsilon, 1+\varepsilon \right) \hat{A}_{i,t} \right] - \beta \mathbb{D}_{KL}\left[ \pi_\theta || \pi_{ref} \right] \right\},$$

$$\hat{A}_{i,t} = \widetilde{r}_i = \frac{r_i - \text{mean}(\mathbf{r})}{\text{std}(\mathbf{r})}$$

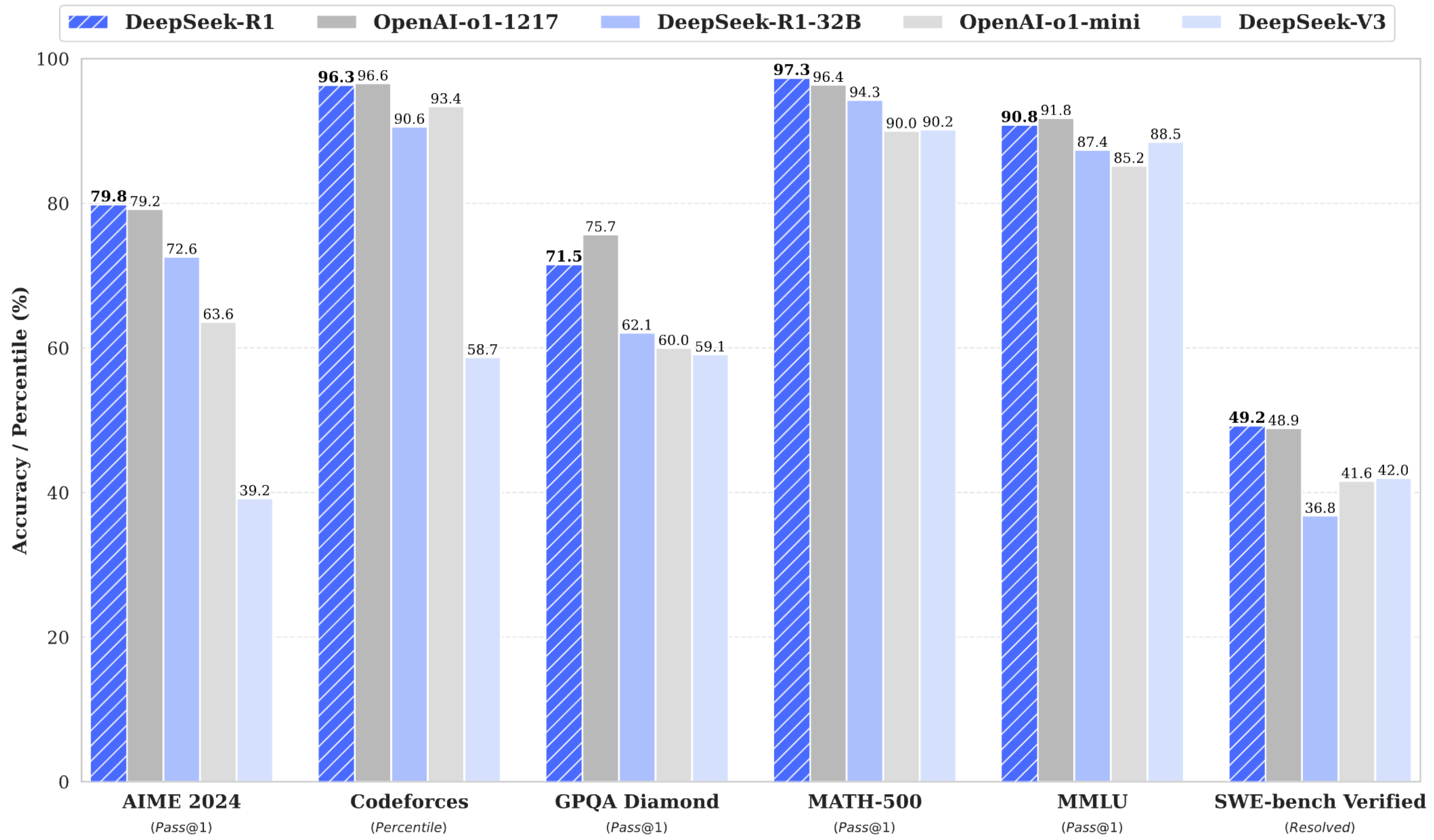DeepSeekMath: Pushing the Limits of Mathematical Reasoning in Open Language Models, Shao etal 2024

**Algorithm 1** Iterative Group Relative Policy Optimization

---

**Input** initial policy model $\pi_{\theta_{\text{init}}}$; reward models $r_\varphi$; task prompts $\mathcal{D}$; hyperparameters $\varepsilon, \beta, \mu$

1:  policy model $\pi_\theta \leftarrow \pi_{\theta_{\text{init}}}$
2:  **for** iteration = $1, \ldots, $ I **do**
3:      reference model $\pi_{ref} \leftarrow \pi_\theta$
4:      **for** step = $1, \ldots, $ M **do**
5:          Sample a batch $\mathcal{D}_b$ from $\mathcal{D}$
6:          Update the old policy model $\pi_{\theta_{old}} \leftarrow \pi_\theta$
7:          Sample $G$ outputs $\{o_i\}_{i=1}^{G} \sim \pi_{\theta_{old}}(\cdot \mid q)$ for each question $q \in \mathcal{D}_b$
8:          Compute rewards $\{r_i\}_{i=1}^{G}$ for each sampled output $o_i$ by running $r_\varphi$
9:          Compute $\hat{A}_{i,t}$ for the $t$-th token of $o_i$ through group relative advantage estimation.
10:         **for** GRPO iteration = $1, \ldots, \mu$ **do**
11:             Update the policy model $\pi_\theta$ by maximizing the GRPO objective (Equation 21)
12:     Update $r_\varphi$ through continuous training using a replay mechanism.

**Output** $\pi_\theta$

---

DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning, Deepseek team 2025

# Step 1: Create a dataset of math puzzles and alike

- Interesting prompts

- Easily verifiable answers

  - Math, Reasoning, Multiple-choice, ...

