CS11-711 Advanced NLP

# Prompting and Complex Reasoning

Graham Neubig

**Carnegie Mellon University**
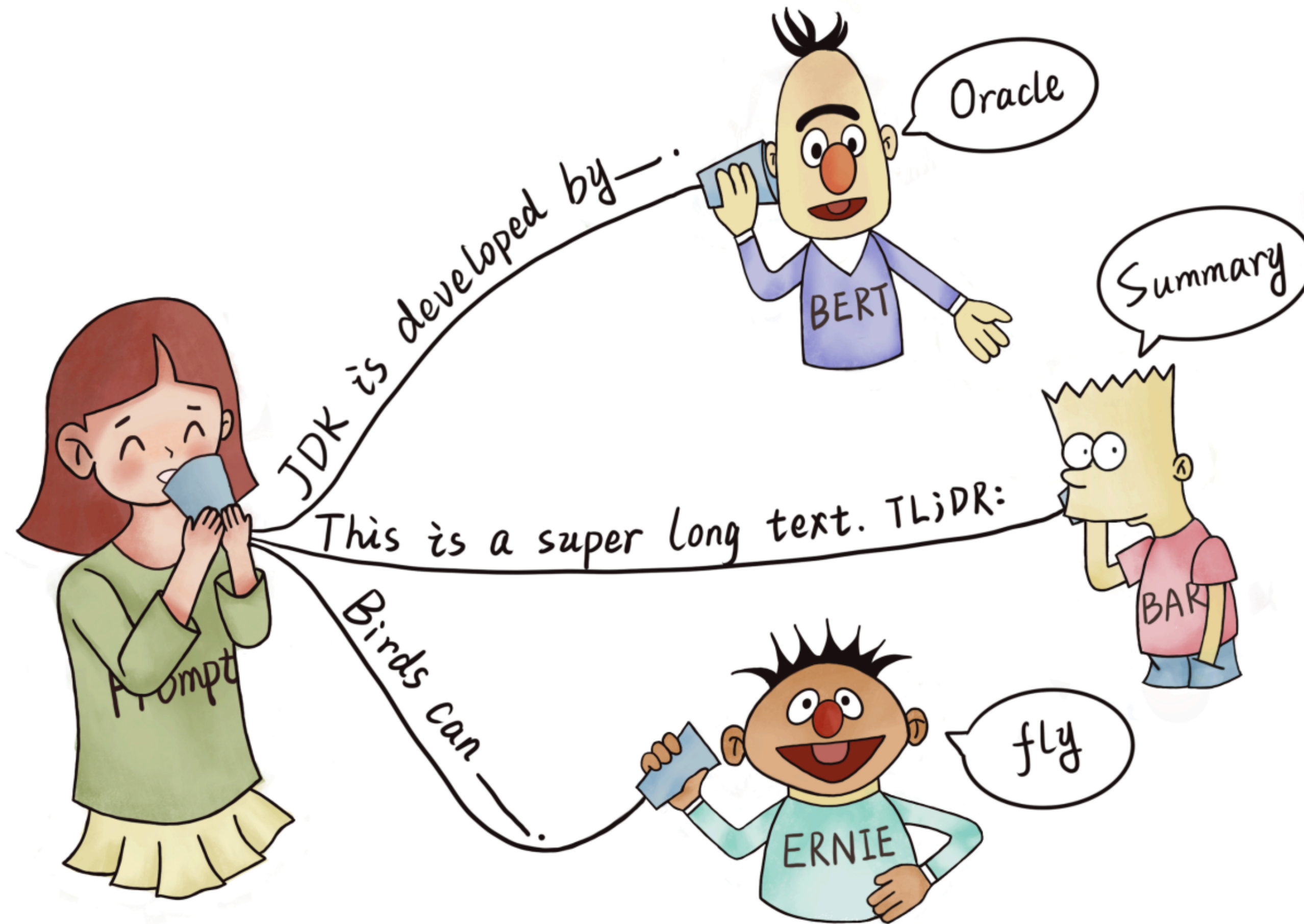**Language Technologies Institute**

https://phontron.com/class/anlp-fall2024/

Some Slides by Pengfei Liu

# What is Prompting ?

☐ Encouraging a pre-trained model to make particular predictions by providing a textual "prompt" specifying the task to be done.

# Prompting Fundamentals

# Basic Prompting (Radford et al. 2018)

- Append a textual string to the beginning of the sequence and complete

x = **When a dog sees a squirrel, it will usually**

(GPT-2 Small)

be afraid of anything unusual. As an exception, that's when a squirrel is usually afraid to bite.

(GPT-2 XL)

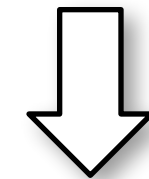lick the squirrel. It will also touch its nose to the squirrel on the tail and nose if it can.

# Standard Prompting Workflow

- Fill a prompt template
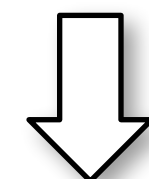- Predict the answer
- Post-process the answer

# Prompt Templates

- A template where you fill in with an actual input

**Input:** x = "I love this movie"

⬇

**Template:** [x] Overall, it was [z]

⬇

**Prompting:** x' = "I love this movie. Overall it was [z]"

# Chat Prompts

- Recently, many models are trained as chatbots
- Usually inputs are specified in OpenAI messages format

```
messages=[
  {
    "role": "system",
    "content": "Please classify movie reviews as 'positive' or 'negative'."
  },
  {
    "role": "user",
    "content": "This movie is a banger."
  },
]
```

- Roles:
  - **"system":** message provided to the system to influence behavior
  - **"user":** message input by the user
  - **"assistant":** message output by the system

# Chat Prompts Behind the Scenes

▪ Behind the scenes, messages are converted to token strings

## LLaMa

Sys.
```
[INST]
<<SYS>>
You are an assistant that …
<</SYS>>
[/INST]
```

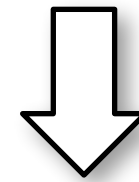User | `[INST]This movie is great.[/INST]`

Asst. | `Positive.`

## Alpaca

Sys. |
```
### Instruction:
You are an assistant that …
```

User |
```
### Instruction:
This movie is great.
```

Asst. |
```
### Response:
Positive.
```

▪ **Software:** See LiteLLM Prompt Templates

# Answer Prediction

- Given a prompt, predict the answer

> **Prompting:** x' = "I love this movie. Overall it was [z]"

⬇

> **Predicting:** x' = "I love this movie. Overall it was fantastic"

- Use any inference algorithms, as in generation class

# Post-processing

- Based on the answer, select the actual output
- For instance:
  - Taking the output as-is
  - Formatting the output for easy visualization
  - Selecting only parts of the output that you want to use
  - Mapping the outputs to other actions

# Output Formatting

## Markdown Rendering

**You**

Please write a table with the last 5 presidents and their birthdates.

**ChatGPT**

Here is a table with the last 5 U.S. Presidents and their birthdates:

| President | Birthdate |
| --- | --- |
| Joe Biden | November 20, 1942 |
| Donald Trump | June 14, 1946 |
| Barack Obama | August 4, 1961 |
| George W. Bush | July 6, 1946 |
| Bill Clinton | August 19, 1946 |

## Code

**You**

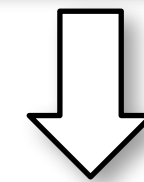Output that in JSON format.

