

# Making Generative AI Better for You:

Fine-tuning and Experimentation  
for Custom Research Solutions

**Shane Storks (he/him)**


PhD Candidate, Computer Science and Engineering  
Situated Language and Embodied Dialogue (SLED) Lab

*MIDAS Generative AI Tutorial Series*

November 29, 2023

# Large Language Models (LLMs)

LLMs like ChatGPT and GPT-4 have recently gained popularity due to their impressive language understanding and reasoning capabilities, making them useful assistants for a variety of language tasks.

 **You**  
When is the next tutorial session based on the below information?


Transforming Your Research with Generative AI – Tutorials  
MIDAS and the Michigan AI Laboratory will jointly offer a training series to researchers across research fields as they start incorporating Generative AI in their research. Each will consist of lectures, demonstrations, and hands-on tutorials for using Generative AI in research.

PAST EVENTS

10/2/23: Writing, Planning and Literature Review: Enhancing Professional Productivity with Generative AI

10/18/23: Code Smarter, Not Harder: Harnessing Generative AI for Research Programming Efficiency

10/27/23: Integrating Generative Image AI into Your Research Workflow

 **ChatGPT**  
The next tutorial session in the series "Transforming Your Research with Generative AI" is titled "Making Generative AI Better for You: Fine-tuning and Experimentation for Custom Research Solutions." This session is scheduled for Wednesday, November 29, 2023, from 2:30 PM to 4:30 PM. It will be held on the 10th Floor of Weiser Hall, located at 500 Church Street, Ann Arbor. The prerequisites for this session include some experience with machine learning.

...

***How can we customize them and apply them to empirical research?***

# Role of LLMs in Research

- LLMs can be helpful *assistants* for tasks like writing and coding, but they can do so much more!
- They can also be useful to automate aspects of:
  - Data annotation
  - Domain-specific content generation
  - Any language-based applications
- May not perform well at specialized tasks like these out of the box
- **How can we *customize* LLMs to adapt them to various specialized language tasks?**



[MIDAS Using Generative AI for Scientific Research User Guide](#)

# Outline

- The Road to LLMs
- Fine-Tuning LLMs
- Prompting LLMs

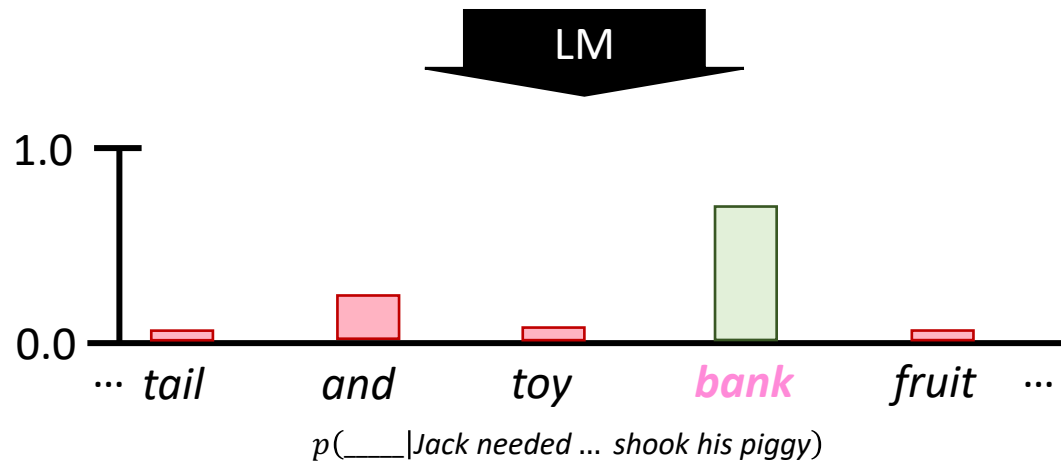
# Outline

- **The Road to LLMs**
- Fine-Tuning LLMs
- Prompting LLMs

# Language Models (LMs)

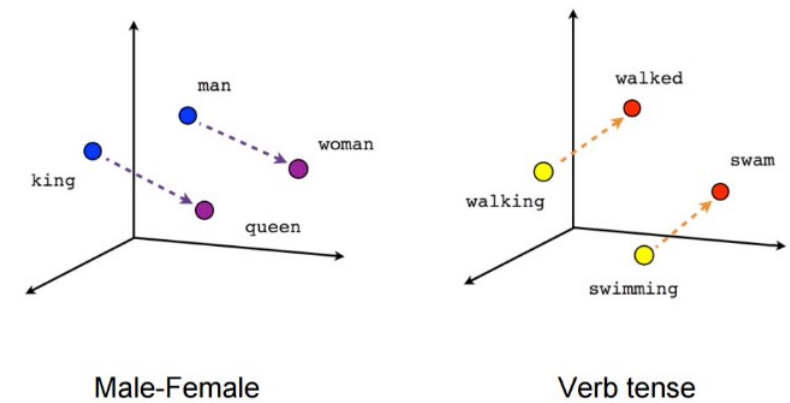
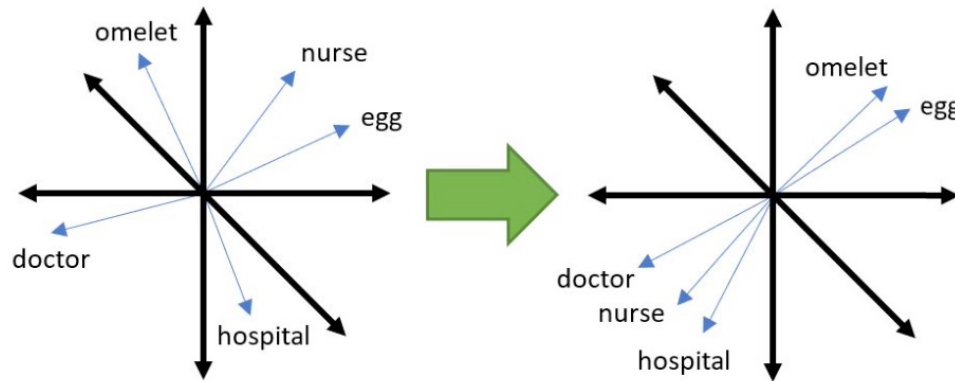
$$p(w_n | w_1, w_2, \dots, w_{n-1})$$

Jack needed some **money**, so he went and shook his **piggy** \_\_\_\_\_

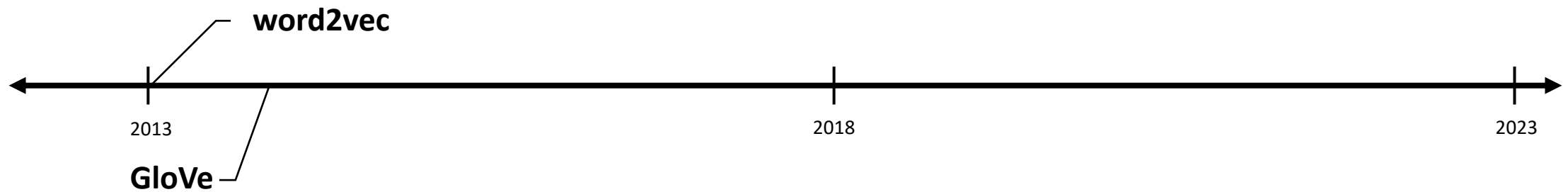


$$p(\text{Jack needed ...shook his piggy bank}) \approx p(\text{bank} | \text{Jack needed ...shook his piggy}) \cdot p(\text{piggy} | \text{Jack needed ...shook his}) \cdot p(\text{his} | \text{Jack needed ...shook}) \cdot \dots$$

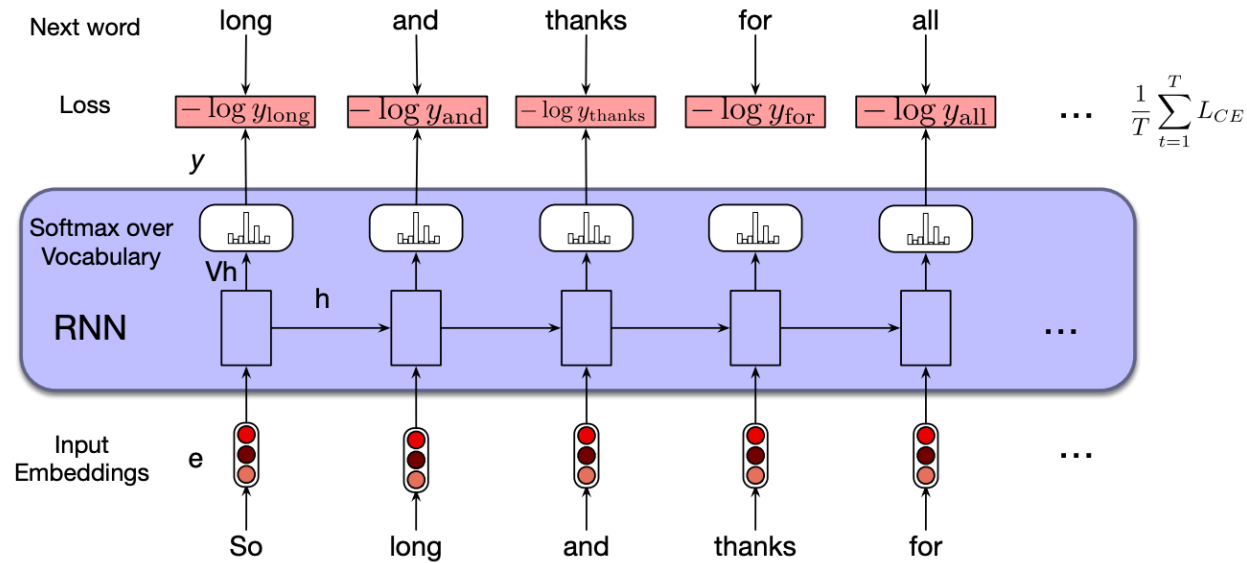
# Vector-Based Word Embeddings



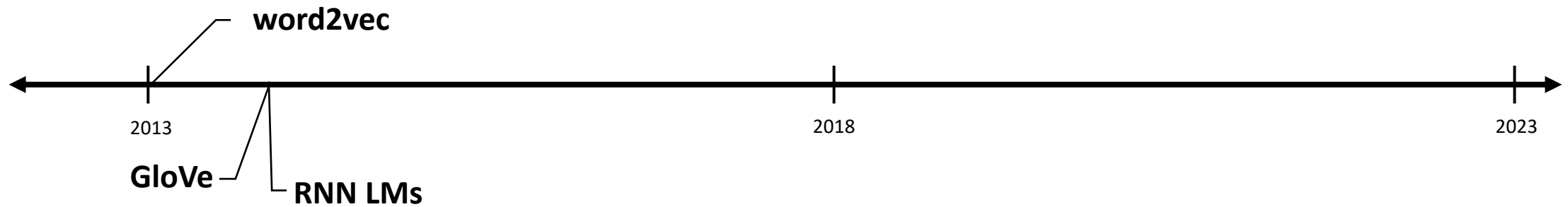
(Image from [TensorFlow docs](#))