- No details given in paper

Question: If $a > 1$, then the sum of the real solutions of $\sqrt{a - \sqrt{a + x}} = x$ is equal to

$$\frac{-1 \pm \sqrt{1 + 4a}}{2}$$

# Step 2: Run GRPO



- No teacher forcing

- Supervise just easily verifiable answers

- Learns reasoning

- R1-Zero

  - From just pre-trained model, no instruction tuning

  - It works

# Step 3: Bootstrap Instruction tuned model

- Use instruction tuning data and R1-Zero data

- Train a chat-bot that can "reason"

| Benchmark (Metric) | Claude-3.5-Sonnet-1022 | GPT-4o 0513 | DeepSeek V3 | OpenAI o1-mini | OpenAI o1-1217 | DeepSeek R1 |
|---|---|---|---|---|---|---|
| Architecture | - | - | MoE | - | - | MoE |
| # Activated Params | - | - | 37B | - | - | 37B |
| # Total Params | - | - | 671B | - | - | 671B |
| **English** MMLU (Pass@1) | 88.3 | 87.2 | 88.5 | 85.2 | **91.8** | 90.8 |
| MMLU-Redux (EM) | 88.9 | 88.0 | 89.1 | 86.7 | - | **92.9** |
| MMLU-Pro (EM) | 78.0 | 72.6 | 75.9 | 80.3 | - | **84.0** |
| DROP (3-shot F1) | 88.3 | 83.7 | 91.6 | 83.9 | 90.2 | **92.2** |
| IF-Eval (Prompt Strict) | **86.5** | 84.3 | 86.1 | 84.8 | - | 83.3 |
| GPQA Diamond (Pass@1) | 65.0 | 49.9 | 59.1 | 60.0 | **75.7** | 71.5 |
| SimpleQA (Correct) | 28.4 | 38.2 | 24.9 | 7.0 | **47.0** | 30.1 |
| FRAMES (Acc.) | 72.5 | 80.5 | 73.3 | 76.9 | - | **82.5** |
| AlpacaEval2.0 (LC-winrate) | 52.0 | 51.1 | 70.0 | 57.8 | - | **87.6** |
| ArenaHard (GPT-4-1106) | 85.2 | 80.4 | 85.5 | 92.0 | - | **92.3** |
| **Code** LiveCodeBench (Pass@1-COT) | 38.9 | 32.9 | 36.2 | 53.8 | 63.4 | **65.9** |
| Codeforces (Percentile) | 20.3 | 23.6 | 58.7 | 93.4 | **96.6** | 96.3 |
| Codeforces (Rating) | 717 | 759 | 1134 | 1820 | **2061** | 2029 |
| SWE Verified (Resolved) | **50.8** | 38.8 | 42.0 | 41.6 | 48.9 | 49.2 |
| Aider-Polyglot (Acc.) | 45.3 | 16.0 | 49.6 | 32.9 | **61.7** | 53.3 |
| **Math** AIME 2024 (Pass@1) | 16.0 | 9.3 | 39.2 | 63.6 | 79.2 | **79.8** |
| MATH-500 (Pass@1) | 78.3 | 74.6 | 90.2 | 90.0 | 96.4 | **97.3** |
| CNMO 2024 (Pass@1) | 13.1 | 10.8 | 43.2 | 67.6 | - | **78.8** |
| **Chinese** CLUEWSC (EM) | 85.4 | 87.9 | 90.9 | 89.9 | - | **92.8** |
| C-Eval (EM) | 76.7 | 76.0 | 86.5 | 68.9 | - | **91.8** |
| C-SimpleQA (Correct) | 55.4 | 58.7 | **68.0** | 40.3 | - | 63.7 |

**Algorithm 1** Iterative Group Relative Policy Optimization

---

**Input** initial policy model $\pi_{\theta_{\text{init}}}$; reward models $r_{\varphi}$; task prompts $\mathcal{D}$; hyperparameters $\varepsilon, \beta, \mu$

1: policy model $\pi_\theta \leftarrow \pi_{\theta_{\text{init}}}$
2: **for** iteration = $1, \ldots, I$ **do**
3:      reference model $\pi_{ref} \leftarrow \pi_\theta$
4:      **for** step = $1, \ldots, M$ **do**
5:          Sample a batch $\mathcal{D}_b$ from $\mathcal{D}$
6:          Update the old policy model $\pi_{\theta_{old}} \leftarrow \pi_\theta$
7:          Sample $G$ outputs $\{o_i\}_{i=1}^{G} \sim \pi_{\theta_{old}}(\cdot \mid q)$ for each question $q \in \mathcal{D}_b$
8:          Compute rewards $\{r_i\}_{i=1}^{G}$ for each sampled output $o_i$ by running $r_{\varphi}$
9:          Compute $\hat{A}_{i,t}$ for the $t$-th token of $o_i$ through group relative advantage estimation.
10:          **for** GRPO iteration = $1, \ldots, \mu$ **do**
11:              Update the policy model $\pi_\theta$ by maximizing the GRPO objective (Equation 21)
12:      Update $r_{\varphi}$ through continuous training using a replay mechanism.

**Output** $\pi_\theta$

---

## Algorithm 1 Iterative Group Relative Policy Optimization

**Input** initial policy model $\pi_{\theta_{init}}$; reward models $r_\varphi$; task prompts $\mathcal{D}$; hyperparameters $\varepsilon, \beta, \mu$

1: policy model $\pi_\theta \leftarrow \pi_{\theta_{init}}$
2: **for** iteration = 1, ..., I **do**
3:      reference model $\pi_{ref} \leftarrow \pi_\theta$
4:      **for** step = 1, ..., M **do**
5:          Sample a batch $\mathcal{D}_b$ from $\mathcal{D}$
6:          Update the old policy model $\pi_{\theta_{old}} \leftarrow \pi_\theta$
7:          Sample $G$ outputs $\{o_i\}_{i=1}^G \sim \pi_{\theta_{old}}(\cdot \mid q)$ for each question $q \in \mathcal{D}_b$
8:          Compute rewards $\{r_i\}_{i=1}^G$ for each sampled output $o_i$ by running $r_\varphi$
9:          Compute $\hat{A}_{i,t}$ for the $t$-th token of $o_i$ through group relative advantage estimation.
10:         ~~for GRPO iteration = 1, ..., $\mu$ do~~
11:          Update the policy model $\pi_\theta$ by maximizing the GRPO objective (Equation 21)
12:     Update $r_\varphi$ through continuous training using a replay mechanism.