**ChatGPT**

Here is the information in JSON format:

```json
[
  {"President":"Joe Biden","Birthdate":"November 20, 1942"},
  {"President":"Donald Trump","Birthdate":"June 14, 1946"},
  {"President":"Barack Obama","Birthdate":"August 4, 1961"},
  {"President":"George W. Bush","Birthdate":"July 6, 1946"},
  {"President":"Bill Clinton","Birthdate":"August 19, 1946"}
]
``` &#8203;`` 【oaicite:0】 ``&#8203;
```

- For user-facing applications, format in a pretty way

# Output Selection

- From a longer response, select the information indicative of an answer

<div style="border:1px solid; border-radius:10px; padding:10px;">

**Predicting:** x' = "I love this movie. Overall it was a movie that was simply fantastic"

⬇

**Extraction: fantastic**

</div>

- Various methods for extraction

  - **Classification:** identify keywords

  - **Regression/numerical problems:** identify numbers

  - **Code:** pull out code snippets in triple-backticks

# Output Mapping

- Given an answer, map it into a class label or continuous value

**Extraction: fantastic**

⇩

**Mapping: fantastic => Positive**

- Often map many extracted words onto a single class

**Positive** → **Interesting**
**Positive** → **Fantastic**
**Positive** → **Happy**

**Negative** → **Boring**
**Negative** → **1-star**

**...**

# Few-shot Prompting / In-context Learning

# Few-shot Prompting (Brown+ 2021)

- Provide a few examples of the task together with the instruction

| `Please classify movie reviews as 'positive' or 'negative'.`

```
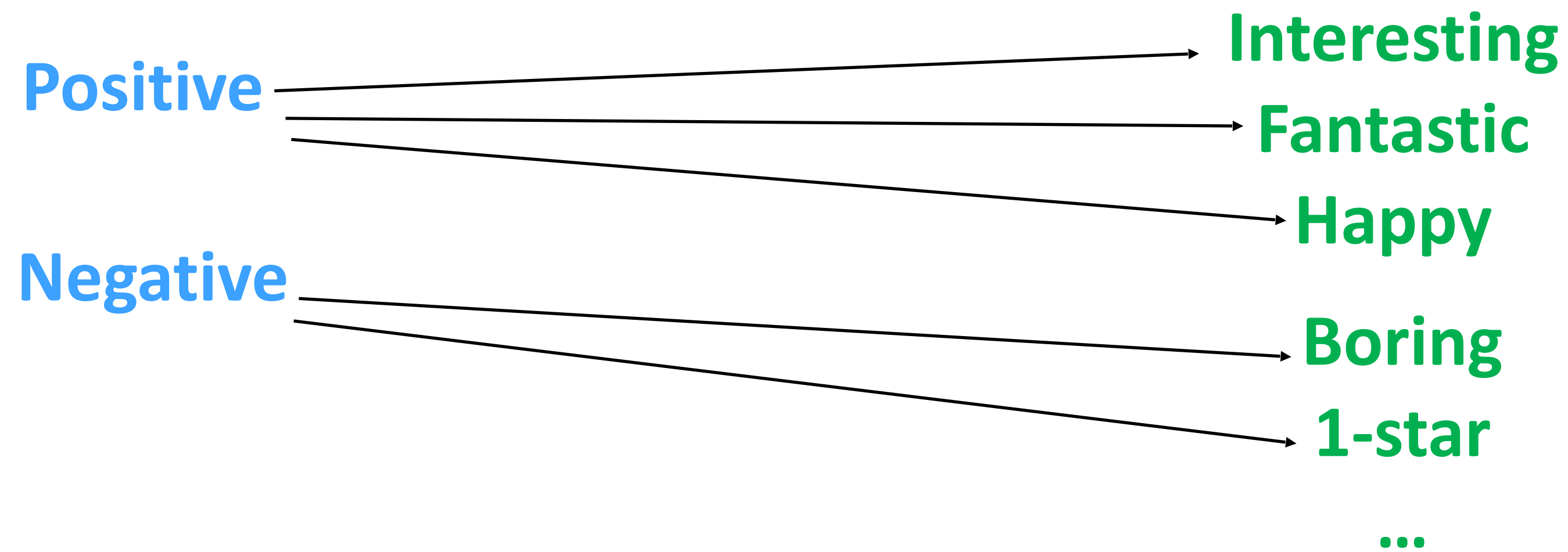Input: I really don't like this movie.
Output: negative

Input: This movie is great!
Output: positive
```

# Few-shot Prompting w/ Chat Prompts

▪ For OpenAI models, add "role": "system" and a "name": "example_assistant" etc.