# Representing Sequences of Words

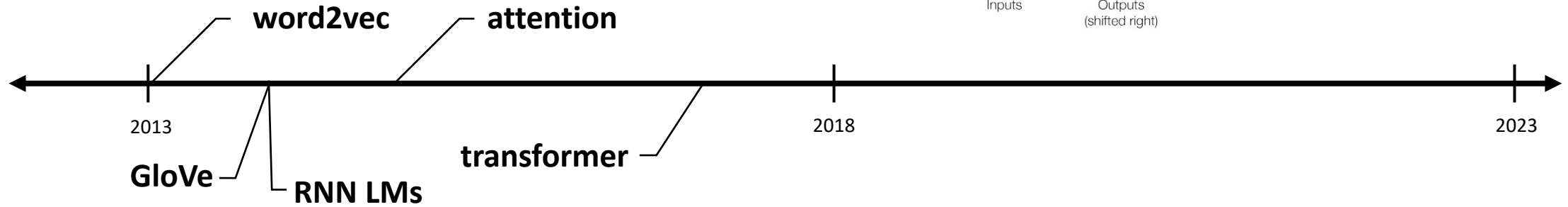
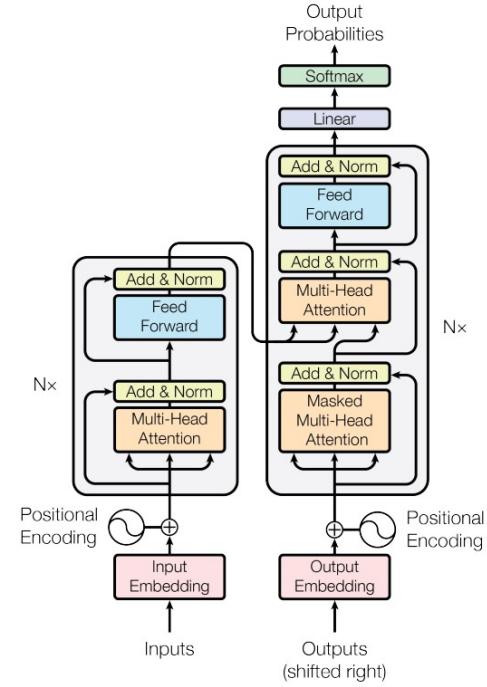
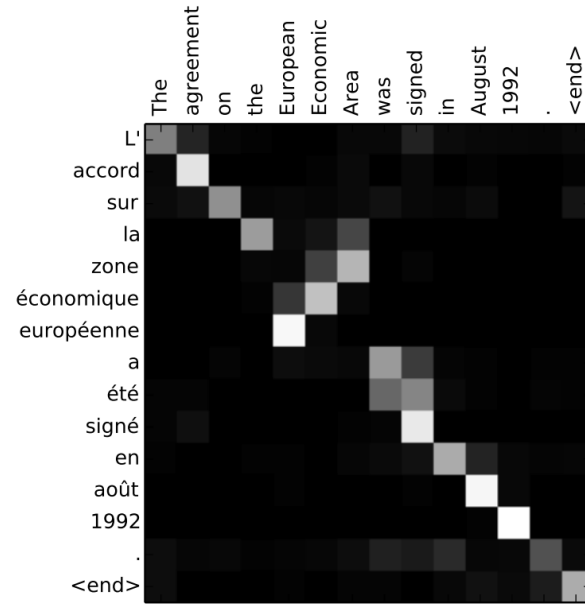


*(figure from Jurafsky and Martin)*



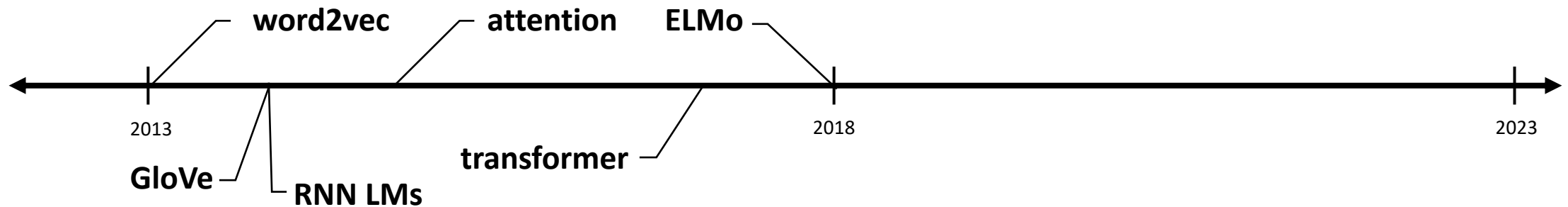
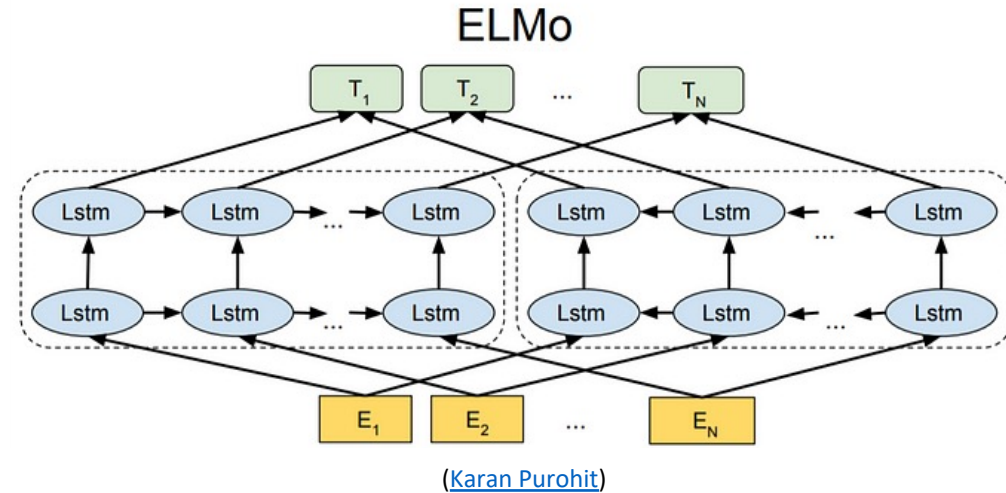


# Attention and Transformers

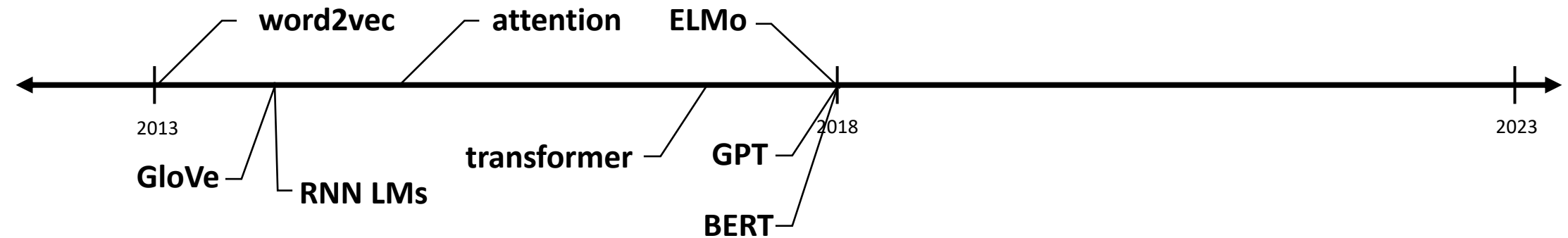
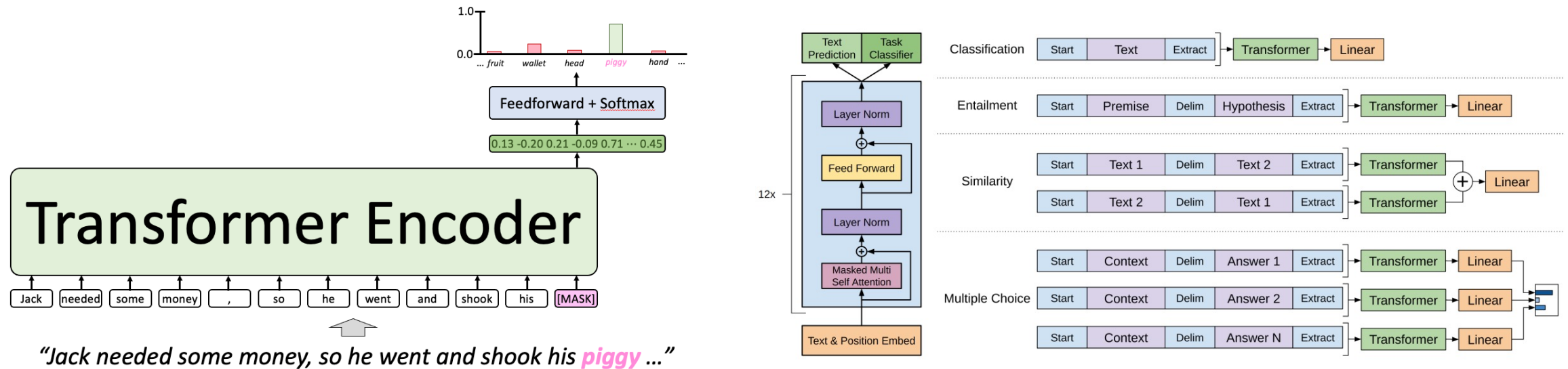