**Output** $\pi_\theta$

$$\hat{A}_{i,t} = \widetilde{r}_i = \frac{r_i - \text{mean}(\mathbf{r})}{\text{std}(\mathbf{r})}$$

max length is set to 1024, and the training batch size is 1024. The policy model only has a single update following each exploration stage. We evaluate DeepSeekMath-RL 7B on benchmarks

$$\mathcal{J}_{GRPO}(\theta) = \mathbb{E}[q \sim P(Q), \{o_i\}_{i=1}^G \sim \pi_{\theta_{old}}(O|q)]$$

$$\frac{1}{G}\sum_{i=1}^G \frac{1}{|o_i|}\sum_{t=1}^{|o_i|}\left\{\min\left[\frac{\pi_\theta(o_{i,t}|q,o_{i,<t})}{\pi_{\theta_{old}}(o_{i,t}|q,o_{i,<t})}\hat{A}_{i,t}, \text{clip}\left(\frac{\pi_\theta(o_{i,t}|q,o_{i,<t})}{\pi_{\theta_{old}}(o_{i,t}|q,o_{i,<t})}, 1-\varepsilon, 1+\varepsilon\right)\hat{A}_{i,t}\right] - \beta\mathbb{D}_{KL}\left[\pi_\theta||\pi_{ref}\right]\right\},$$

GRPO = RLOO with advantage normalization

DeepSeekMath: Pushing the Limits of Mathematical Reasoning in Open Language Models, Shao etal 2024

# RLOO vs GRPO

## RLOO

Initialize $\theta$

for ever:

    Sample (or iterate over) $\mathbf{c}$

$$\mathbf{x}_k \sim p_\theta(\,\cdot\mid\mathbf{c})\ \text{for k=1...K}$$

$$b(\mathbf{c}) = \frac{1}{K}\sum_{k=1}^{K} R(\mathbf{c}, \mathbf{x}_k)$$

$$\theta \leftarrow \theta + \epsilon\frac{1}{K}\sum_{k=1}^{K}(R(\mathbf{c}, \mathbf{x}_k) - b(\mathbf{x}))\,\nabla\log p_\theta(\mathbf{x}_k\mid\mathbf{c})$$

## GRPO (in practice)

Initialize $\theta$

for ever:

    Sample (or iterate over) $\mathbf{c}$

$$\mathbf{x}_k \sim p_\theta(\,\cdot\mid\mathbf{c})\ \text{for k=1...K}$$

$$b(\mathbf{c}) = \frac{1}{K}\sum_{k=1}^{K} R(\mathbf{c}, \mathbf{x}_k)$$

$$\theta \leftarrow \theta + \epsilon\frac{1}{K}\sum_{k=1}^{K}\frac{R(\mathbf{c}, \mathbf{x}_k) - b(\mathbf{x})}{\text{std}(\mathbf{x})}\,\nabla\log p_\theta(\mathbf{x}_k\mid\mathbf{c})$$

# Interactive Digital Agents

- Train LLMs that interact with API's on the users behalf



**User:** I owe Alice and Bob money for the trip to Maui. They texted me to let me know the amount. Venmo them what I owe.

**Agent:**
```python
# Let's find Alice's recent messages
messages = phone.search_text_messages("Alice", "$")
print(messages)
```

**API:**
```
[
    {"from": "Alice A", "content": "Gym at 8am?", ...},
    {"from": "Alice B", "content": "You owe me for a
third of the $330 hotel bill.", "number": ...}
    ...
]
```

**Agent:**
```python
# It looks like we owe Alice B one third of $330
alice_amount = 330.0 / 3
alice_phone = messages[1]["number"]

# Let's send Alice her money
alice_venmo = venmo.search_users(alice_phone)
venmo.create_transaction(alice_venmo, alice_amount)
...
```

**API:** $110.0 sent to Alice B

**Agent:** # Let's find Bob's recent messages ...

⋮

**Agent:** Done. I sent Alice B $110 and Bob C $50.

Reinforcement Learning for Long-Horizon Interactive LLM Agents, Chen … Krähenbühl 2025

# Interactive Digital Agents

## LOOP

- Use simulator for API interactions

  - AppWorld

- Training data = 24 scenarios (simple request, initial simulator state, test cases)

- Trained using simple combination of Leave-One-Out estimator and PPO

**Algorithm 1** Leave-One-Out Proximal Policy Optimization

**Input:** Policy $p_\theta$, dataset of tasks and initial states $\mathcal{D}$
**Output:** Policy $p_\theta$ maximizing $\mathbb{E}_{\mathbf{s}_0,\mathbf{c}\sim\mathcal{D}}\left[L_\theta(\mathbf{s}_0,\mathbf{c})\right]$ (Eq. 7)

1: **for** iteration $= 1, 2, \ldots$ **do**
2:    $\mathbf{B} \leftarrow \{\}$       ▷ *Initialize rollout buffer*
3:    **for** $(\mathbf{s}_0, \mathbf{c}) \sim \mathcal{D}$ **do**     ▷ *Rollout collection*
4:       Collect $K$ rollouts $\mathbf{x}_1, \ldots, \mathbf{x}_K \overset{\text{i.i.d.}}{\sim} \rho_\theta(\cdot|\mathbf{s}_0, \mathbf{c})$
5:       Estimate advantages $A_1, \ldots, A_K$ using Eq. 3
6:       $\mathbf{B} \leftarrow \mathbf{B} \cup \{(\mathbf{x}_1, A_1), \ldots, (\mathbf{x}_K, A_K)\}$
7:    **for** epoch $= 1, \ldots, N_{\text{epoch}}$ **do**   ▷ *Policy update*
8:       **for** mini-batch $\{(\mathbf{x}_i, A_i)\}_{i=1}^M \sim \mathbf{B}$ **do**
9:          Update policy using PPO gradient (Eq. 5)

$$A(\mathbf{c}, \mathbf{x}_k) = \frac{K}{K-1}\left(R(\mathbf{c}, \mathbf{x}_k) - \frac{1}{K}\sum_{i=1}^{K} R(\mathbf{c}, \mathbf{x}_i)\right). \quad (3)$$
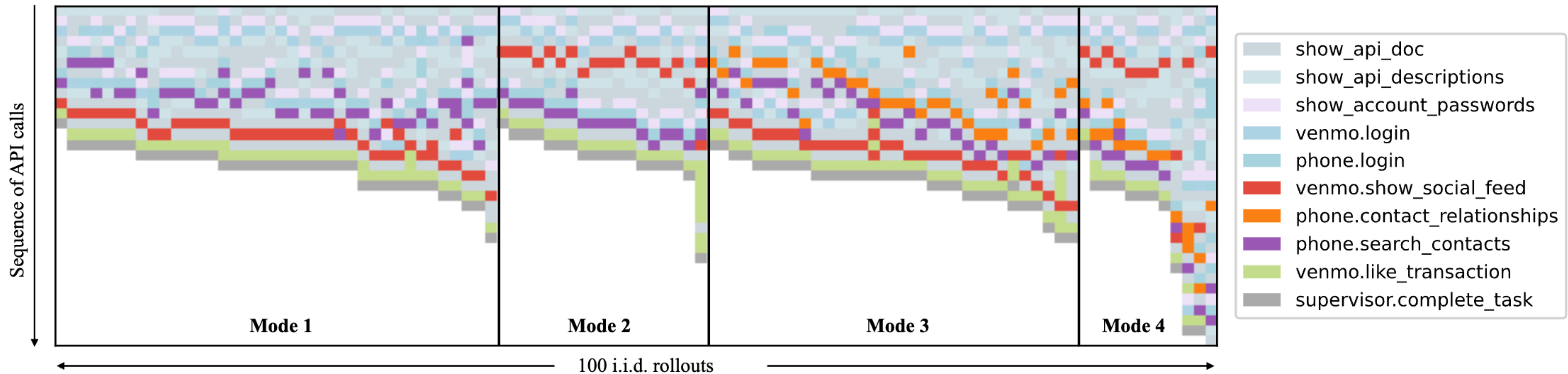
$$L_\theta^{\text{MDP}}(\mathbf{c}) =$$
$$\mathbb{E}_{\mathbf{x}\sim p_\psi(\cdot|\mathbf{c})}\left[\frac{1}{|\mathbf{x}|}\sum_{t=1}^{|\mathbf{x}|}\min\left(\frac{p_\theta(x_t|\mathbf{c},x_{1:t-1})}{p_\psi(x_t|\mathbf{c},x_{1:t-1})}A(\mathbf{c},\mathbf{x}), g_\epsilon(A(\mathbf{c},\mathbf{x}))\right)\right].$$
$$(5)$$