```
messages=[
  {
    "role": "system",
    "content": "You are an assistant that translates corporate jargon into plain English."
  },
  {
    "role": "system",
    "name":"example_user",
    "content": "New synergies will help drive top-line growth."
  },
  {
    "role": "system",
    "name": "example_assistant",
    "content": "Things working well together will increase revenue."
  },
  ...,
  {
    "role": "user",
    "content": "This late pivot means we don't have time to boil the ocean for the client deliverable."
  },
]
```

# LMs are Sensitive to Small Changes in In-context Examples

- Example ordering (Lu et al. 2021)

- Label balance (Zhang et al. 2022)



Figure 1: Four-shot performance for 24 different sample orders across different sizes of GPT-family models (GPT-2 and GPT-3) for the SST-2 and Subj datasets.

- Label coverage (Zhang et al. 2022)

# But Effects are Sometimes Counter-intuitive (Min et al. 2022)

- Replacing correct labels with random labels sometimes barely hurts accuracy



- More demonstrations can sometimes hurt accuracy

# Prompt Engineering

# Design of Prompts

- Manual
  - □ Configure a manual template based on the characteristics of the task

- Automated search
  - □ Search in discrete space

  - □ Search in continuous space

# Manual Engineering: Format

- Make sure that the format matches that of a trained model (e.g. chat format)

- This can have a large effect on models! (Sclar et al. 2023)

# Manual Engineering: Instructions

- Instructions should be clear, concise and easy to understand
- Good examples: https://www.promptingguide.ai/introduction/tips

Less Precise:
Explain the concept prompt engineering. Keep the explanation short, only a few sentences, and don't be too descriptive.

More Precise:
Use 2-3 sentences to explain the concept of prompt engineering to a high school student.

- Similar to humans, but (right now) LMs don't complain when you're vague

# Methods for Automatic Prompt Engineering

- Prompt paraphrasing