# Contextual Language Representations

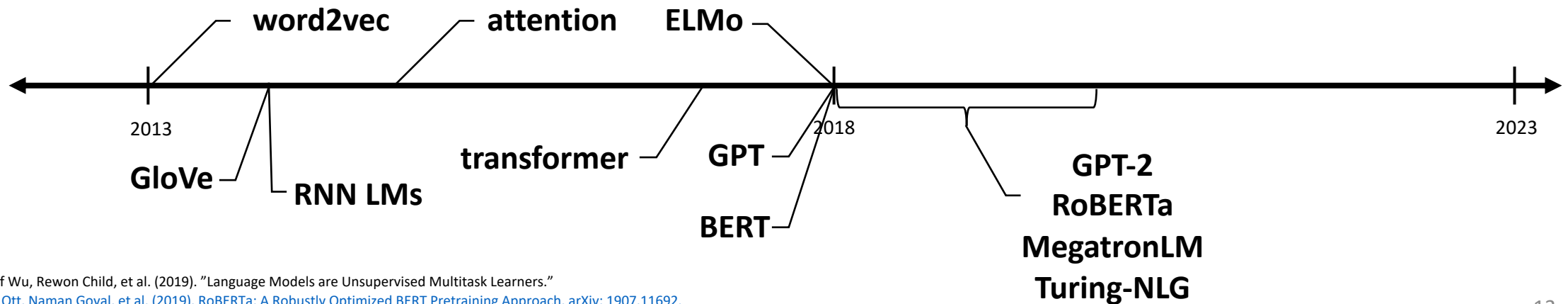
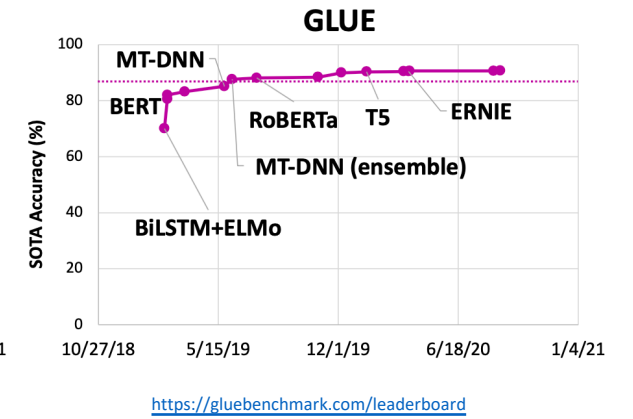
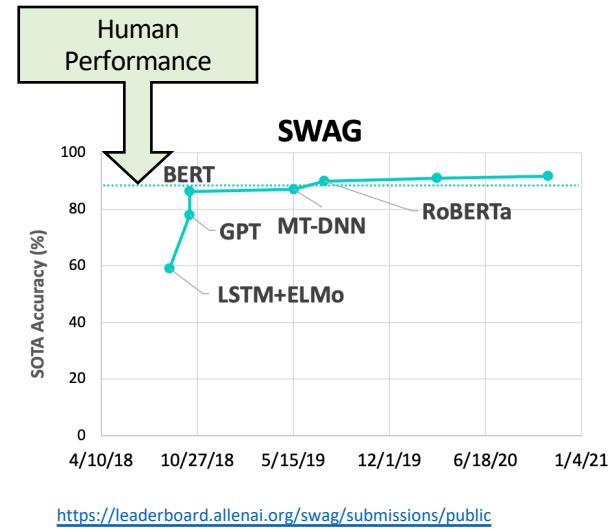
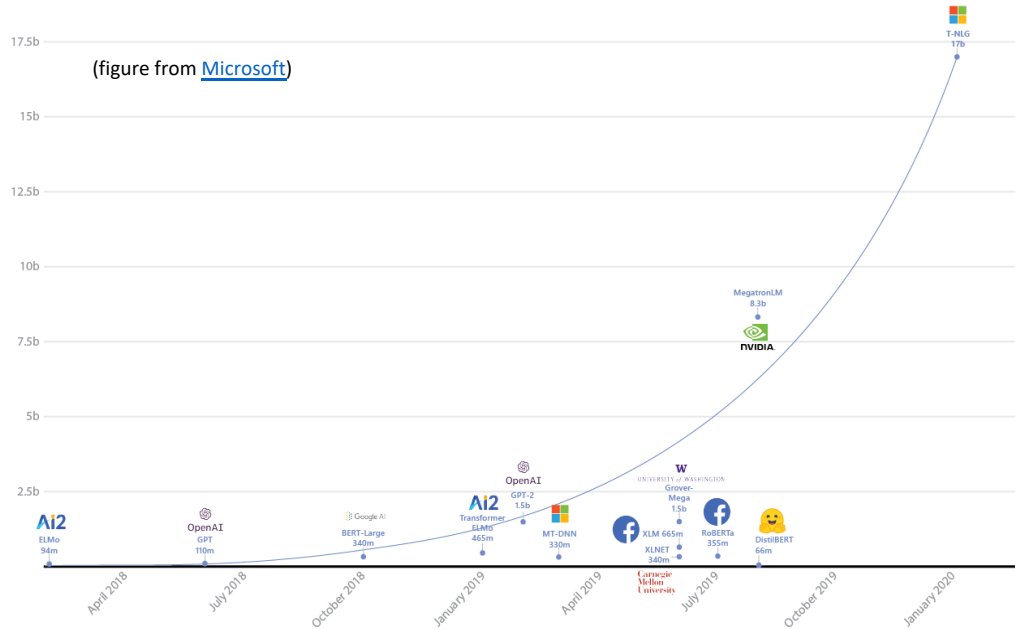
Source	Nearest Neighbors
GloVe play	playing, game, games, played, players, plays, player, Play, football, multiplayer
ELMo Chico Ruiz made a spectacular <u>play</u> on Alusik 's grounder {...}	Kieffer , the only junior in the group , was commended for his ability to hit in the clutch , as well as his all-round excellent play .
ELMo Olivia De Havilland signed to do a Broadway <u>play</u> for Garson {...}	{...} they were actors who had been handed fat roles in a successful play , and had talent enough to fill the roles competently , with nice understatement .



# Self-Supervision and Transfer Learning in LMs



# Bigger Data & Bigger Models -> LLMs



Alec Radford, Jeff Wu, Rewon Child, et al. (2019). "Language Models are Unsupervised Multitask Learners."

[Yinhan Liu, Myle Ott, Naman Goyal, et al. \(2019\). RoBERTa: A Robustly Optimized BERT Pretraining Approach. arXiv: 1907.11692.](#)

[Mohammad Shoeybi, Mostofa Patwary, Raul Puri, et al. \(2019\). Megatron-LM: Training Multi-Billion Parameter Language Models Using Model Parallelism. arXiv: 1909.08053.](#)

[Corby Rosset. \(2020\). Turing-NLG: A 17-billion-parameter language model by Microsoft.](#)

# Prompting & In-Context Learning

### Zero-shot

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.

```

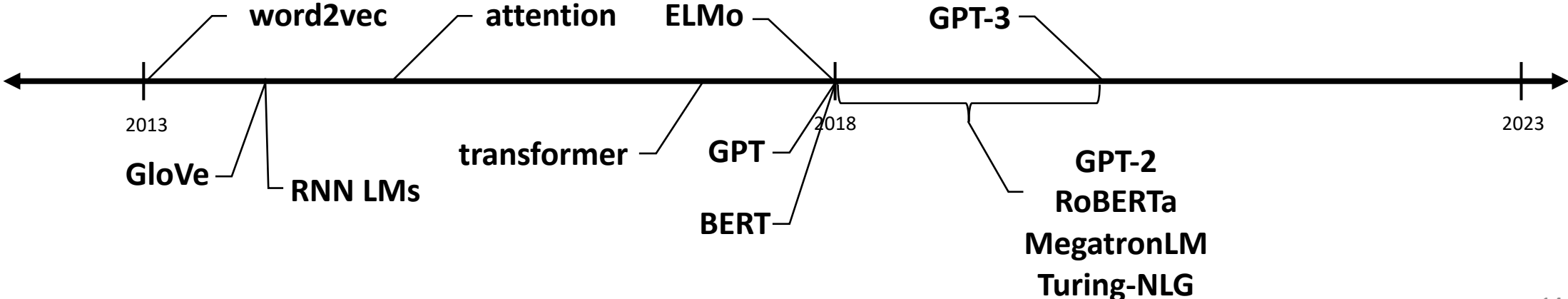
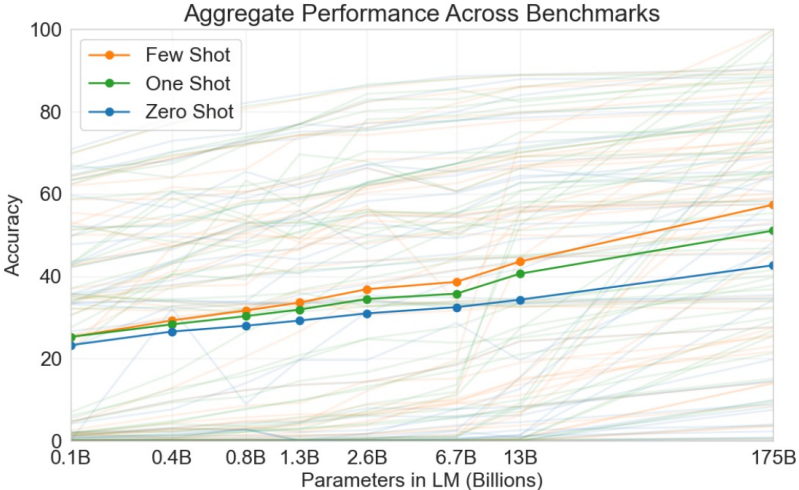
1 Translate English to French: ← task description
2 cheese => ..... ← prompt
    
```

### Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.

```

1 Translate English to French: ← task description
2 sea otter => loutre de mer ← examples
3 peppermint => menthe poivrée ← examples
4 plush girafe => girafe peluche ← examples
5 cheese => ..... ← prompt
    
```



# Instruction Tuning

## Finetune on many tasks ("instruction-tuning")

**Input (Commonsense Reasoning)**

Here is a goal: Get a cool sleep on summer days.  
How would you accomplish this goal?

OPTIONS:  
-Keep stack of pillow cases in fridge.  
-Keep stack of pillow cases in oven.

**Target**  
keep stack of pillow cases in fridge

**Input (Translation)**

Translate this sentence to Spanish:  
The new office building was built in less than three months.

**Target**  
El nuevo edificio de oficinas se construyó en tres meses.

Sentiment analysis tasks

Coreference resolution tasks

...

## Inference on unseen task type

**Input (Natural Language Inference)**

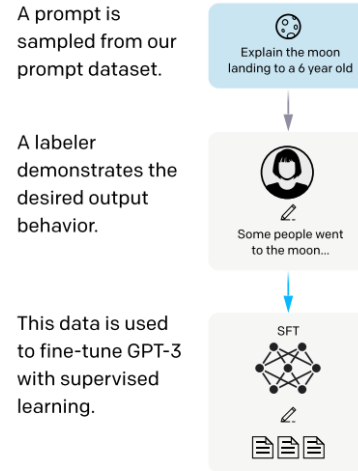
Premise: At my age you will probably have learnt one lesson.  
Hypothesis: It's not certain how many lessons you'll learn by your thirties.

Does the premise entail the hypothesis?

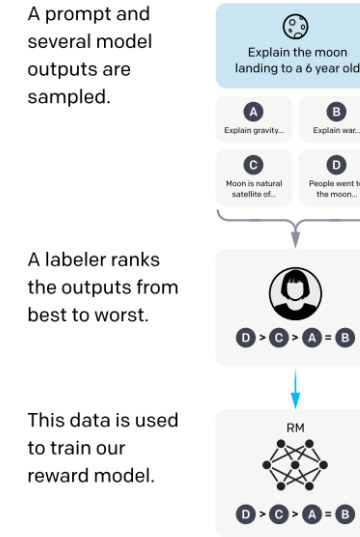
OPTIONS:  
-yes -it is not possible to tell -no

**FLAN Response**  
It is not possible to tell

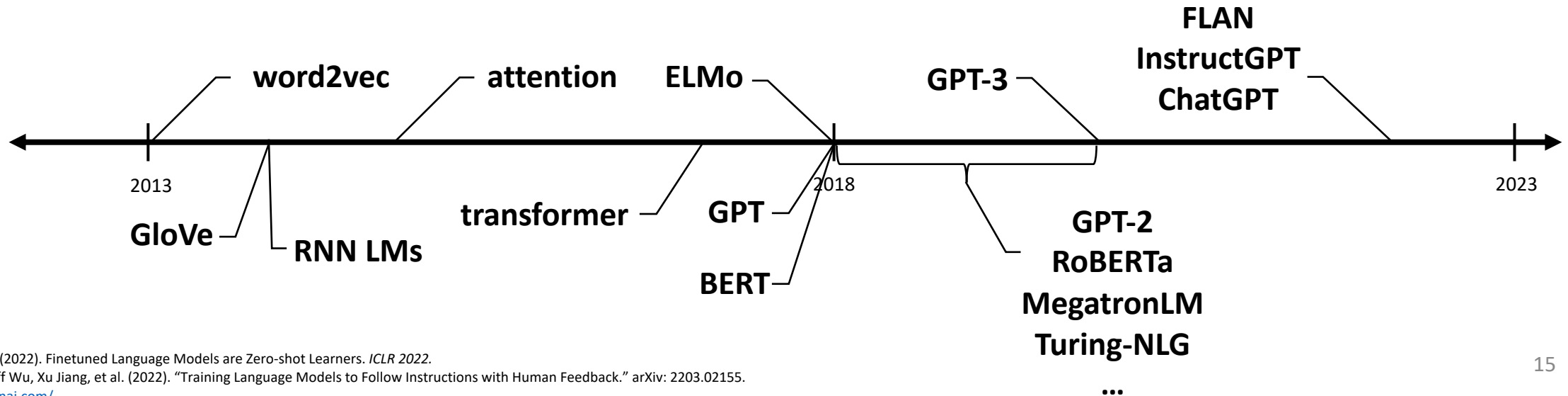
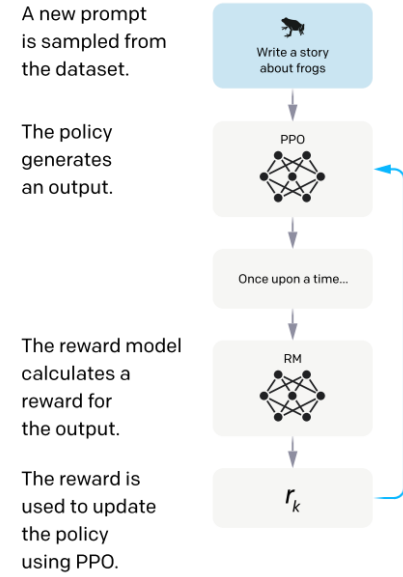
Step 1  
Collect demonstration data, and train a supervised policy.