Reinforcement Learning for Long-Horizon Interactive LLM Agents, Chen ... Krähenbühl 2025

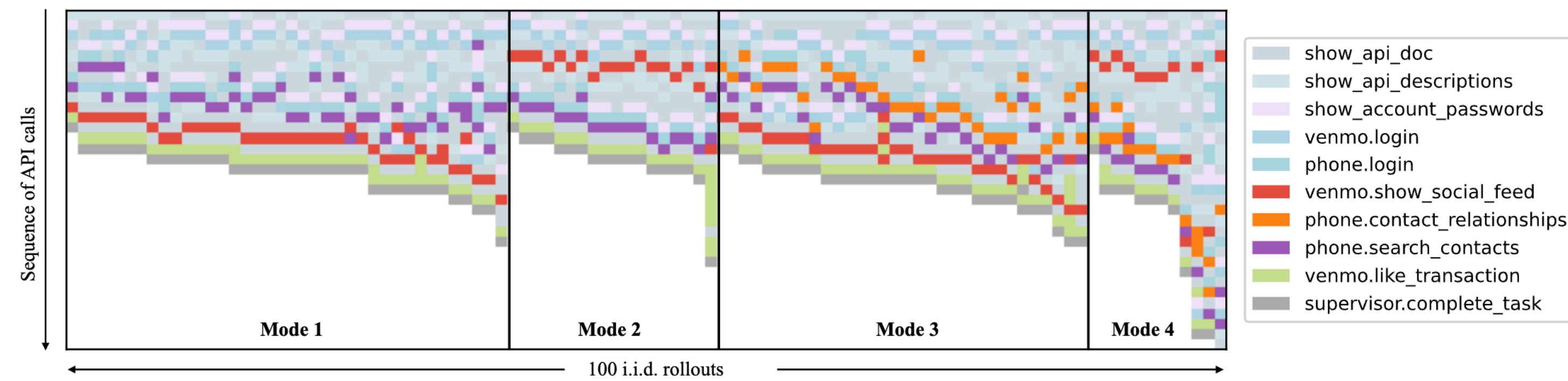| Type | Algorithm | Action | Strictly on-policy | Normalized reward | Test Normal (Test-N) | | Test Challenge (Test-C) | |
|------|-----------|--------|--------------------|-------------------|----------------------|----|------------------------|----|
| | | | | | TGC | SGC | TGC | SGC |
| NFT | GPT-4o | – | – | – | 48.8 | 32.1 | 30.2 | 13 |
| NFT | OpenAI o1 | – | – | – | 61.9 | 41.1 | 36.7 | 19.4 |
| NFT | Llama 3 70B | – | – | – | 24.4 | 17.9 | 7.0 | 4.3 |
| NFT | Qwen 2.5 32B | – | – | – | $39.2 \pm 3.5$ | $18.6 \pm 2.0$ | $21.0 \pm 1.4$ | $7.5 \pm 1.2$ |
| SFT | SFT-GT | – | – | – | $6.2 \pm 0.7$ | $1.8 \pm 0.0$ | $0.8 \pm 0.2$ | $0.1 \pm 0.3$ |
| SFT | RFT | – | – | – | $47.9 \pm 3.7$ | $26.4 \pm 2.3$ | $26.4 \pm 1.8$ | $11.4 \pm 2.3$ |
| SFT | EI | – | – | – | $58.3 \pm 2.8$ | $36.8 \pm 6.0$ | $32.8 \pm 0.7$ | $17.6 \pm 1.3$ |
| DPO | DPO-MCTS | – | – | – | $57.0 \pm 1.5$ | $31.8 \pm 4.2$ | $31.8 \pm 1.3$ | $13.7 \pm 1.5$ |
| DPO | DMPO | – | – | – | $59.0 \pm 1.2$ | $36.6 \pm 4.7$ | $36.3 \pm 1.8$ | $18.4 \pm 2.3$ |
| RL | PPO (learned critic) | token | | | $50.8 \pm 3.7$ | $28.9 \pm 7.9$ | $26.4 \pm 0.5$ | $10.5 \pm 2.1$ |
| RL | RLOO | traj | ✓ | | $57.2 \pm 2.6$ | $35.7 \pm 2.9$ | $36.7 \pm 1.6$ | $17.4 \pm 1.4$ |
| RL | GRPO | token | ✓[3] | ✓ | $58.0 \pm 1.8$ | $36.8 \pm 3.9$ | $39.5 \pm 1.9$ | $22.4 \pm 0.8$ |
| RL | GRPO no kl | token | ✓[3] | ✓ | $59.0 \pm 1.4$ | $35.7 \pm 2.9$ | $42.7 \pm 1.3$ | $21.3 \pm 1.7$ |
| RL | LOOP (bandit) | traj | | | $53.3 \pm 3.4$ | $33.6 \pm 3.2$ | $27.7 \pm 1.5$ | $13.0 \pm 0.9$ |
| RL | LOOP (turn) | turn | | | $64.1 \pm 2.2$ | $43.5 \pm 3.5$ | $40.8 \pm 1.5$ | $26.5 \pm 2.4$ |
| RL | **LOOP (token)** | token | | | $\mathbf{71.3 \pm 1.3}$ | $\mathbf{53.6 \pm 2.2}$ | $\mathbf{45.7 \pm 1.3}$ | $\mathbf{26.6 \pm 1.5}$ |
| RL | LOOP RwNorm (token) | token | | ✓ | $61.9 \pm 4.0$ | $44.1 \pm 7.8$ | $39.8 \pm 1.3$ | $20.4 \pm 2.1$ |

# RL vs SFT

- Generations from LLM are very diverse even after RL training

# RL vs SFT

# RL vs SFT



- Generations from LLM are very diverse even after RL training

  - 98 / 100 solutions are correct

  - Only 3 repeat overall structure of commands (no exact repetition)