- Gradient-based discrete prompt search
- Prompt tuning
- Prefix tuning

# Prompt Paraphrasing

- Paraphrase an existing prompt to get other candidates (Jiang et al. 2019)

[X] shares a border with [Y]. ⟶ **Paraphrasing Model** ⟶ [X] has a common border with [Y].
[X] adjoins [Y].
......

- Can be done through iterative paraphrasing (Zhou et al. 2021)

# Gradient-based Search (Shin et al. 2020)

- Automatically optimize arbitrary prompts based on existing

# Prompt Tuning (Lester et al. 2021)

- Optimize the embeddings of a prompt, instead of the words.

# Prefix Tuning (Li and Liang 2021)

- "Prompt Tuning" optimizes only the embedding layer
- "Prefix Tuning" optimizes the prefix of all layers

# Prompting and Fine Tuning

# Connection to Other Training Methods

- We will be covering LoRA, Adapters, and other parameter-efficient training methods in a future class

- But prefix tuning is actually very similar, with a few design differences! (He et al. 2021)



(a) Adapter   (b) Prefix Tuning   (c) LoRA   (d) Parallel Adapter   (e) Scaled PA

# A Taxonomy of Prompting Methods

# Prompting as a Prior

- Prompts are a human-interpretable method for providing a prior over the parameters of the model!

- This can be used in place of fine-tuning, or combined with fine-tuning (e.g. Schick and Schütze 2020)

  - More in future classes!

# Reasoning in LMs

# What is Reasoning?

- Using **evidence** and **logic** to arrive at conclusions and make judgments
- Can be
  - □ **Formal:** focusing on strict truth values
  - □ **Informal:** based on experience and intuition

# Types of Reasoning
**(examples: Huang and Chang 2023)**

- Using **evidence** and **logic** to arrive at conclusions and make judgments (Huang and Chang 2023)

**Deductive:** Use logic to go from premise to firm conclusion.

> Premise: All mammals have kidneys.
> Premise: All whales are mammals.
> Conclusion: All whales have kidneys.