Step 2  
Collect comparison data, and train a reward model.



Step 3  
Optimize a policy against the reward model using reinforcement learning.



# Vision & Multimodality



What is the common thing about these three images?

They are all flamingos.

What is the difference between these three images?

The first one is a cartoon, the second one is a real flamingo, and the third one is a 3D model of a flamingo.

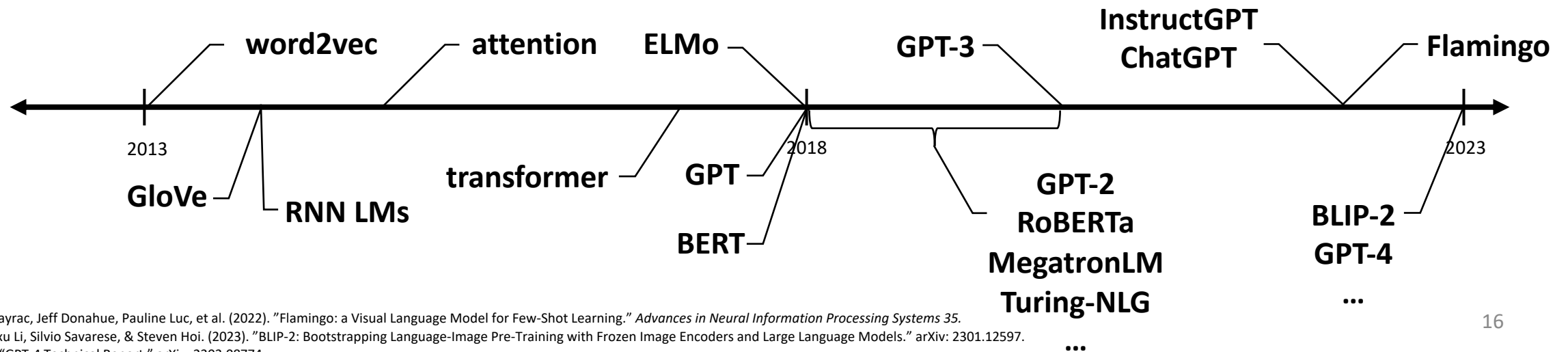
## GPT-4 visual input example, Extreme Ironing:

User What is unusual about this image?



Source: <https://www.barnorama.com/wp-content/uploads/2016/12/03-Confusing-Pictures.jpg>

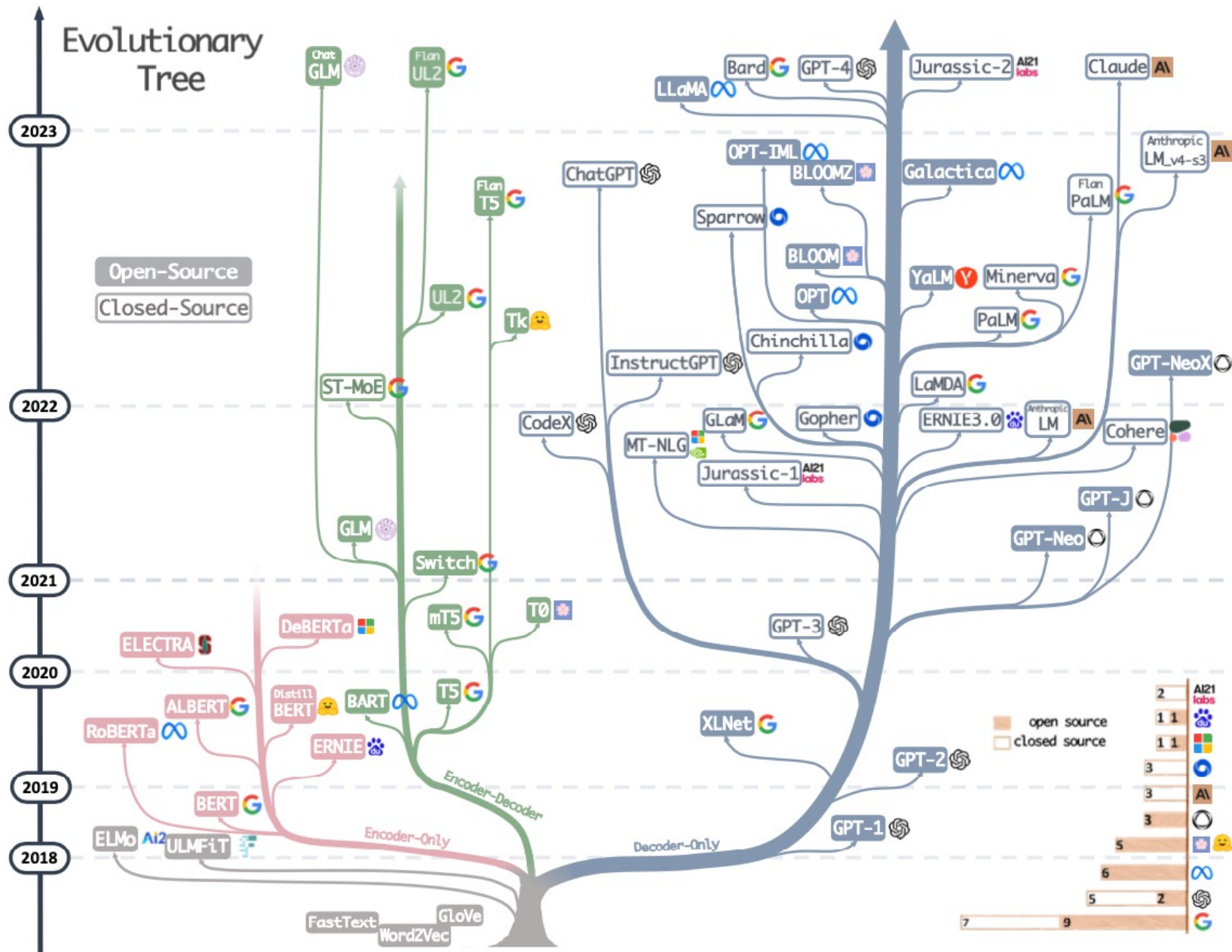
GPT-4 The unusual thing about this image is that a man is ironing clothes on an ironing board attached to the roof of a moving taxi.



Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, et al. (2022). "Flamingo: a Visual Language Model for Few-Shot Learning." *Advances in Neural Information Processing Systems* 35.

Junnan Li, Dongxu Li, Silvio Savarese, & Steven Hoi. (2023). "BLIP-2: Bootstrapping Language-Image Pre-Training with Frozen Image Encoders and Large Language Models." arXiv: 2301.12597.

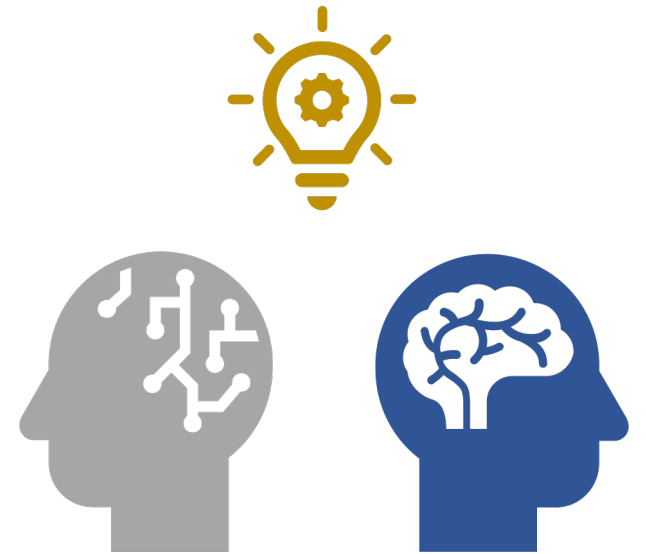
OpenAI. (2023). "GPT-4 Technical Report." arXiv: 2303.08774.





# Limitations of LLMs

- Despite these advancements and impressive capabilities, LLMs have some key limitations that cause undesirable behaviors
- In order to effectively and responsibly apply them in research, we need to be mindful of these limitations!



# Limitations of LLMs: Spurious Cues

Karen was assigned a roommate her first year of college. Her roommate asked her to go to a nearby city for a concert. Karen agreed happily. The show was absolutely exhilarating.

**How does the story end?**

Karen became good friends with her roommate.



Karen hated her roommate.



# Limitations of LLMs: Data Contamination

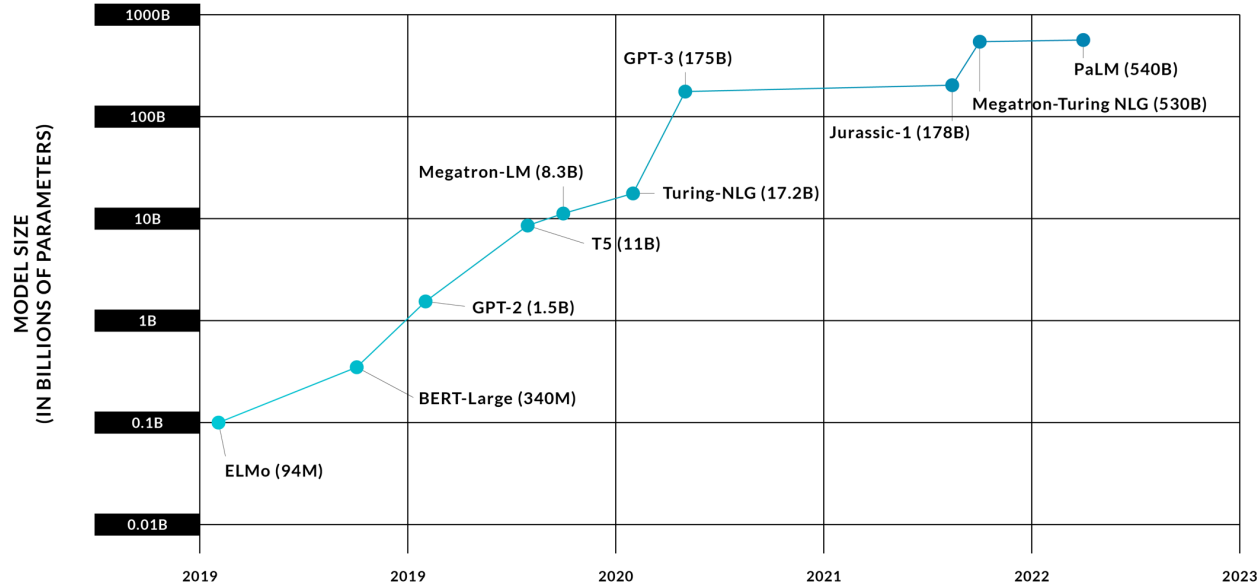
- LLMs have seen so much data in pre-training
- They may have been trained on benchmark datasets...
- Training on the test data is not an objective evaluation!