- Early in training: Awesome exploration

- Late in training: No collapse / overfitting

Reinforcement Learning for Long-Horizon Interactive LLM Agents, Chen ... Krähenbühl 2025

# Reinforcement Learning and LLMs

- Easy to implement

  - RLOO, LOOP are just generation + reward computation + mean subtraction + training

  - Great open-source tools exist for all of this

- Much more flexible

- Less data-hungry

- It is here to stay

Initialize $\theta$
for ever:
    Sample (or iterate over) $\mathbf{c}$
$$\mathbf{x}_k \sim p_\theta(\,\cdot\,|\,\mathbf{c})\ \text{ for k=1...K}$$
$$b(\mathbf{c}) = \frac{1}{K}\sum_{k=1}^{K} R(\mathbf{c}, \mathbf{x}_k)$$
$$\theta \leftarrow \theta + \epsilon\frac{1}{K}\sum_{k=1}^{K} (R(\mathbf{c}, \mathbf{x}_k) - b(\mathbf{x}))\,\nabla \log p_\theta(\mathbf{x}_k|\,\mathbf{c})$$

# References

- Back to Basics: Revisiting REINFORCE Style Optimization for Learning from Human Feedback in LLMs, Ahmadian etal. 2024

- Buy 4 REINFORCE Samples, Get a Baseline for Free!, Kohl etal 2019

- DeepSeekMath: Pushing the Limits of Mathematical Reasoning in Open Language Models, Shao etal. 2024

- DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning, DeepSeek-AI 2025

- Reinforcement Learning for Long-Horizon Interactive LLM Agents, Chen etal 2025

# Limitations of LLMs

Philipp Krähenbühl, UT Austin

# Politics of LLM research

- Many different camps

  - With conflicting often hidden motives

**Model Builders**

Develop new models

Make $$$, fame, glory,
(Invent AGI)

**AI Safety research**

Study limitations,
biases, and dangers

Concerns about
societal impacts of
LLMs, fame

**External Analyses**

Bring tools from other sciences
into LLM world

Study LLMs as "creatures",
More scientific approach,
fame

# ChatGPT is bullshit

- LLMs generate falsehoods

  - AKA Hallucinations

- **Bullshit** (general): Any utterance produced where a speaker has indifference towards the truth of the utterance.

- **Hard** bullshit: Bullshit produced with the **intention to mislead** the audience about the utterer's agenda.

- **Soft** bullshit: Bullshit produced **without the intention to mislead** the hearer regarding the utterer's agenda.

ChatGPT is bullshit, Hicks etal 2024



ON BULLSHIT
*Harry G. Frankfurt*

Bullshitters misrepresent themselves to their audience not as liars do, that is, by deliberately making false claims about what is true. Rather, bullshitters seek to convey a certain impression of themselves without being concerned about whether anything at all is true. - Frankfurt
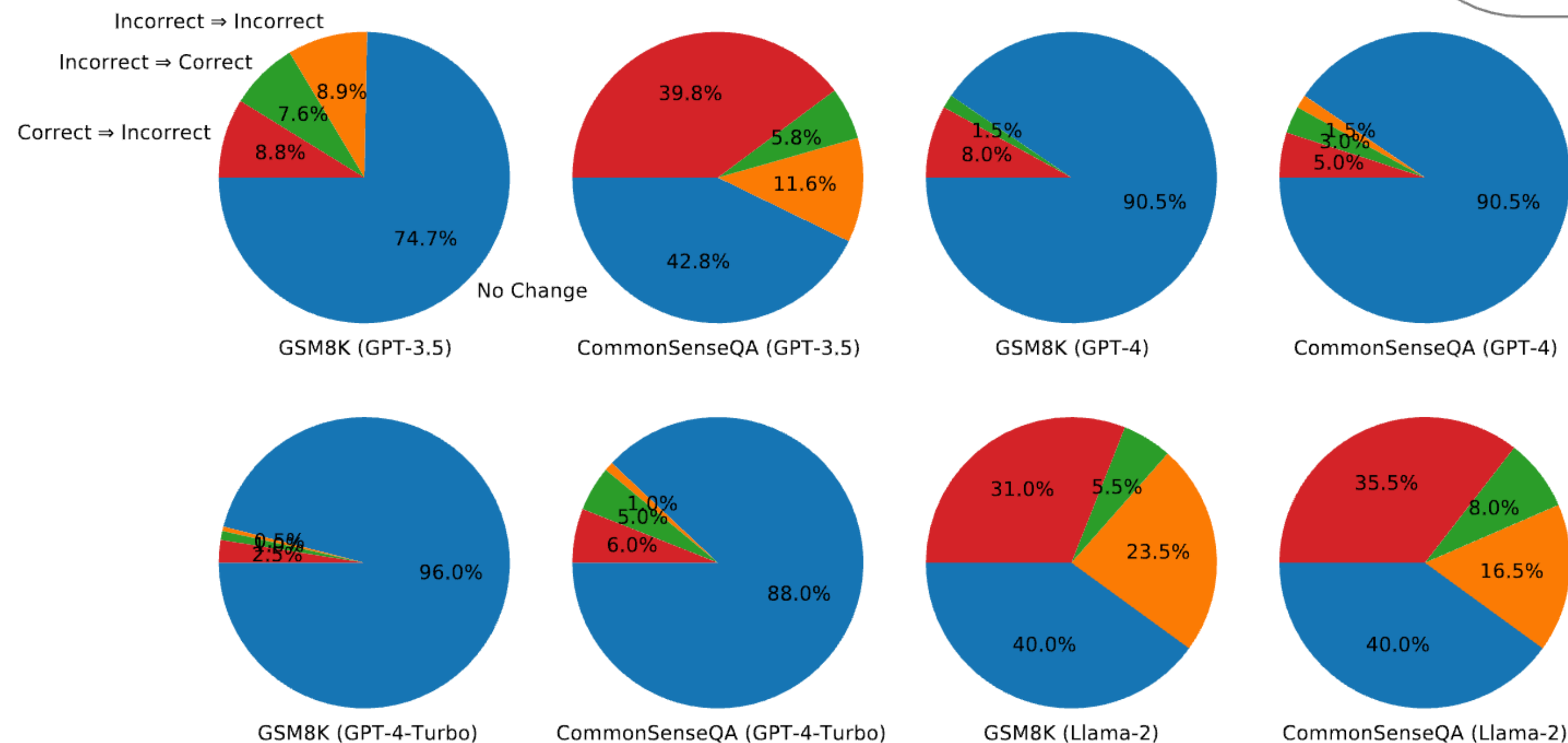
# ChatGPT is bullshit

- "[LLMs] are in an important way indifferent to the truth of their outputs."