**Inductive:** From observation, predict a likely conclusion.

> Observation: When we see a creature with wings, it is usually a bird.
> Observation: We see a creature with wings.
> Conclusion: The creature is likely to be a bird.

**Abductive:** From observation, predict the most likely explanation.

> Observation: The car cannot start and there is a puddle of liquid under the engine.
> Likely Explanation: The car has a leak in the radiator

# Symbolic Reasoning

- Reasoning is something that traditional semantic representations are good at! (e.g. Prolog)

$$\forall x \text{DIE}(x) \rightarrow \text{DIE}(\text{MIA}) \wedge \text{DIE}(\text{ZED})$$

| | | |
|---|---|---|
| 1 | $F(\forall x \text{DIE}(x) \rightarrow \text{DIE}(\text{MIA}) \wedge \text{DIE}(\text{ZED}))$ | |
| 2 | $T \forall x \text{DIE}(x)$ | $1, F_\rightarrow$ |
| 3 | $F(\text{DIE}(\text{MIA}) \wedge \text{DIE}(\text{ZED}))$ | $1, F_\rightarrow$ |
| 4 | $T \text{DIE}(\text{MIA})$ | $2, T_\forall$ |
| 5 | $T \text{DIE}(\text{ZED})$ | $2, T_\forall$ |

6   $F \text{DIE}(\text{MIA})$   $3, F_\wedge$         7   $F \text{DIE}(\text{ZED})$   $3, F_\wedge$

- See "Representation and Inference for Natural Language" (Blackburn & Bos 1999)

# Chain-of-thought Prompting

# Chain of Thought Prompting (Wei et al. 2022)

- Get the model to explain its reasoning before making an answer

## 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. ✓

Figure 1: Chain-of-thought prompting enables large language models to tackle complex arithmetic, commonsense, and symbolic reasoning tasks. Chain-of-thought reasoning processes are highlighted.

- Provides the model with adaptive computation time

# Unsupervised Chain-of-thought Prompting (Kojima et al. 2022)

- Just adding a prompt that encourages the model to explain decisions can induce reasoning

### (a) Few-shot

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

---

(Output) The answer is 8. X

### (b) Few-shot-CoT

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

---