## **4 Measuring and Preventing Memorization Of Benchmarks**

Since our training dataset is sourced from the internet, it is possible that our model was trained on some of our benchmark test sets. Accurately detecting test contamination from internet-scale datasets is a new area of research without established best practices. While it is common practice to train large models without investigating contamination, given the increasing scale of pretraining datasets, we believe this issue is becoming increasingly important to attend to.

# Limitations of LLMs: Interpretability

Language Model Sizes Over Time



(figure from [Vinay Iyengar](#))

## GPT-4 Technical Report

OpenAI\*


...

### 2 Scope and Limitations of this Technical Report

This report focuses on the capabilities, limitations, and safety properties of GPT-4. GPT-4 is a Transformer-style model [39] pre-trained to predict the next token in a document, using both publicly available data (such as internet data) and data licensed from third-party providers. The model was then fine-tuned using Reinforcement Learning from Human Feedback (RLHF) [40]. Given both the competitive landscape and the safety implications of large-scale models like GPT-4, this report contains no further details about the architecture (including model size), hardware, training compute, dataset construction, training method, or similar.

# Limitations of LLMs: Hallucination

- **Hallucination:** generation of text that is factually incorrect, nonsensical, unfaithful to inputs, or otherwise incoherent

 **ChatGPT**

Dr. Joyce Chai is a well-known researcher in the field of artificial intelligence and natural language processing (NLP). As of my last knowledge update in January 2022, here's some information about her:

1. Background: Joyce Chai is a professor in the Department of Electrical Engineering and Computer Science at the University of Michigan. She has made significant contributions to the areas of natural language processing, machine learning, and human-robot interaction.
2. Research Focus: Her research interests include dialogue systems, sentiment analysis, emotion recognition, and the development of AI systems that can understand and generate human language. Her work is particularly relevant to the development of conversational AI systems and chatbots.
3. Contributions: Dr. Chai has published numerous research papers and articles in top conferences and journals in the field of NLP and AI. Her research has helped advance our understanding of how AI systems can better understand and respond to human language in natural and meaningful ways.
4. Education: She holds a Ph.D. in Computer Science from the University of Washington.

## Lawyer cites fake cases generated by ChatGPT in legal brief

The high-profile incident in a federal case highlights the need for lawyers to verify the legal insights generated by AI-powered tools.

Published May 30, 2023



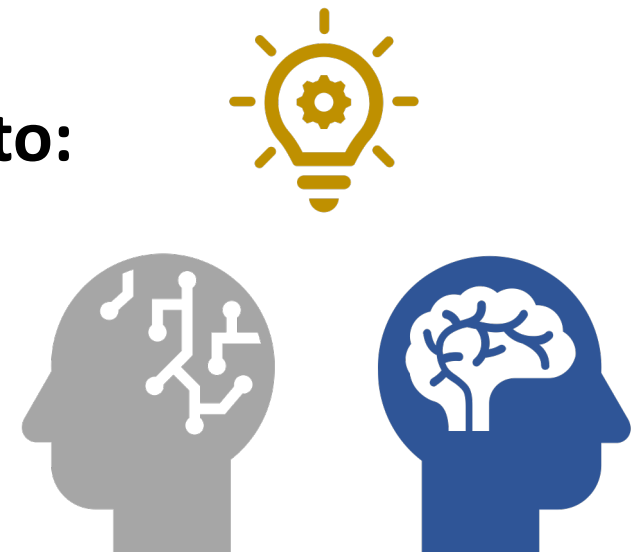
Lyle Moran  
Reporter



[Legal Dive](#)

# Summary

- LLMs' are remarkably useful for many language tasks, but these limitations make them impossible to trust consistently
- Verifying LLM outputs is important:
  - Automated metrics
  - Human evaluation
- **We must be mindful that LLMs are primarily trained to:**
  - **Generate fluent-sounding language (pre-training)**
  - **Satisfy users' requests (instruction-tuning)**



# 2 Ways to Customize LLMs

## Fine-Tuning:

Small hardware requirements

Host locally (private, more flexible)

Optimized for specific task

Technical skills, engineering effort

Large amount of training data

Hard to adapt once trained

## Prompting:

Larger hardware requirements

Best LMs behind proprietary APIs

Requires prompt engineering

User-friendly language interface

No training data needed

Generalizable and adaptable

# Outline

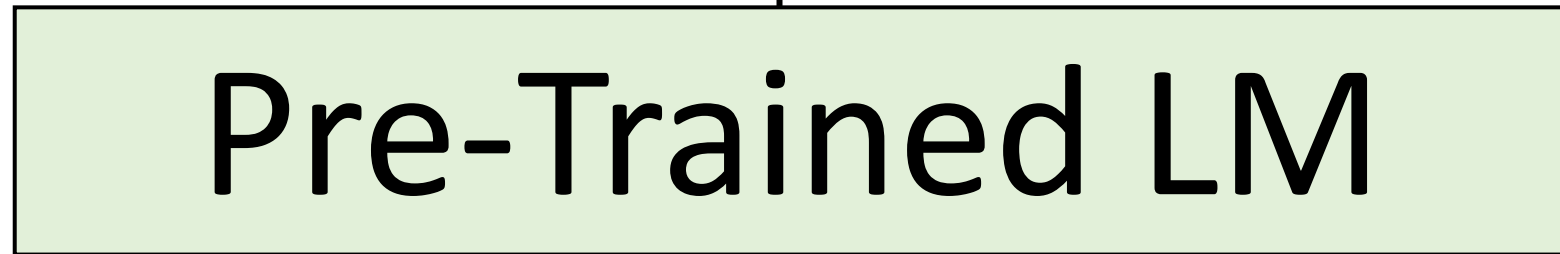
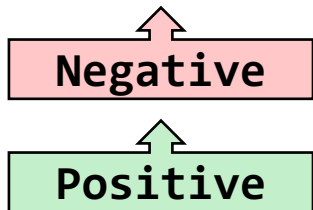
- The Road to LLMs
- **Fine-Tuning LLMs**
- Prompting LLMs



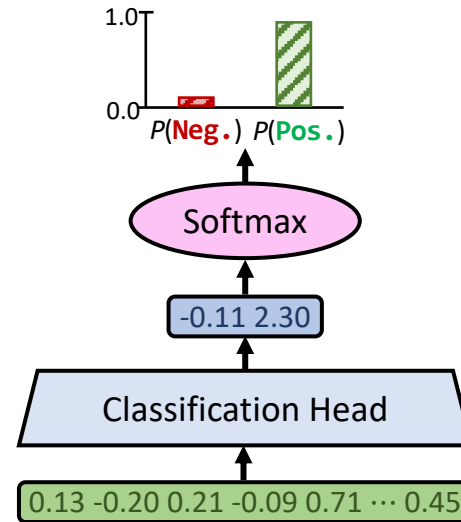
# Fine-Tuning: Text Classification

What is the sentiment of this text?

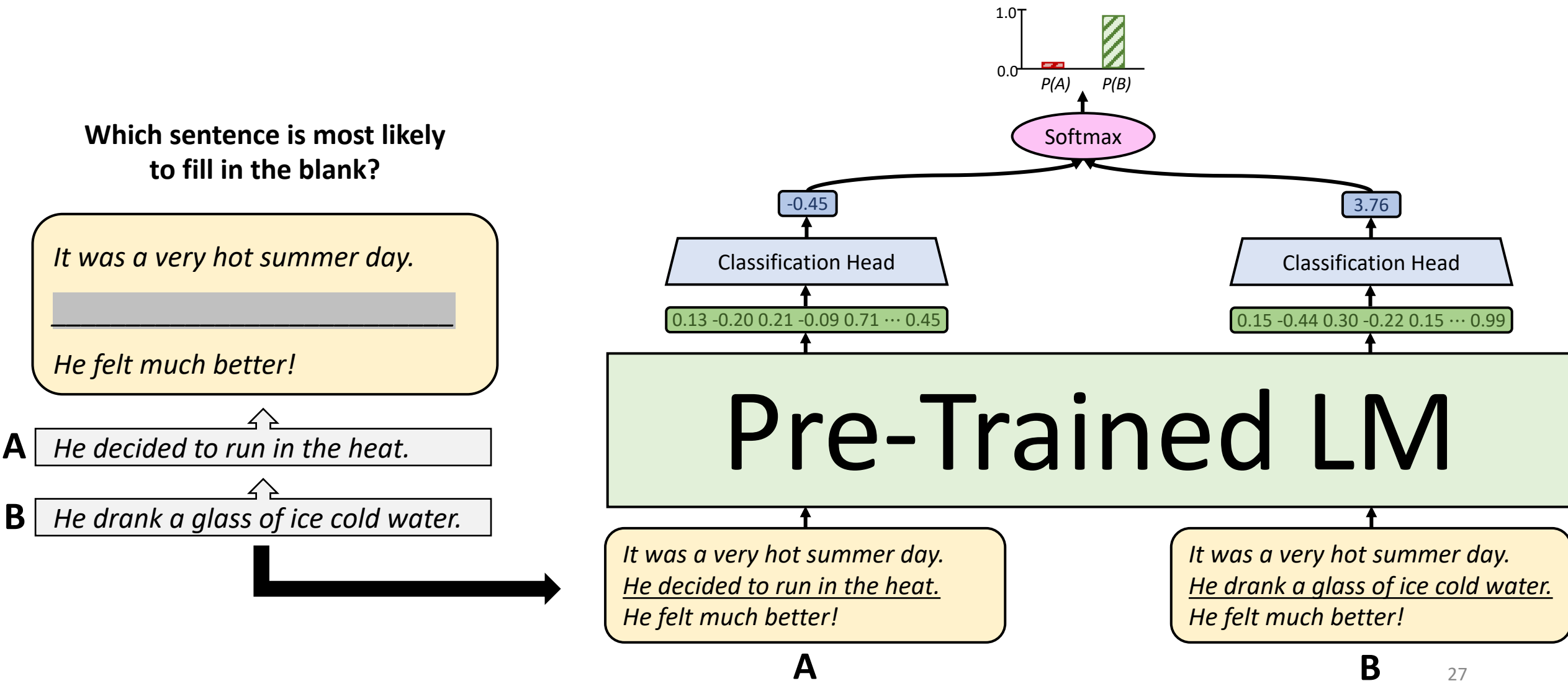
*The film was a charming and affecting journey.*