  - Training: Reproduce plausible text

  - Human cognition: goals and behaviors

- Paper makes a few assumptions

  - Disregards alignment

The problem here isn't that large language models hallucinate, lie, or misrepresent the world in some way. It's that <u>they are not designed to represent the world at all</u>; instead, they are designed to convey convincing lines of text.

ChatGPT is at minimum a soft bullshitter or a bullshit machine, because if it is not an agent then it can neither hold any attitudes towards truth nor towards deceiving hearers about its (or, perhaps more properly, its users') agenda

ChatGPT is bullshit, Hicks etal 2024

# No self-correction Yet

- LLMs rarely change their mind

  - If they do, they make things worse

Christina is planning a birthday party and needs .75 gift bags per invited guest, because 1/4 of attendees don't show up. She invited 16 friends. Gift bags are $2 each. How much will she spend?

... Therefore, Christina will spend $18 on gift bags. Answer: \boxed{18}. (Incorrect)

Review your previous answer and find problems with your answer.

Upon reviewing my previous answer, I recognize a mistake. I incorrectly stated that Christina needs 0.75 gift bags per invited guest. The given information states that 1/4 of the attendees won't show up…

Based on the problems you found, improve your answer.

... Therefore, Christina will spend $24 on gift bags. Answer: \boxed{24}. (Correct)

Terry eats 2 yogurts a day. They are currently on sale at 4 yogurts for $5.00. How much does he spend on yogurt over 30 days?

... So, Terry spends $\boxed{75.00}$ on yogurt over 30 days. (Correct)

Review your previous answer and find problems with your answer.

Upon reviewing my previous answer, I realized that I made a mistake in calculating Terry's spending on yogurt over 30 days. I incorrectly stated that Terry spends $2.50 per day for 30 days, resulting in a total of $75.00…

Based on the problems you found, improve your answer.

... Therefore, the final answer is Terry spends $\boxed{37.50}$ on yogurt over 30 days. (Incorrect)



Incorrect ⇒ Incorrect
Incorrect ⇒ Correct
Correct ⇒ Incorrect

GSM8K (GPT-3.5): 8.9%, 7.6%, 8.8%, 74.7% (No Change)

CommonSenseQA (GPT-3.5): 39.8%, 5.8%, 11.6%, 42.8%

GSM8K (GPT-4): 1.5%, 8.0%, 90.5%

CommonSenseQA (GPT-4): 1.5%, 3.0%, 5.0%, 90.5%

GSM8K (GPT-4-Turbo): 0.5%, 2.5%, 96.0%

CommonSenseQA (GPT-4-Turbo): 1.0%, 5.0%, 6.0%, 88.0%

GSM8K (Llama-2): 31.0%, 5.5%, 23.5%, 40.0%

CommonSenseQA (Llama-2): 35.5%, 8.0%, 16.5%, 40.0%

Large Language Models Cannot Self-Correct Reasoning Yet, Huang etal 2023
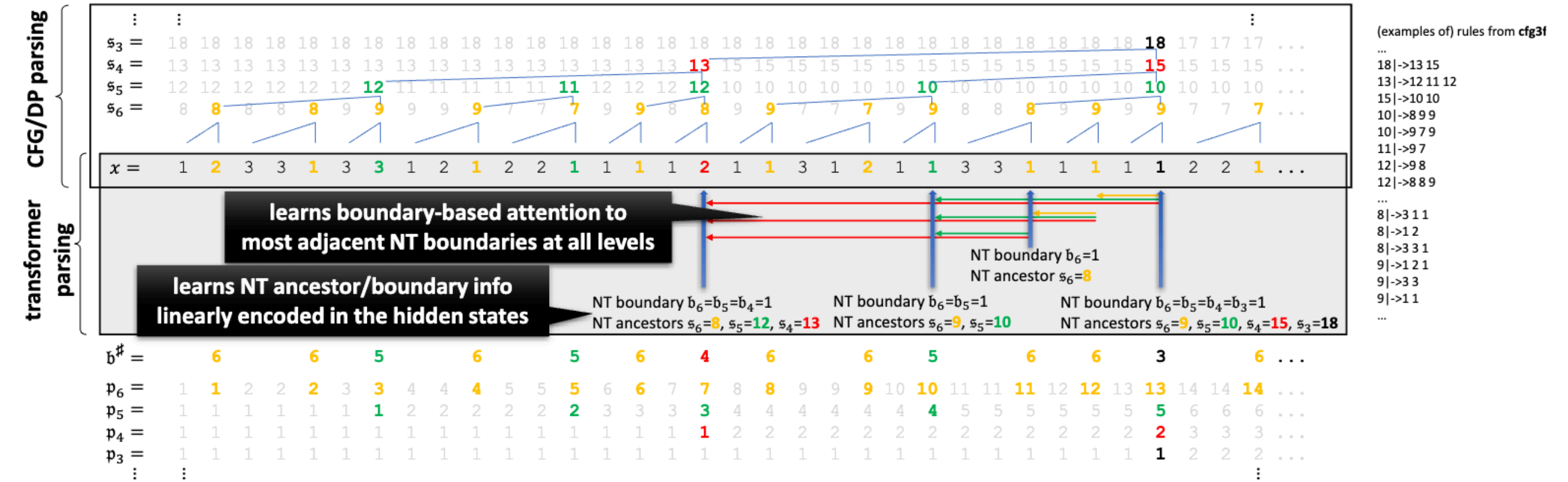
# Conflation of language and thought

- "good at language -> good at thought" fallacy

- formal vs. functional linguistic

- Fairly balanced analysis of current models



**SELECT FORMAL COMPETENCE SKILLS** | **EXAMPLES OF GOOD AND BAD FORMS**

**FORMAL COMPETENCE** getting the form of language right

**phonology**
e.g., rules governing valid wordforms

*blick* could be a valid English word | *\*bnick* could not be a valid English word

**morphology**
e.g., morpheme ordering constraints, rules governing novel morphemic combinations

*Lady Gaga-esque-ness* | *\*Lady Gaga-ness-esque*

**lexical semantics**
e.g., parts of speech, lexical categories, word meanings

I'll take my coffee with cream and *sugar*. | *I'll take my coffee with cream and *red.*

**syntax**
e.g., agreement, word order constraints, constructional knowledge…

The key to the cabinets *is* on the table. | *The key to the cabinets *are* on the table.

**SELECT FUNCTIONAL COMPETENCE SKILLS** | **SUCCESSES/FAILURES IN EACH DOMAIN**

**FUNCTIONAL COMPETENCE** using language to do things in the world

**formal reasoning**
e.g., logic, math, planning

Fourteen birds were sitting on a tree. Three left, one joined. There are now **twelve** birds. | Fourteen birds were sitting on a tree. Three left, one joined. There are now **eleven** birds.