(Output) The juggler can juggle 16 balls. Half of the balls are golf balls. So there are 16 / 2 = 8 golf balls. Half of the golf balls are blue. So there are 8 / 2 = 4 blue golf balls. The answer is 4. ✓

### (c) Zero-shot

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?
A: The answer (arabic numerals) is

---

(Output) 8 X

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

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?
A: **Let's think step by step.**

---

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

- Note: GPT models reason even w/o specific instructions now (probably due to instruction tuning)

# Reasoning is an "Emergent" Ability (Wei et al. 2022)

- Emergent abilities — only appear when models are very large

- Note: emergent abilities are somewhat an artifact of how we measure accuracy (Schaeffer et al. 2023)

# Reliability of Explanations (Ye and Durrett 2022)

- Examine the relationship between
  - accuracy of the answer
  - factuality of the explanation
  - consistency of the answer+explanation

- E-P: explain then predict (CoT)
- P-E: predict then explain

|  |  | Acc | Fac | Con | Acc=Fac | Acc=Con |
|---|---|---|---|---|---|---|
|  | *reliability of explanations generated by InstructGPT* | | | | | |
| InstructGPT | SYNTH (E-P) | 58.4 | 72.8 | 64.8 | 66.5 | 68.8 |
|  | SYNTH (P-E) | 54.8 | 51.6 | 95.2 | **89.6** | 57.2 |
|  | ADVHP (E-P) | 62.0 | 79.6 | 91.2 | **80.0** | 68.4 |
|  | ADVHP (P-E) | 54.0 | 69.2 | 82.0 | **77.6** | 67.2 |
|  | E-SNLI (P-E) | 62.0 | — | 98.8 | — | 62.0 |
|  | *reliability of explanations generated by other LLMs on* SYNTH | | | | | |
| OPT (175B) | SYNTH (E-P) | 30.0 | 77.2 | 47.2 | 45.6 | 58.8 |
|  | SYNTH (P-E) | 39.6 | 64.0 | 81.2 | **69.2** | 49.6 |
| GPT-3 | SYNTH (E-P) | 46.8 | 59.2 | 64.8 | **66.8** | 61.2 |
|  | SYNTH (P-E) | 52.4 | 52.4 | 83.2 | **78.4** | 58.0 |
| text-davinci-002 | SYNTH (E-P) | 86.0 | 91.6 | 85.2 | **91.2** | 84.8 |
|  | SYNTH (P-E) | 81.6 | 83.2 | 96.4 | **95.8** | 82.8 |