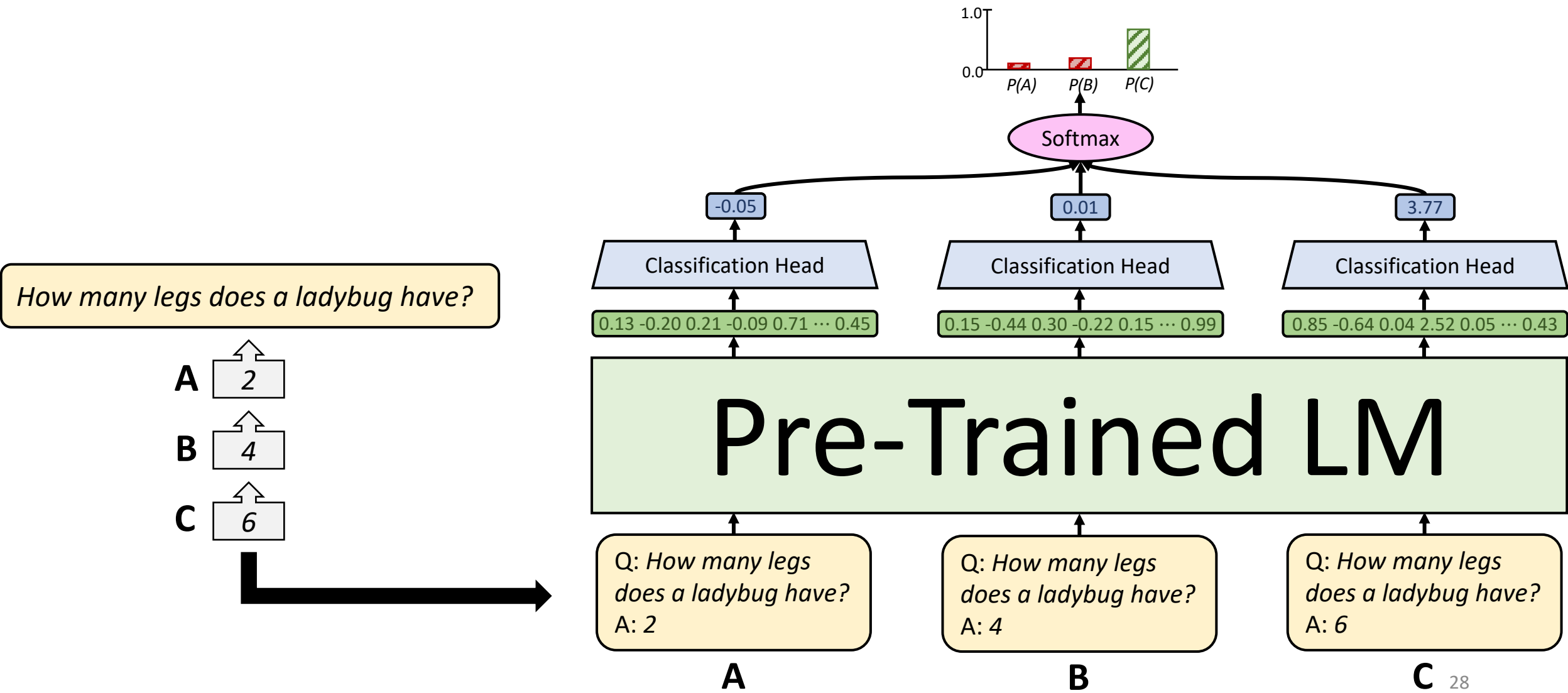
*The film was a charming and affecting journey.*



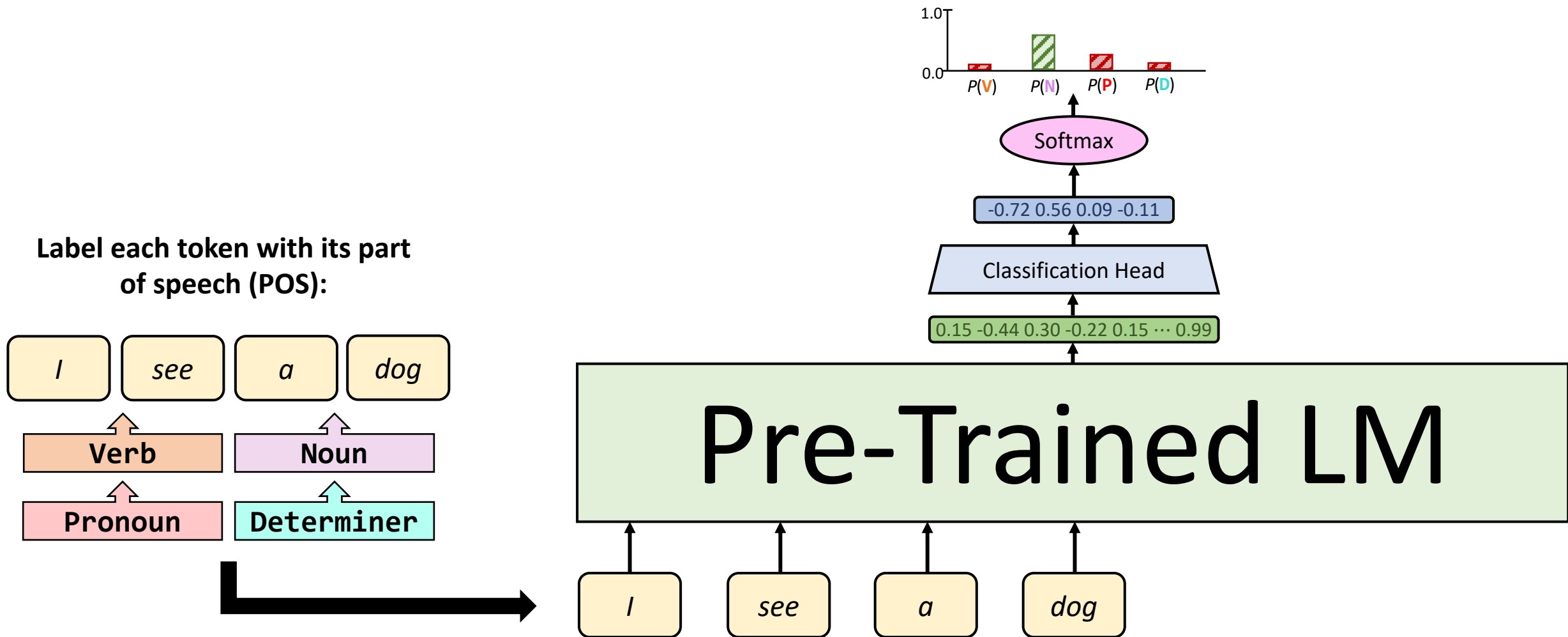
# Fine-Tuning: Multiple Choice Completion



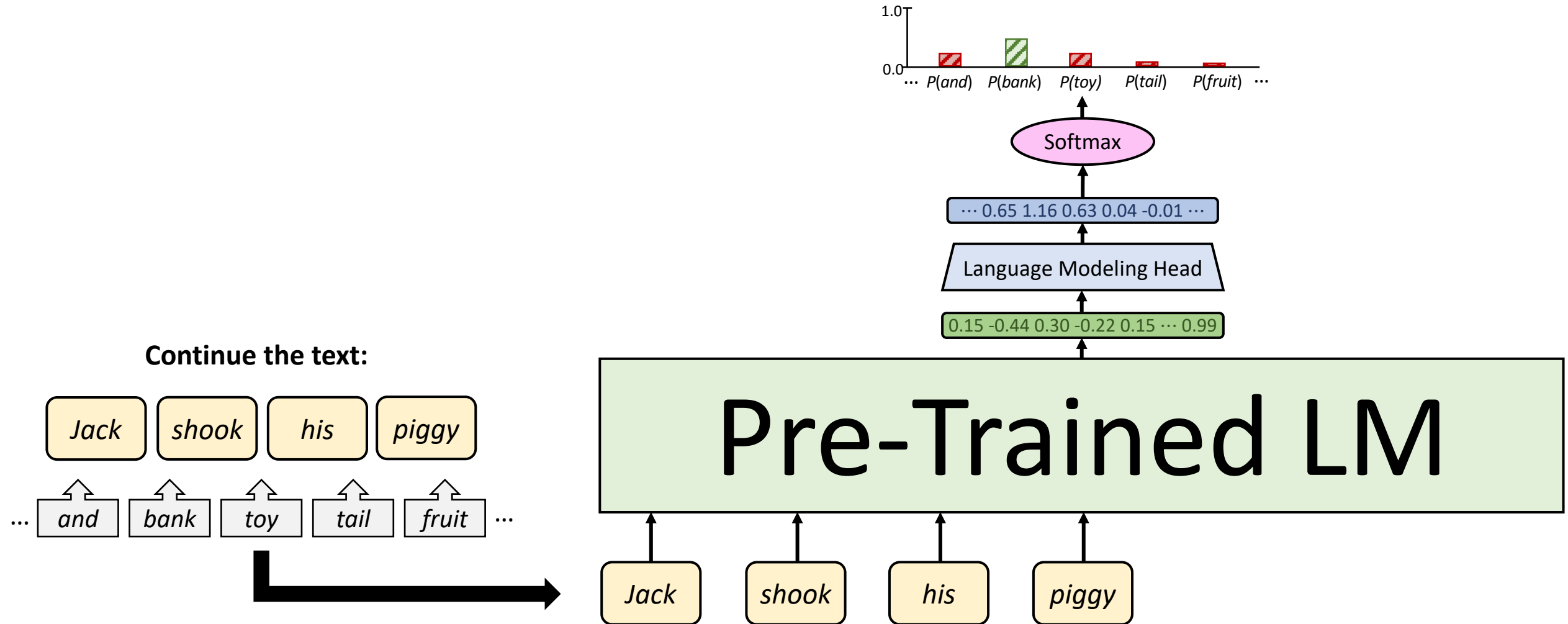
# Fine-Tuning: Multiple Choice QA



# Fine-Tuning: Token Classification



# Fine-Tuning: Text Generation



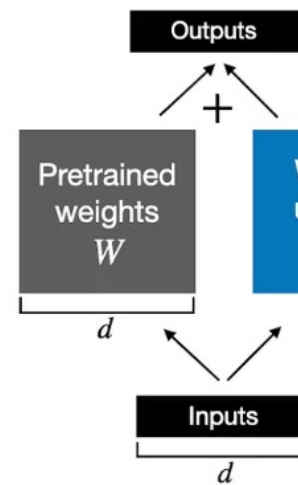
# Parameter-Efficient Fine-Tuning (PEFT)

- While fine-tuning LMs is generally more feasible when we have less available compute, there are still some problems:
  - Fine-tuning on a large amount of data can take a long time
  - The size of LM we can fine-tune is limited by compute
  - Updating all weights of the LM during fine-tuning is expensive and inefficient
- Creates a need for **parameter-efficient fine-tuning (PEFT)** methods!