**world knowledge**
e.g., facts, concepts, common sense

The trophy did not fit into the suitcase because **the suitcase** was too small. | The trophy did not fit into the suitcase because **the trophy** was too small.

**situation modeling**
e.g., discourse coherence, narrative structure

Sally owns a dog. **The dog** is black. | Sally doesn't own a dog. **The dog** is black.

**social reasoning**
e.g., pragmatics, theory of mind

Lu put the toy in the box and left. Bo secretly moved it to the closet. Lu now thinks the toy is in the **box.** | Lu put the toy in the box and left. Bo secretly moved it to the closet. Lu now thinks the toy is in the **closet.**

Dissociating language and thought in LLMs, Mahowald etal 2023

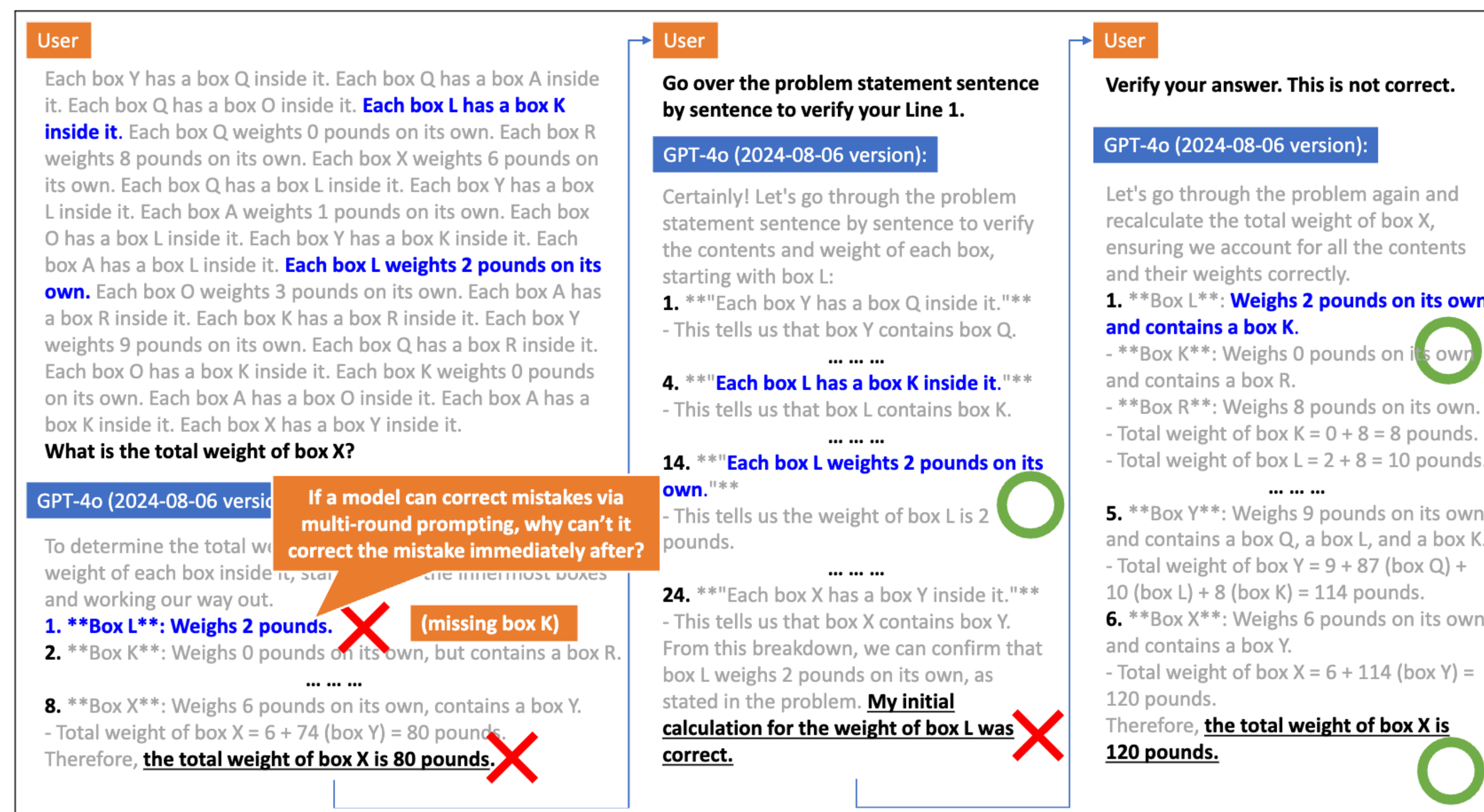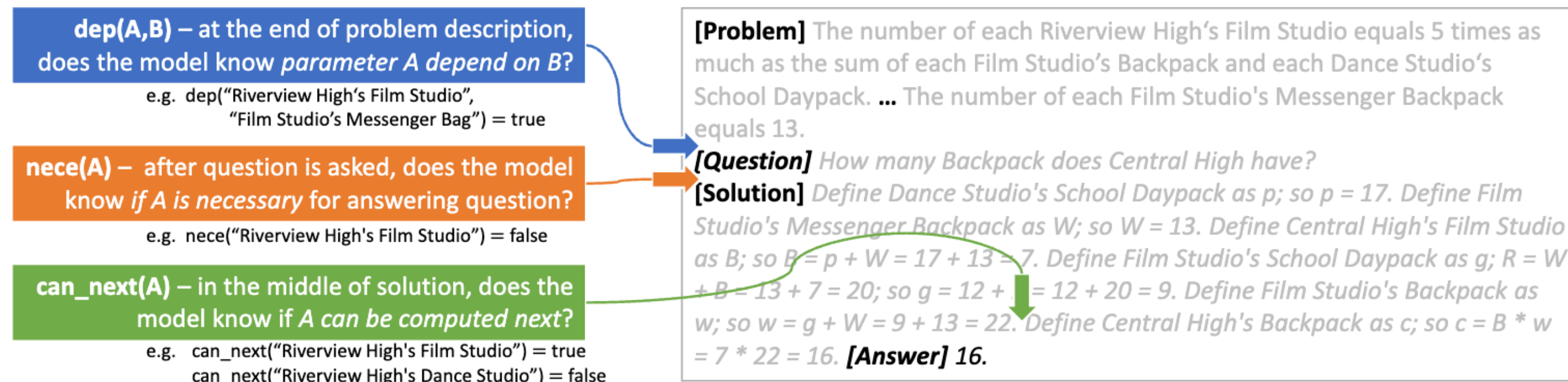# Physics of LLMs