# **Complexity-based Prompting** (Fu et al. 2022)

- Interestingly, for some tasks a larger number of reasoning steps is indicative of improved accuracy



GSM8K — In-distribution: 43 (2, Simple), 52.5 (4), 52 (6), 58.5 (9, Complex)

MathQA — Noisy-labeled: 36 (2, Simple), 41 (5), 41.8 (10, Complex)

MultiArith — Transfered prompt: 82 (2, Simple), 83.5 (4), 93 (6), 86 (9, Complex)

- This can be used to simply improve accuracy:

  - Sample multiple reasoning paths

  - Perform self-consistency over the longer reasoning paths

# Training for Chain of Thought

# ORCA: Training Small Models for Reasoning (Mukherjee et al. 2024)

- Generates a large and diverse chain-of-thought dataset from GPT-3.5 and GPT-4

- 5M complex instructions + CoT explanations

- Corresponding high accuracy compared to other datasets



Evaluation with GPT-4

- Replicated in OpenORCA dataset

# Chain of Thought Reward Models (Lightman et al. 2023)

- Get human supervision on the steps

- Train reward model, give feedback on each sentence



44

# Prompting and Programs

# Structuring Outputs as Programs can Help (Madaan et al. 2022)

- When predicting a structured output, using a programming language instead of natural language often increases accuracy
- Why? Programs are highly-structured *and* included in pre-training data
- Asking the model to generate JSON can help formatting problems



(a) The script $\mathcal{G}$

(b) $\mathcal{G}$ converted to Python code $\mathcal{G}_c$ using our approach

(c) Straightforward encodings of the graph using the "DOT"

(d) Text format, or as a list of edges (node pairs)

# **Program-aided Language Models** (Gao et al. 2022)

- Using a program to generate outputs can be more precise than asking the LM to do so

- Especially useful for numeric questions

- See ChatGPT code interpreter, Bard code execution

- (More on agents/tools later)



**Chain-of-Thought (Wei et al., 2022)**

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 tennis balls. 2 cans of 3 tennis balls each is 6 tennis balls. 5 + 6 = 11. The answer is 11.