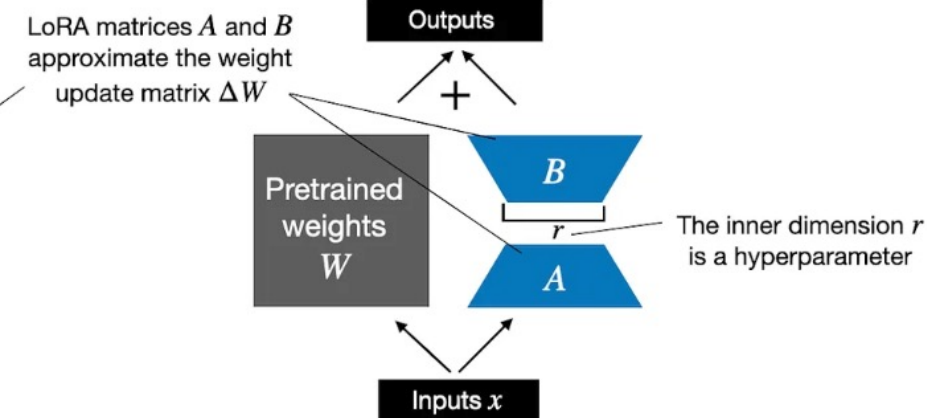
# Low-Resource Adaptation (LoRA)

- Instead of updating weights  $W$  directly during fine-tuning, learn the weight update  $\Delta W$
- Approximate  $\Delta W$  by a decomposition  $AB$ :
  - Reducing the number of learned parameters
  - Faster training with lower GPU memory requirements!
  - Pre-compute  $\Delta W$  in deployment for fast inference

Weight update in **regular finetuning**



Weight update in **LoRA**



(figure from [Sebastian Raschka](#))

# Outline

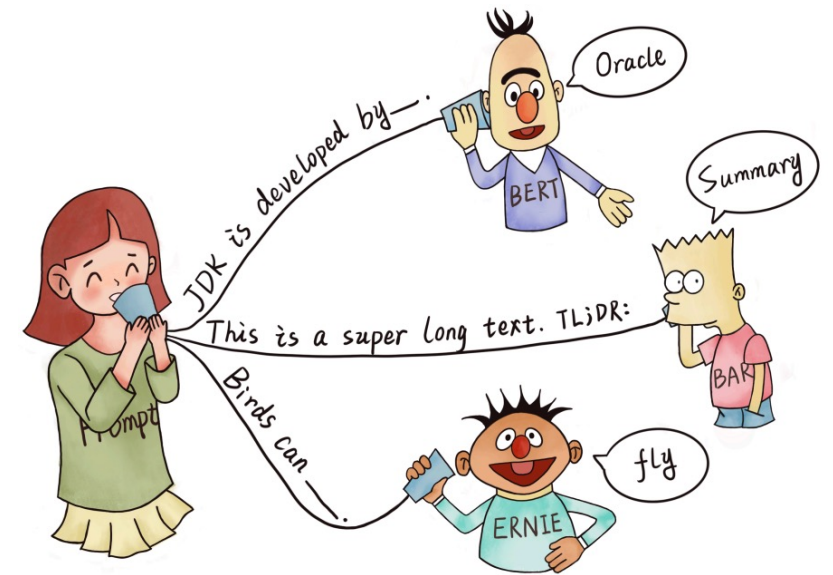
- The Road to LLMs
- Fine-Tuning LLMs
- **Prompting LLMs**



# Prompting LMs

To customize an LLM for your problem through prompting, need to make a few choices (**prompt engineering**):

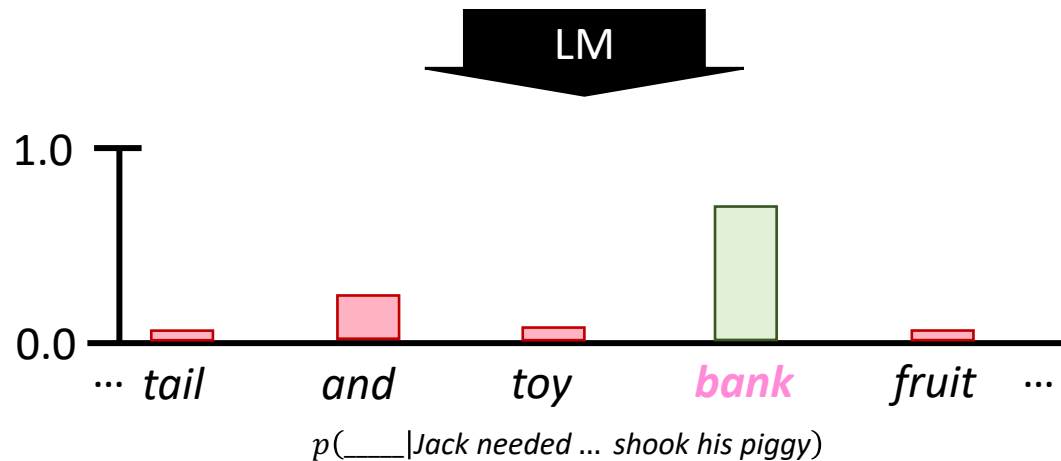
1. Prompt template
2. Answer mapping
3. In-context demonstration



# Language Models (LMs)

$$p(w_n | w_1, w_2, \dots, w_{n-1})$$

Jack needed some **money**, so he went and shook his **piggy** \_\_\_\_\_



$$p(\text{Jack needed ...shook his piggy bank}) \approx p(\text{bank} | \text{Jack needed ...shook his piggy}) \cdot p(\text{piggy} | \text{Jack needed ...shook his}) \cdot p(\text{his} | \text{Jack needed ...shook}) \cdot \dots$$

# Prompt Templates

If filling a blank from a few possible choices, can use a **cloze prompt**:

---

Task	Inputs ([X])	Template	Answer ([Z])
------	--------------	----------	--------------

---

# Prompt Templates

When completing a prompt or generating text, use a **prefix prompt**:

---

Task	Inputs ([X])	Template	Answer ([Z])
------	--------------	----------	--------------

---

# Prompt Templates

When completing a prompt or generating text, use a **prefix prompt**:

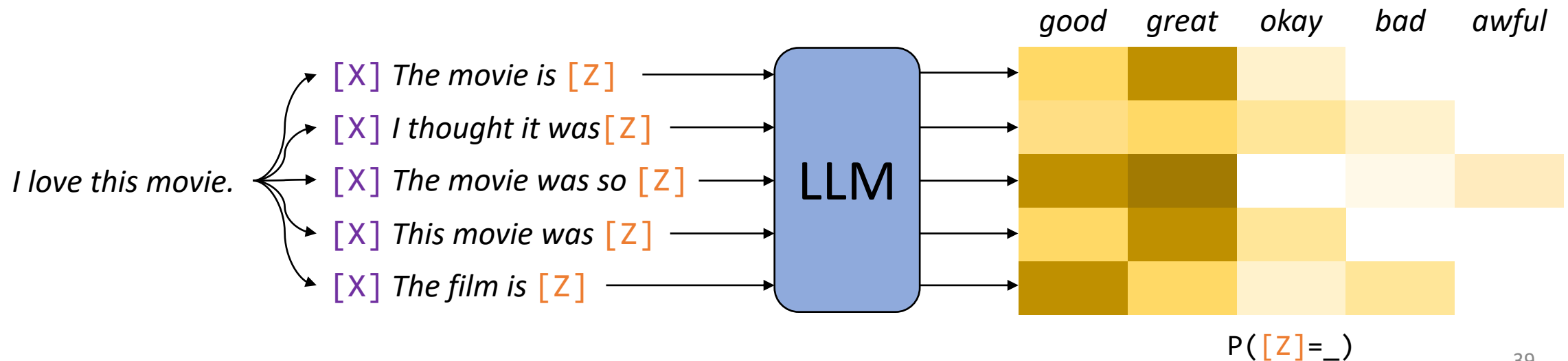
---

Task	Inputs ([X])	Template	Answer ([Z])
------	--------------	----------	--------------

---

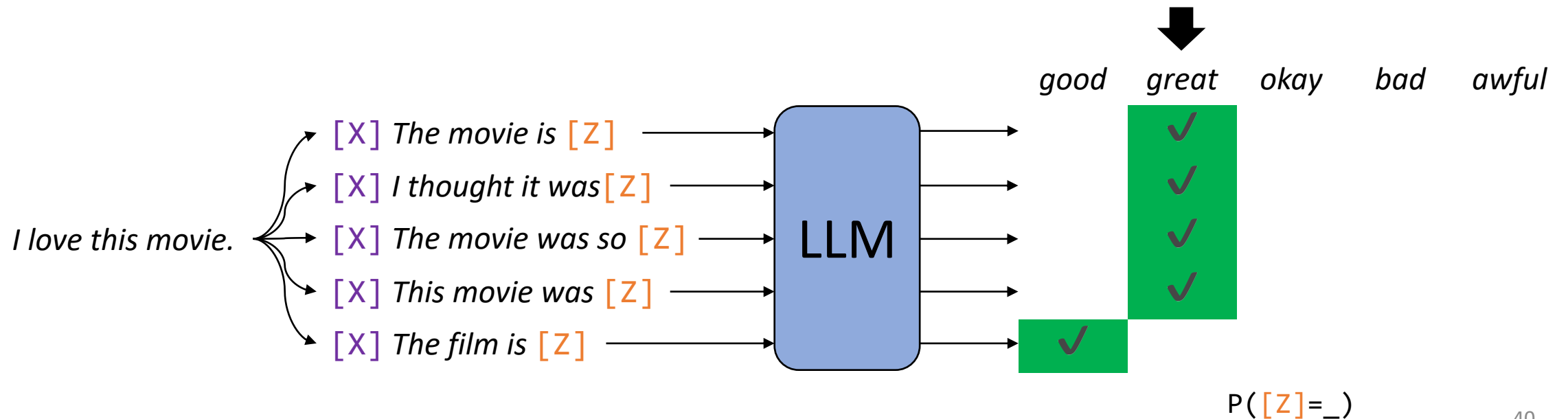
# Finding the Best Template and Answers

- Different prompts can yield different results
- May take extra work to find the best prompt
  - Trial and error
  - Ensembling templates



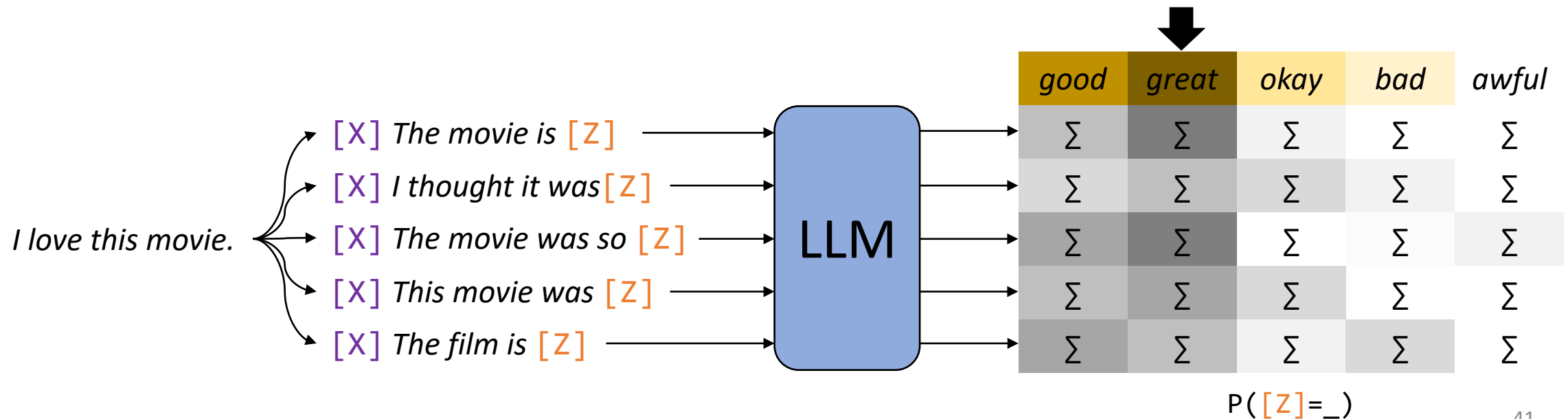
# Finding the Best Template and Answers

- Different prompts can yield different results
- May take extra work to find the best prompt
  - Trial and error
  - Ensembling templates