## Limits and capabilities of LLMs



- Large **synthetic** data experiments

- Causal LLMs can learn to parse CFGs

  - Internally use Dynamic Programming-like algorithm

- Bi-directional architectures cannot

# Physics of LLMs
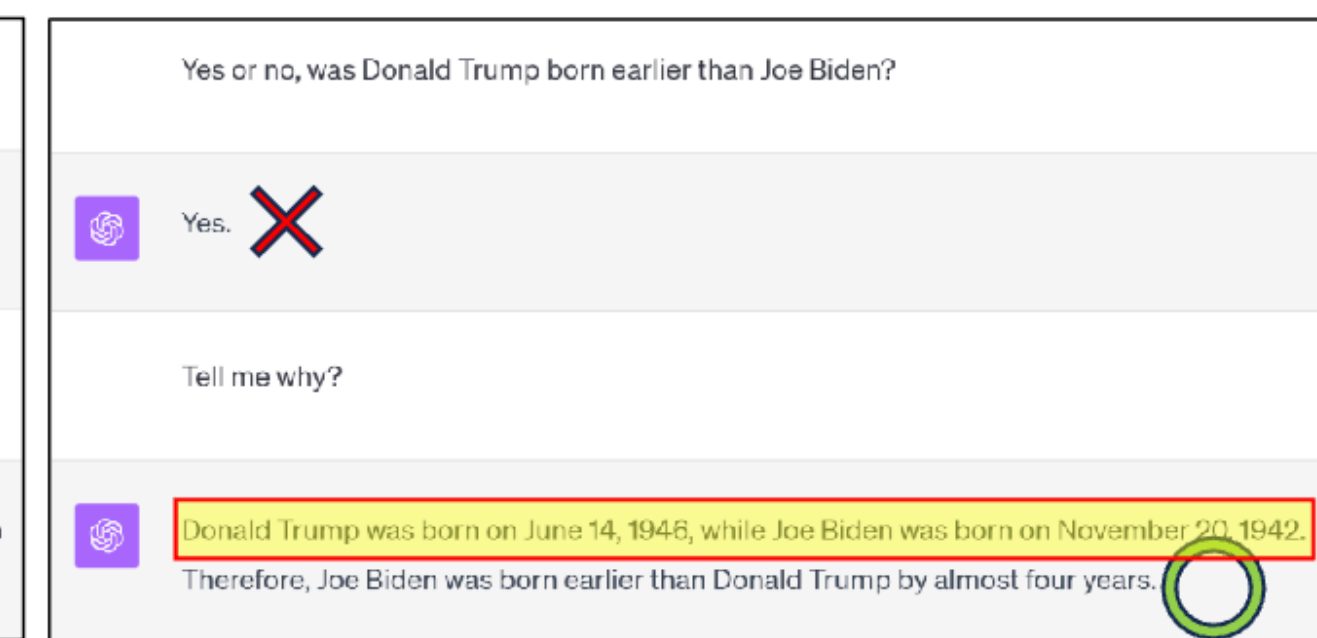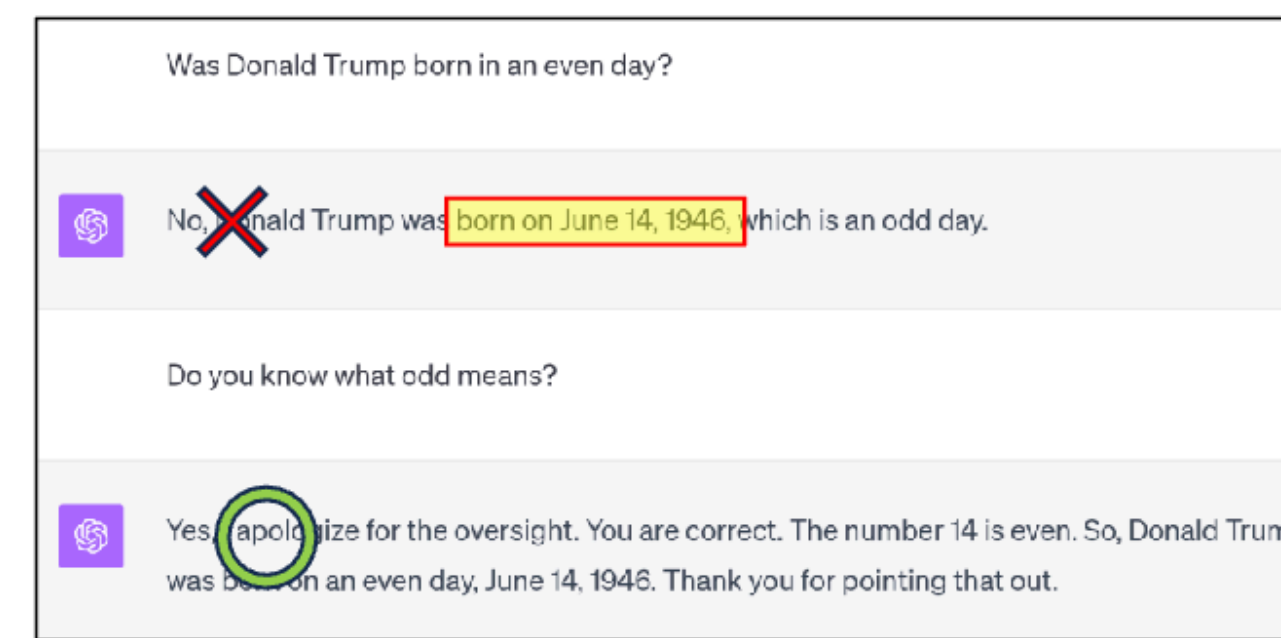## Limits and capabilities of LLMs

- Large **synthetic** data experiments

- LLMs can learn mathematical reasoning (not just memorization)

- LLMs can learn to solve math problems like humans

- Depth (#layers) matters for mathematical reasoning

- LLMs can learn from mistakes if seen during **pre-**training

Physics of Language Models, Allen-Zhu 2023-2024

# Physics of LLMs

## Limits and capabilities of LLMs



- Large **synthetic** data experiments

- Causal LLMs

  - 2 bits of knowledge per parameter, even when quantized to int8

  - Order of knowledge matters (inverse knowledge search often fails)

  - Replication of knowledge in pre-training data is important

- Bi-directional architectures cannot

Physics of Language Models, Allen-Zhu 2023-2024

# Limitations of LLMs

- LLMs are not perfect

- Neither is their analysis

**Model Builders**

Develop new models

Make $$$, fame, glory,
(Invent AGI)

**AI Safety research**

Study limitations,
biases, and dangers

Concerns about
societal impacts of
LLMs, fame

**External Analyses**

Bring tools from other sciences
into LLM world

Study LLMs as "creatures",
More scientific approach,
fame

# References

- [1] ChatGPT is bullshit, Hicks etal 2024

- [2] Large Language Models Cannot Self-Correct Reasoning Yet, Huang etal 2023

- [3] Dissociating language and thought in LLMs, Mahowald etal 2023

- [4] Physics of Language Models, Allen-Zhu 2023-2024