Q: The bakers at the Beverly Hills Bakery baked 200 loaves of bread on Monday morning. They sold 93 loaves in the morning and 39 loaves in the afternoon. A grocery store returned 6 unsold loaves. How many loaves of bread did they have left?

Model Output

A: The bakers started with 200 loaves. They sold 93 in the morning and 39 in the afternoon. So they sold 93 + 39 = 132 loaves. The grocery store returned 6 loaves. So they had 200 - 132 - 6 = 62 loaves left.
The answer is 62. ❌

**Program-aided Language models (this work)**

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 tennis balls.
tennis_balls = 5
2 cans of 3 tennis balls each is
bought_balls = 2 * 3
tennis balls. The answer is
answer = tennis_balls + bought_balls

Q: The bakers at the Beverly Hills Bakery baked 200 loaves of bread on Monday morning. They sold 93 loaves in the morning and 39 loaves in the afternoon. A grocery store returned 6 unsold loaves. How many loaves of bread did they have left?

Model Output

A: The bakers started with 200 loaves
loaves_baked = 200
They sold 93 in the morning and 39 in the afternoon
loaves_sold_morning = 93
loaves_sold_afternoon = 39
The grocery store returned 6 loaves.
loaves_returned = 6
The answer is
answer = loaves_baked - loaves_sold_morning
       - loaves_sold_afternoon + loaves_returned

>>> print(answer)
74 ✓

# Questions?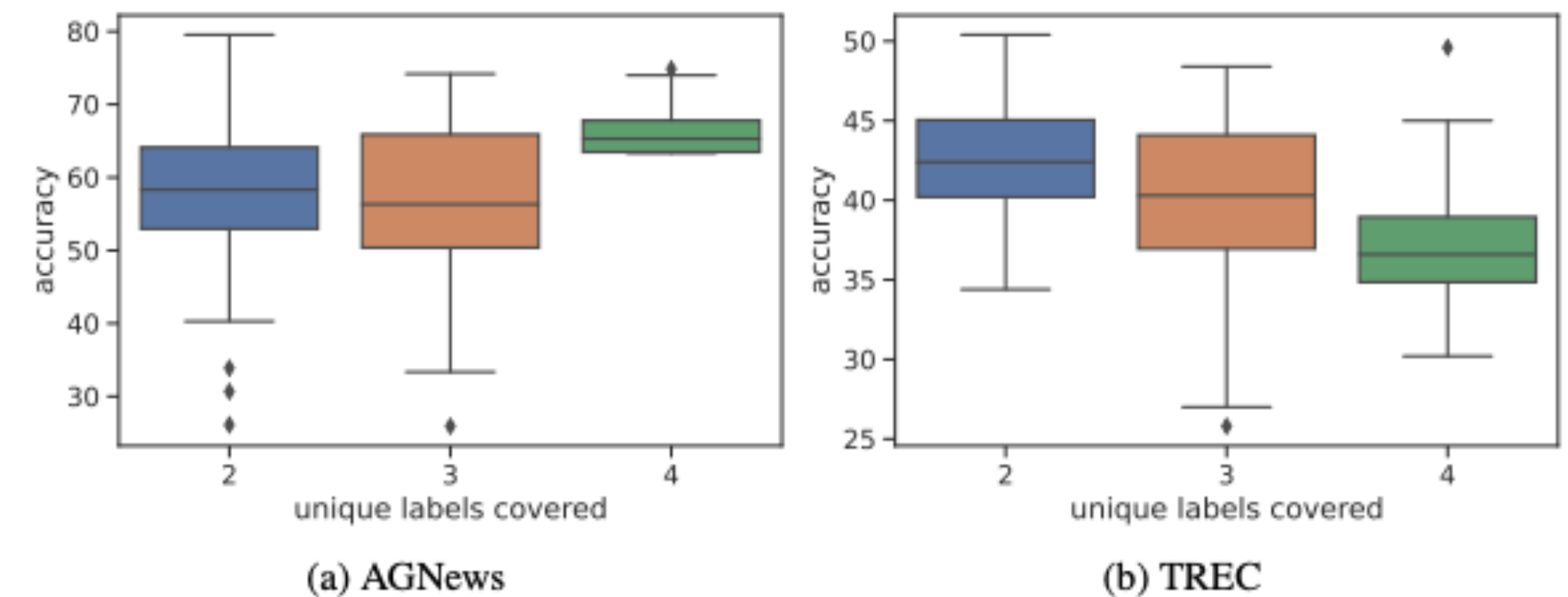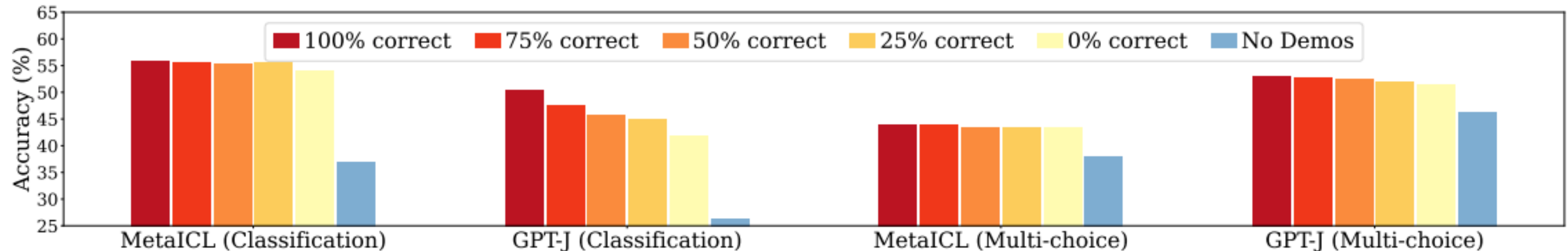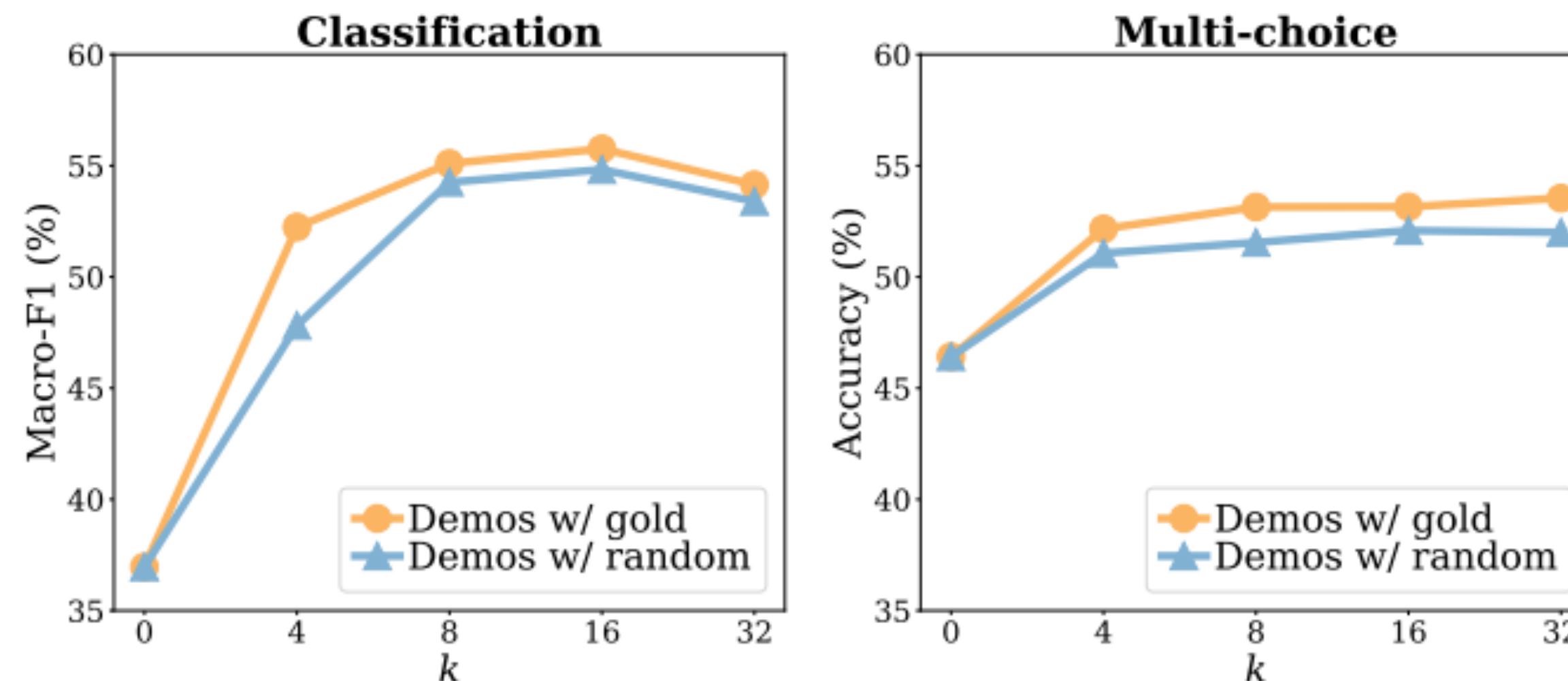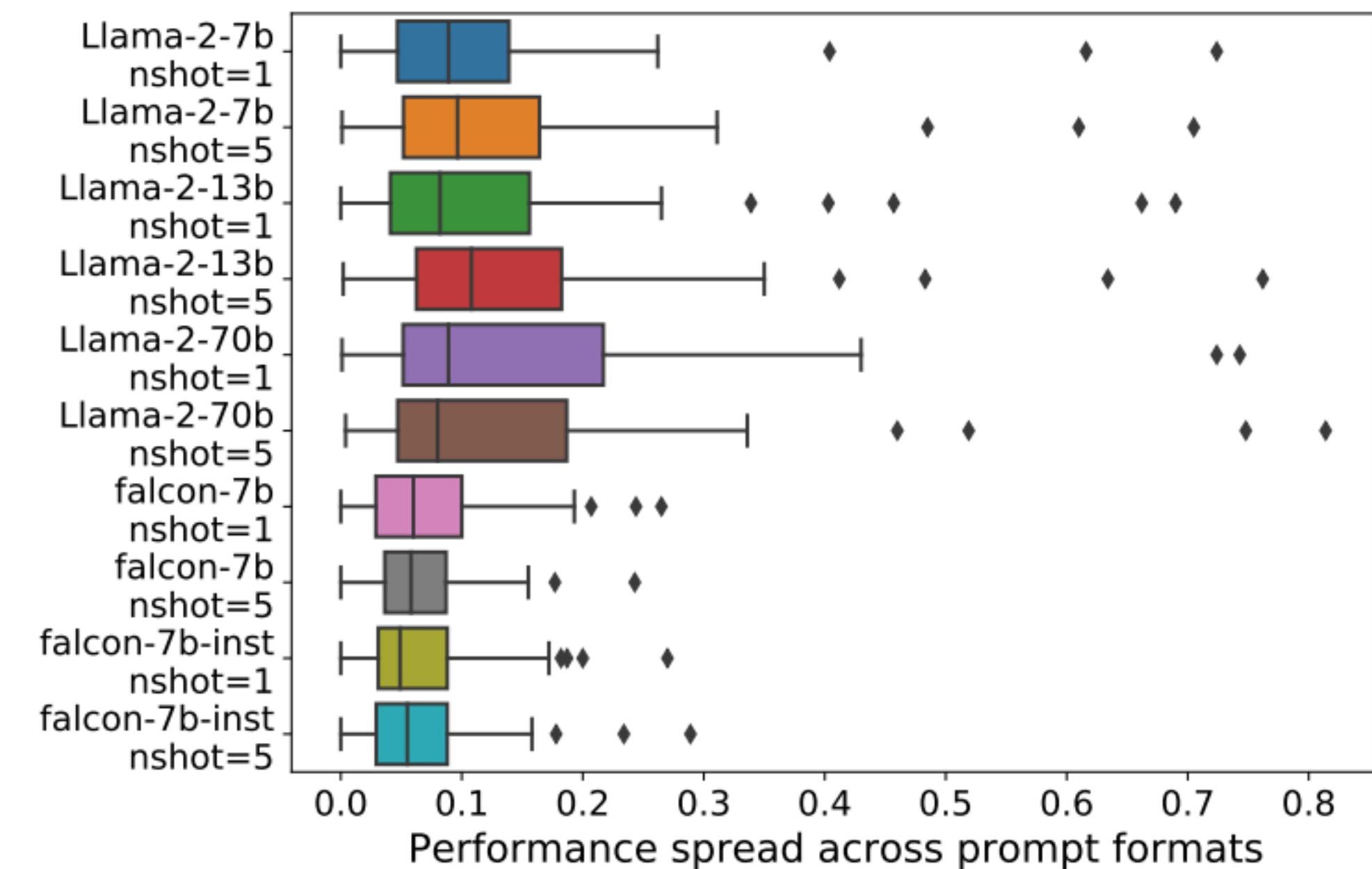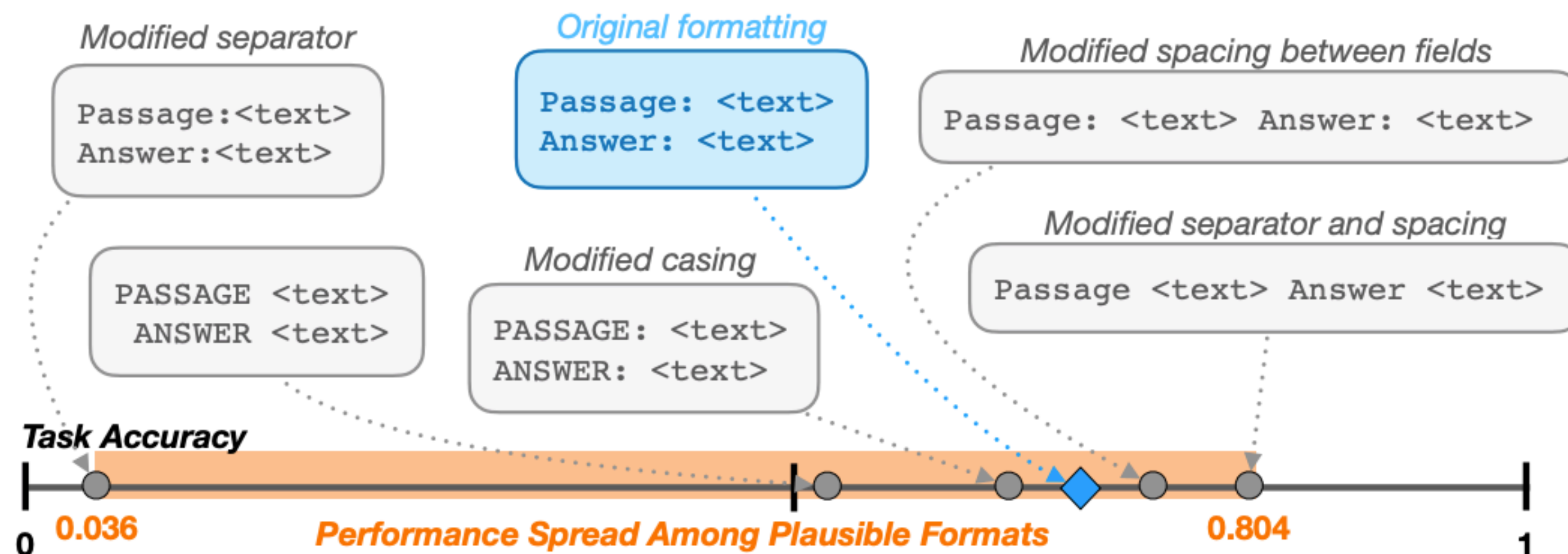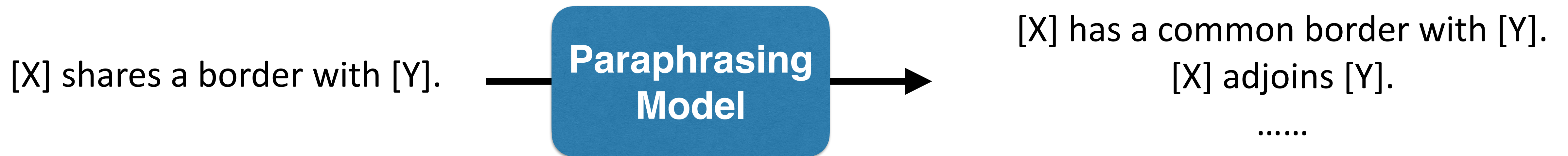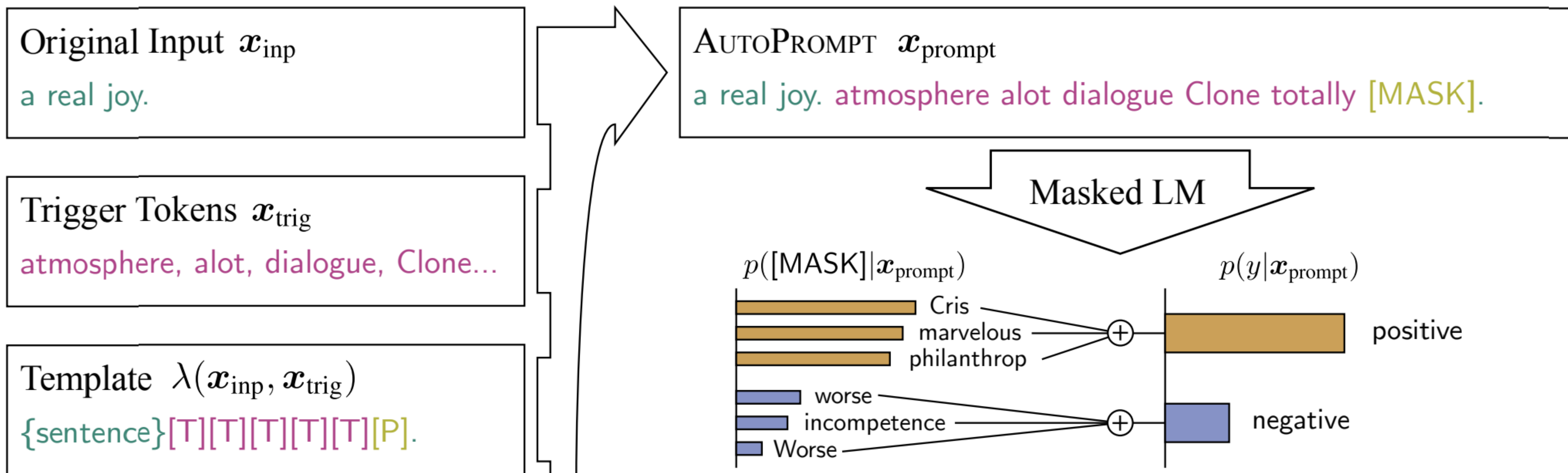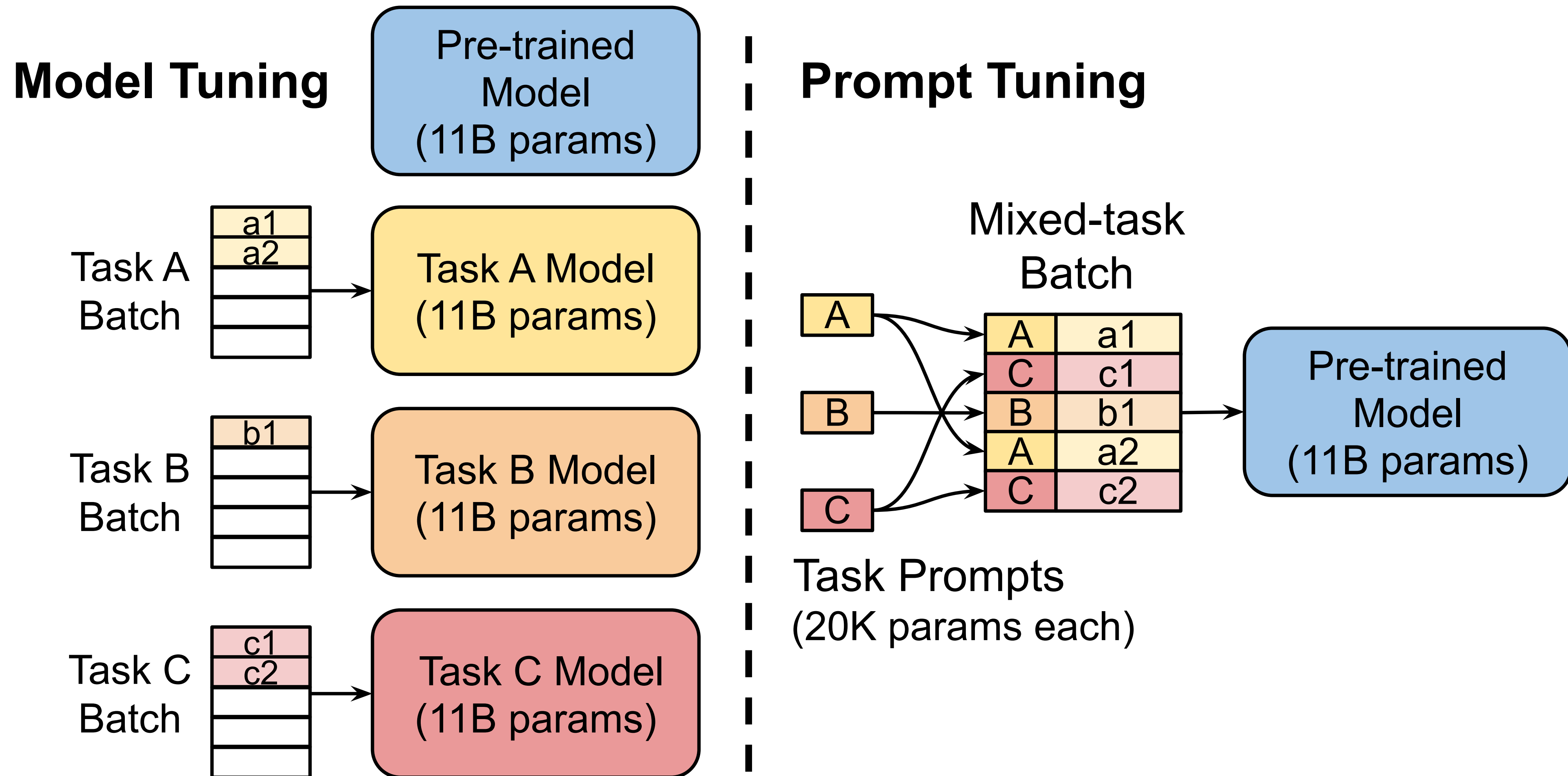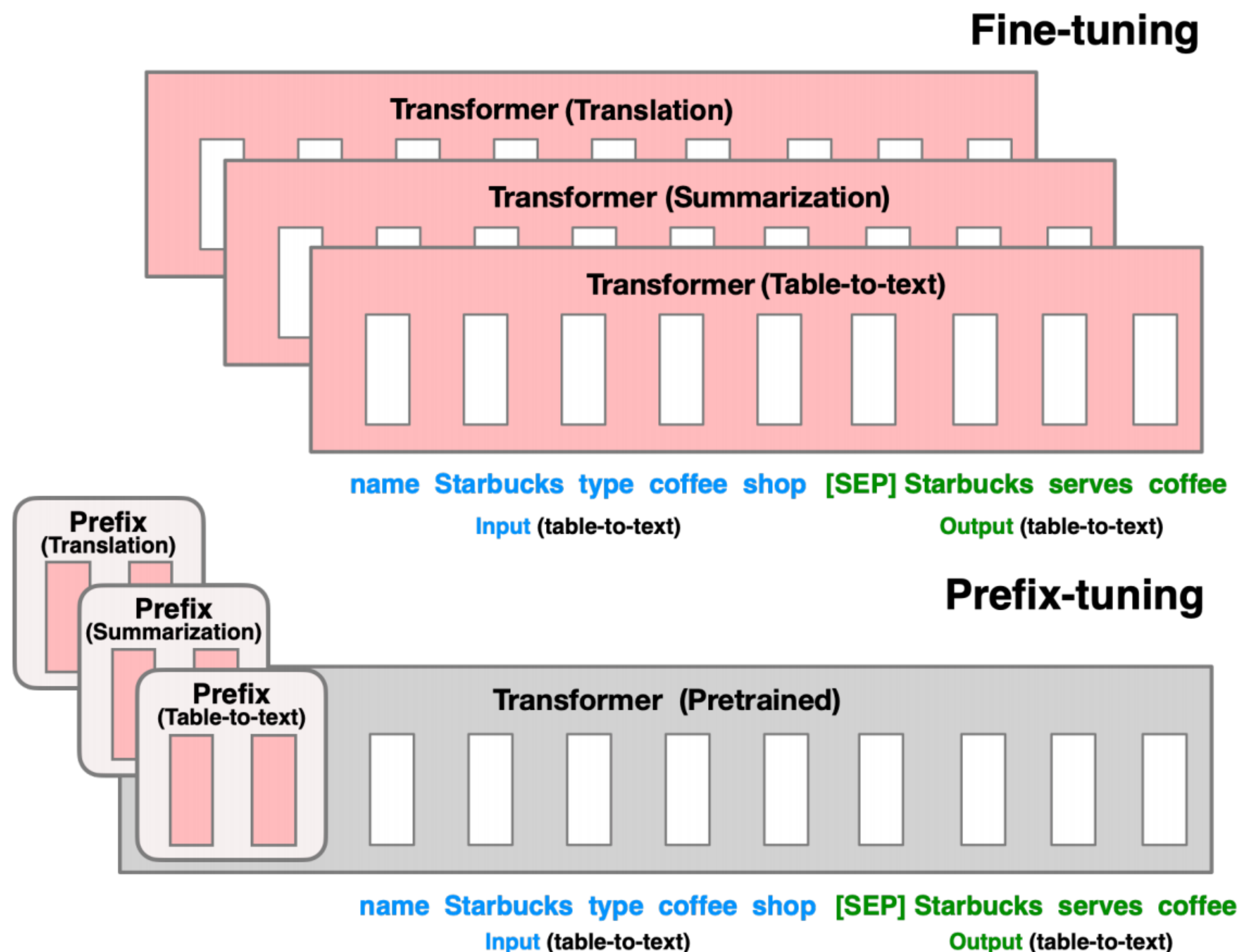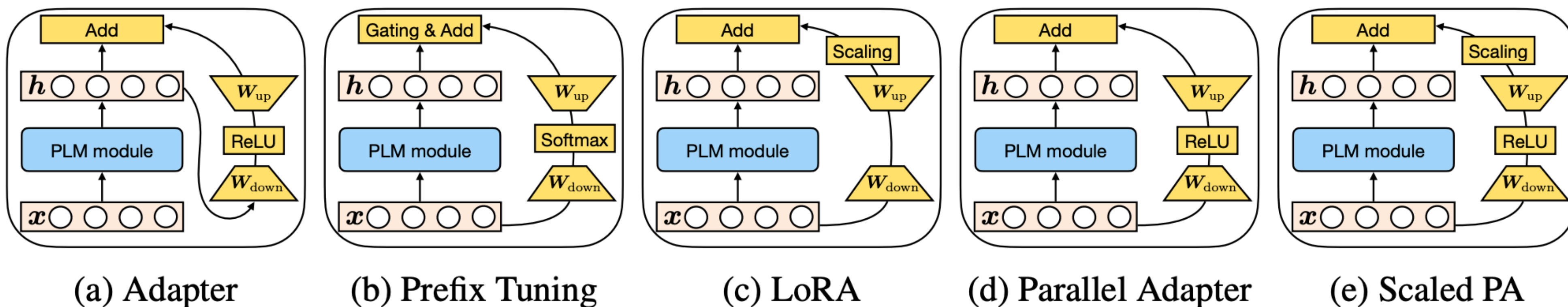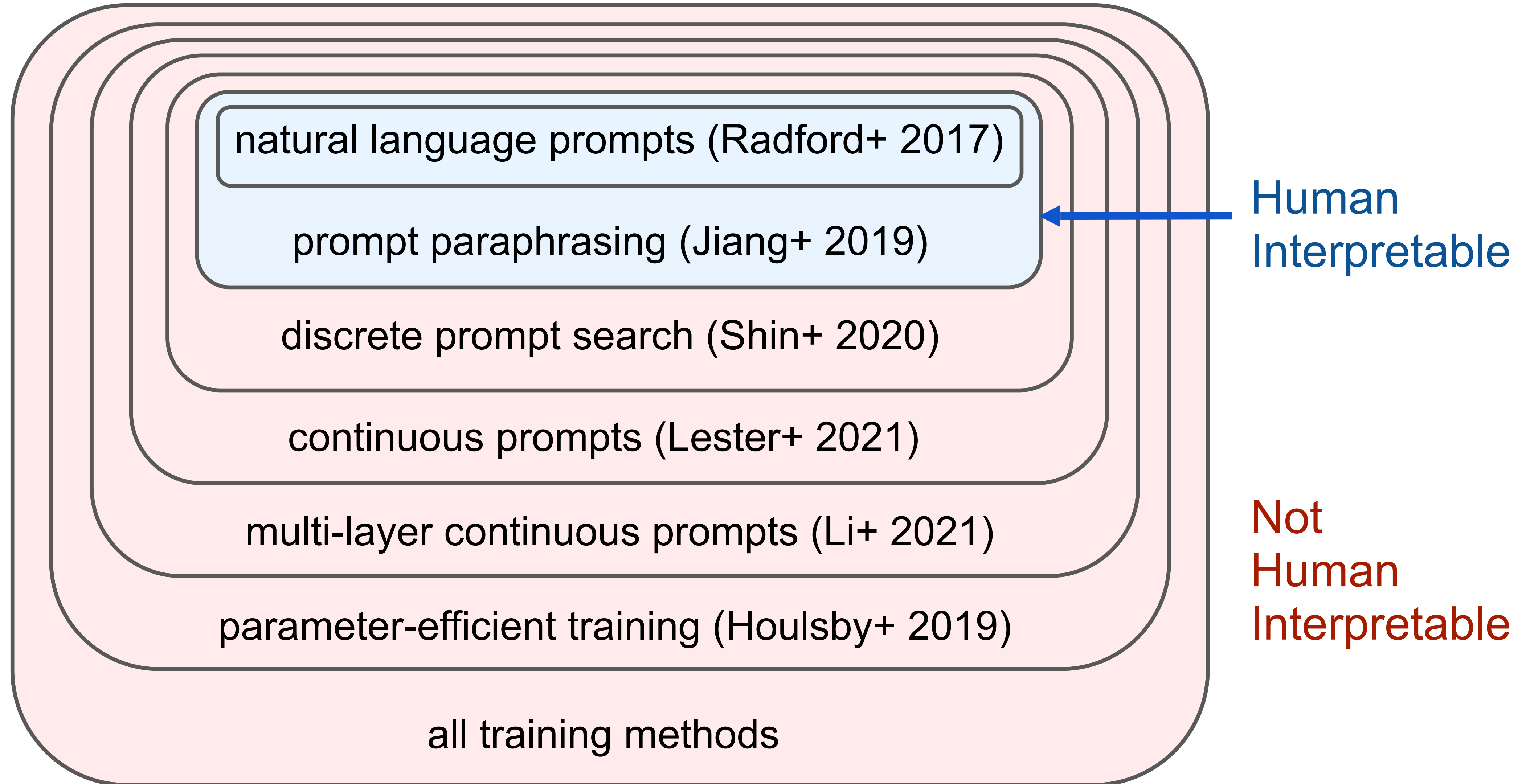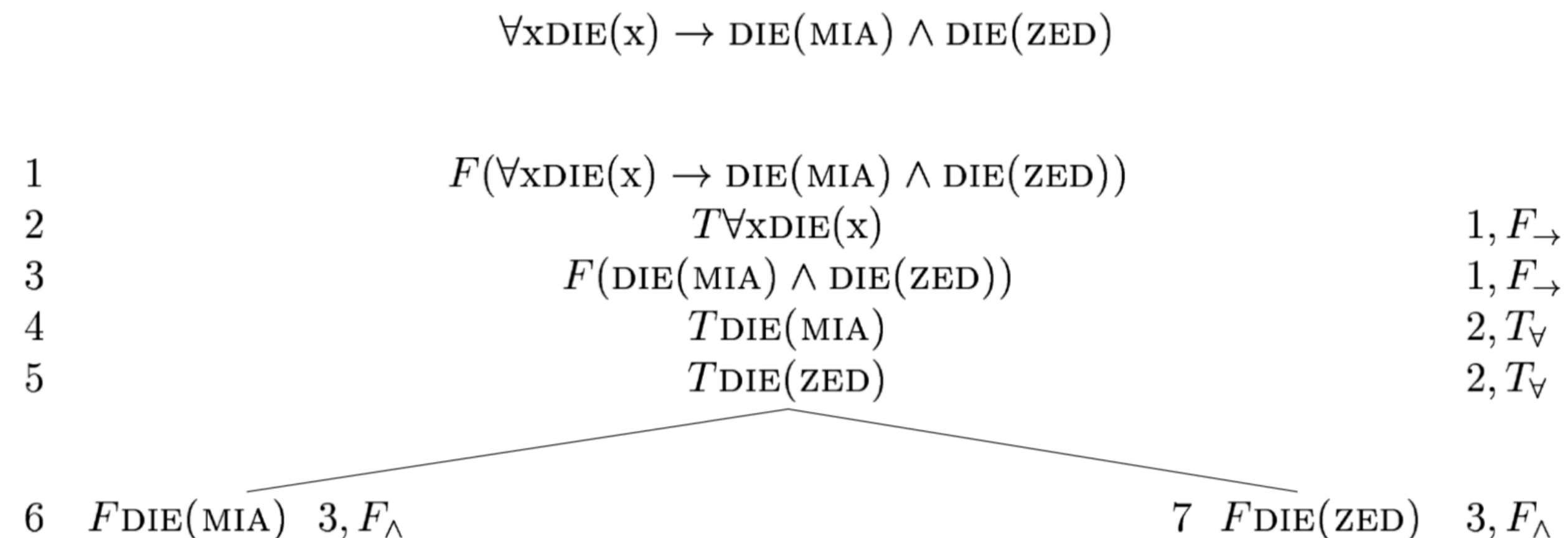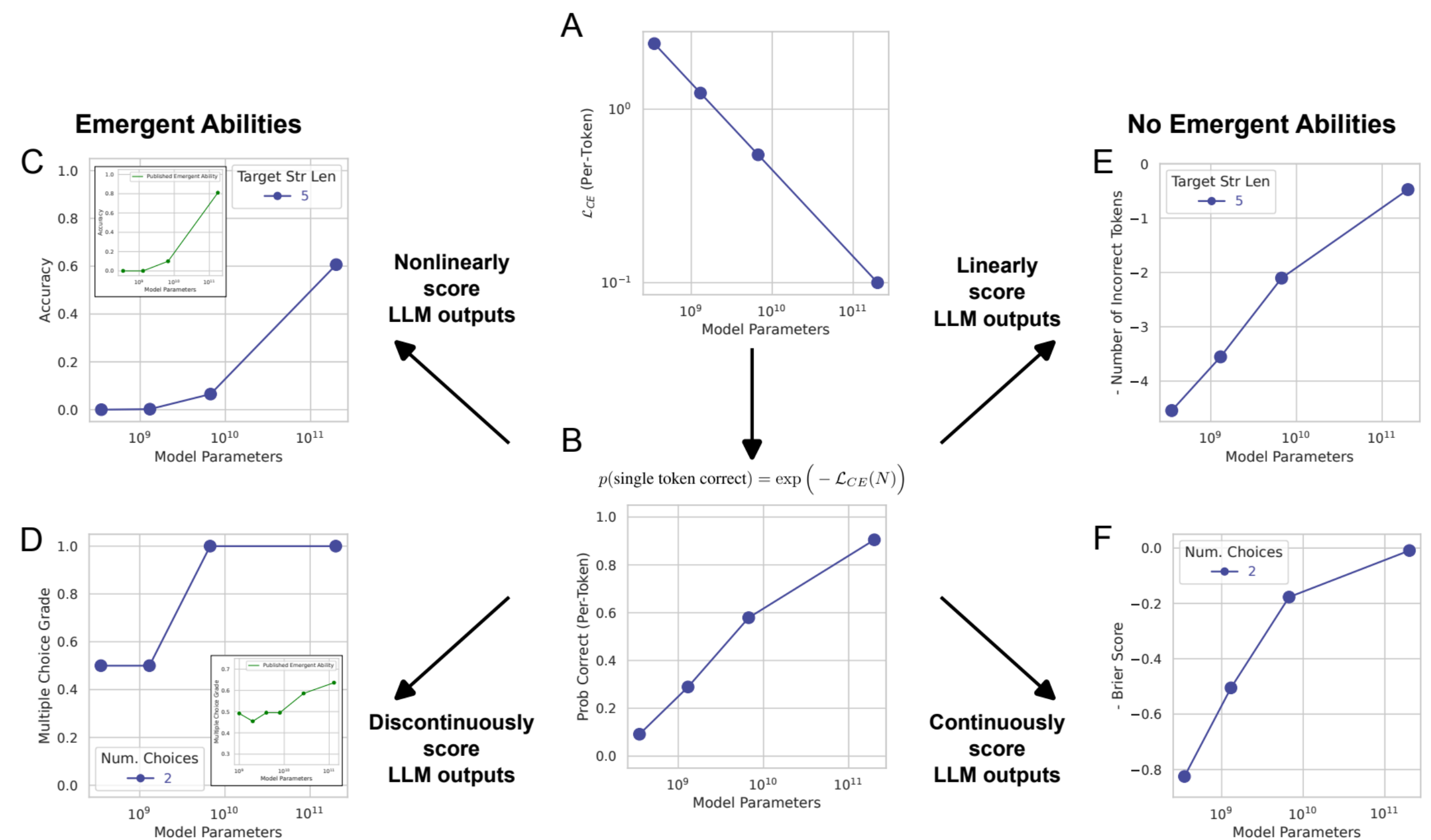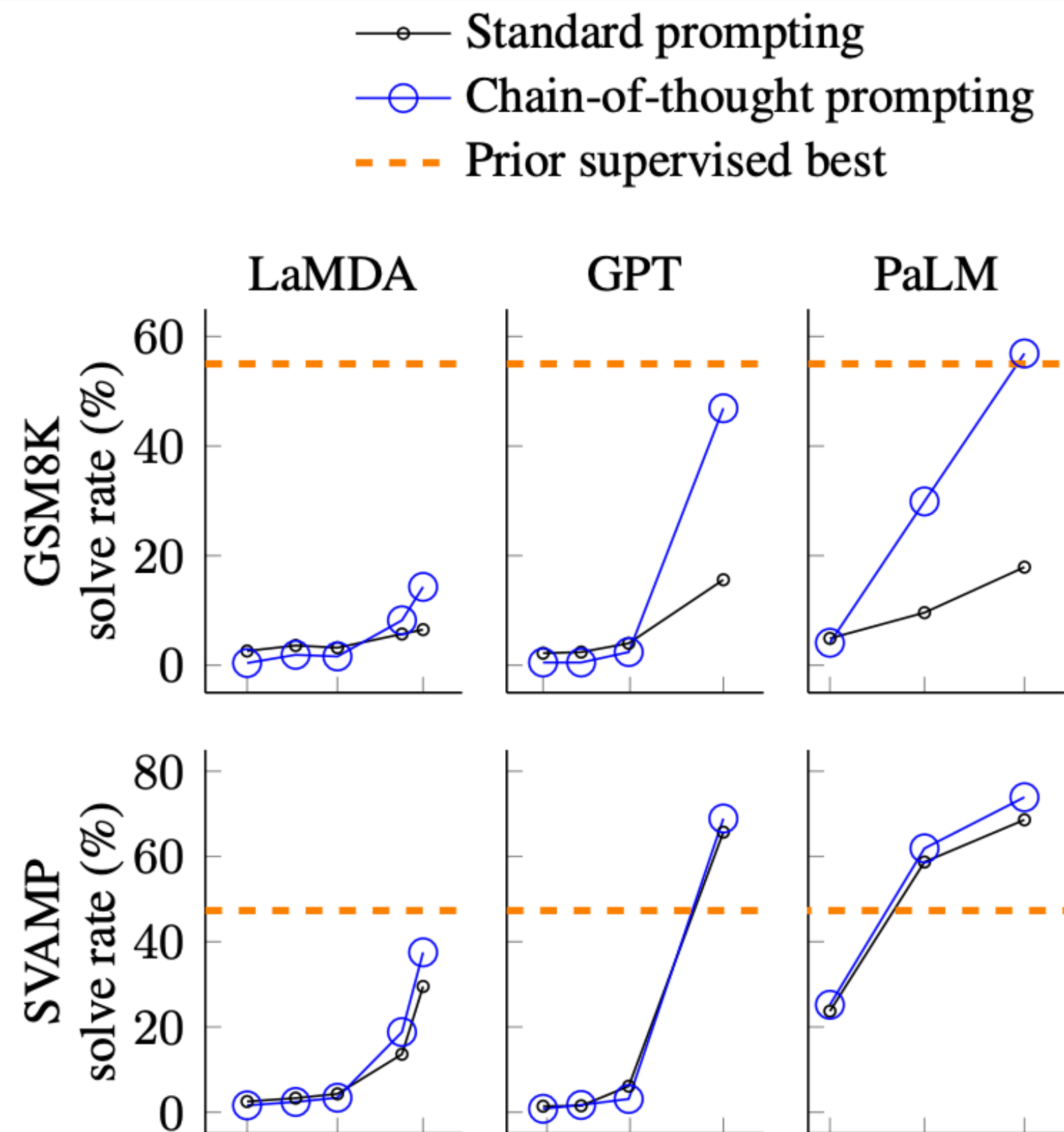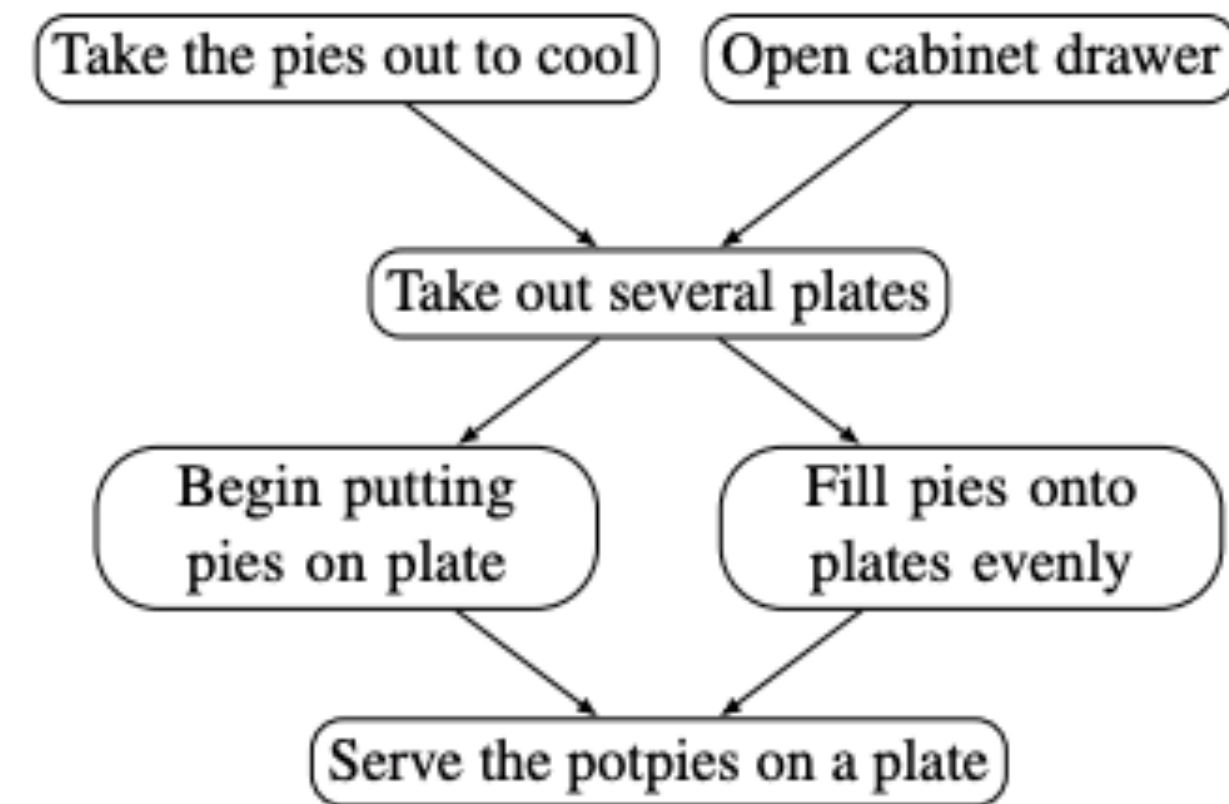