# Finding the Best Template and Answers

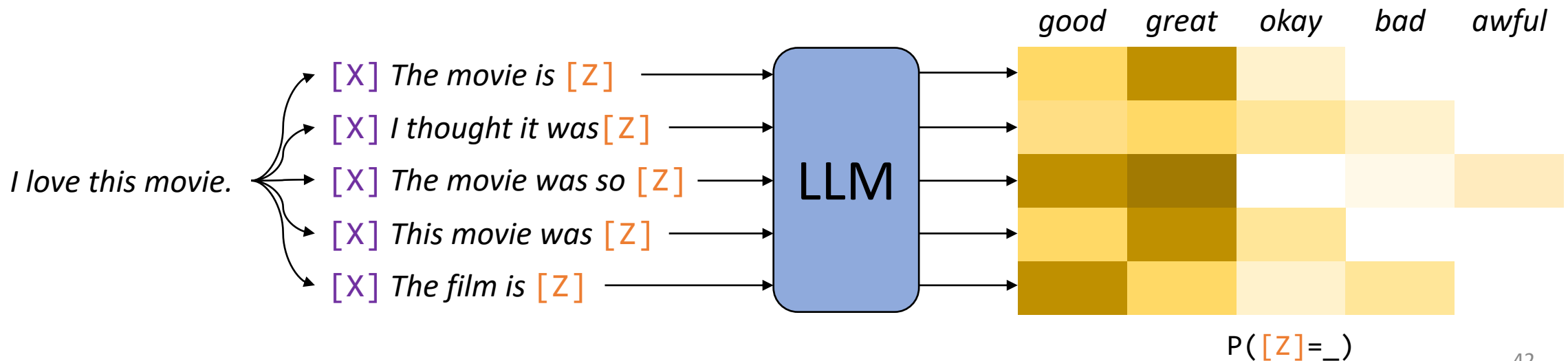
- Different prompts can yield different results
- May take extra work to find the best prompt
  - Trial and error
  - Ensembling templates





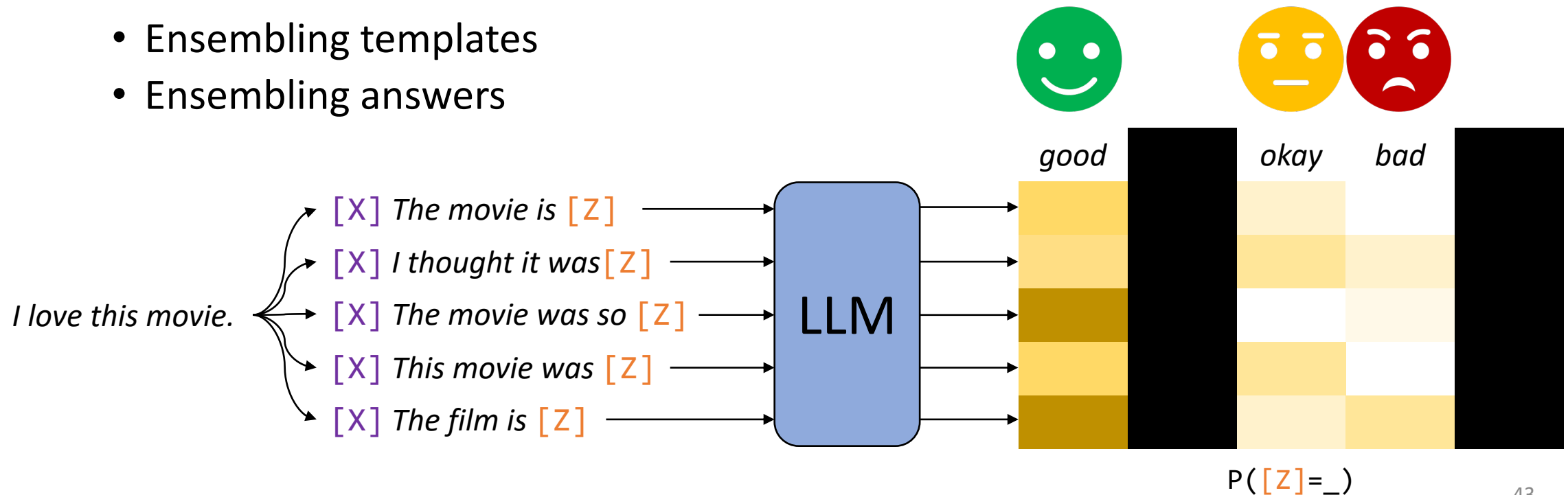
# Finding the Best Template and Answers

- Different prompts can yield different results
- May take extra work to find the best prompt
  - Trial and error
  - Ensembling templates
  - Ensembling answers



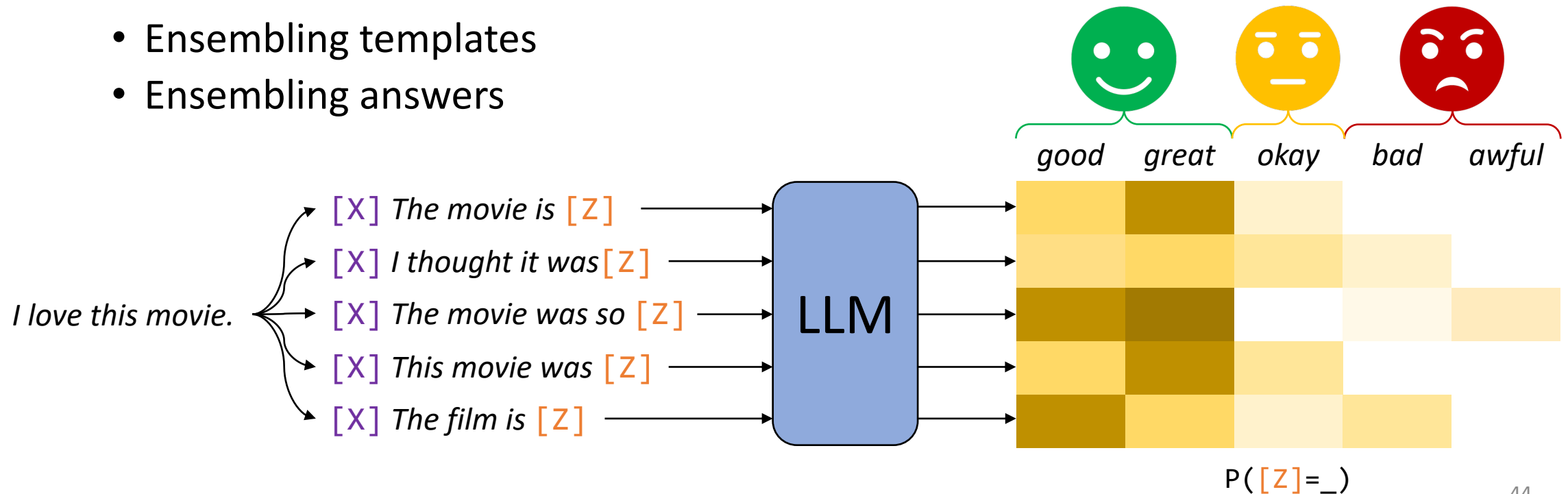
# Finding the Best Template and Answers

- Different prompts can yield different results
- May take extra work to find the best prompt
  - Trial and error
  - Ensembling templates
  - Ensembling answers



# Finding the Best Template and Answers

- Different prompts can yield different results
- May take extra work to find the best prompt
  - Trial and error
  - Ensembling templates
  - Ensembling answers



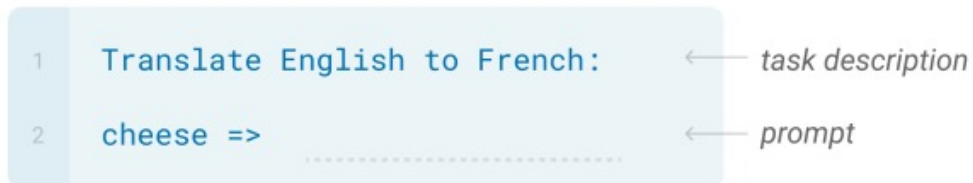
# Managing Randomness in LLMs

- LLM decoding algorithms may incorporate some randomness by default to increase the diversity of generation
- Some solutions:
  - Generate multiple times and average results
  - Greedy decoding

# In-Context Learning

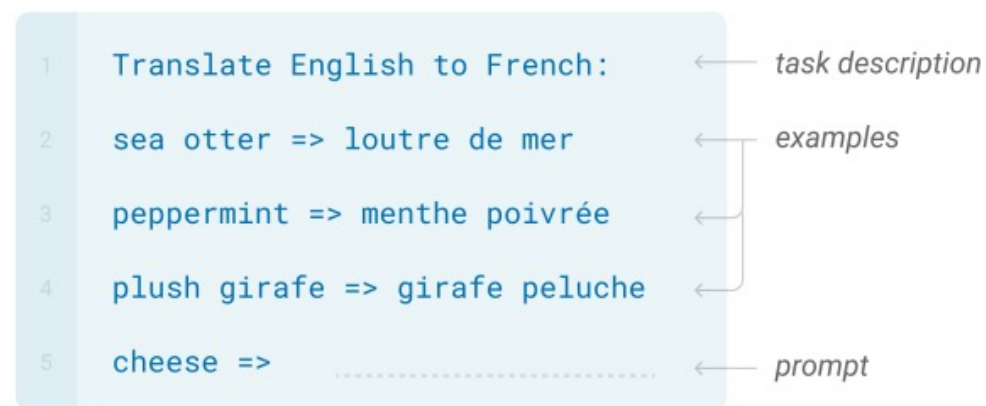
## Zero-shot

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.



## Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.



# Chain-of-Thought Prompting

## Standard Prompting

### 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

### Input

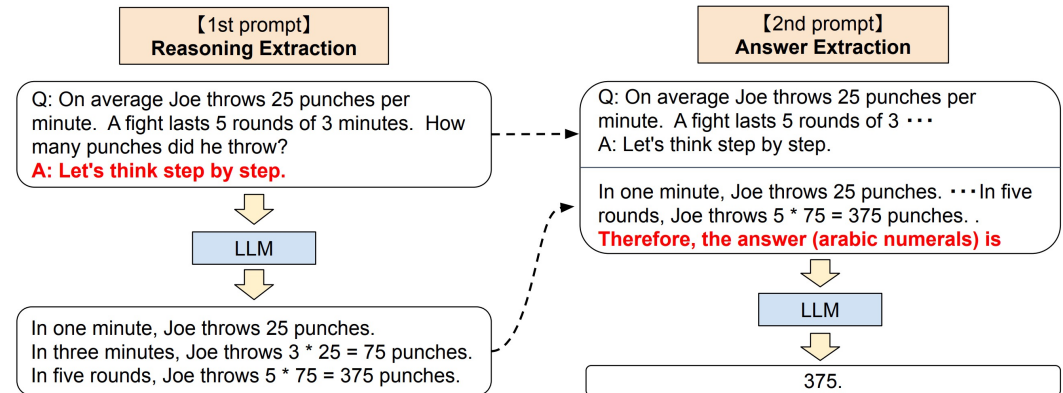
Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

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. ✅



# *Next: From Theory to Practice!*

I'm on the job market for academic and industry positions!



@shanestorks



[www.shanestorks.com](http://www.shanestorks.com)