

# In-context Learning

CSCI 601 471/671  
NLP: Self-Supervised Models

<https://self-supervised.cs.jhu.edu/sp2023/>

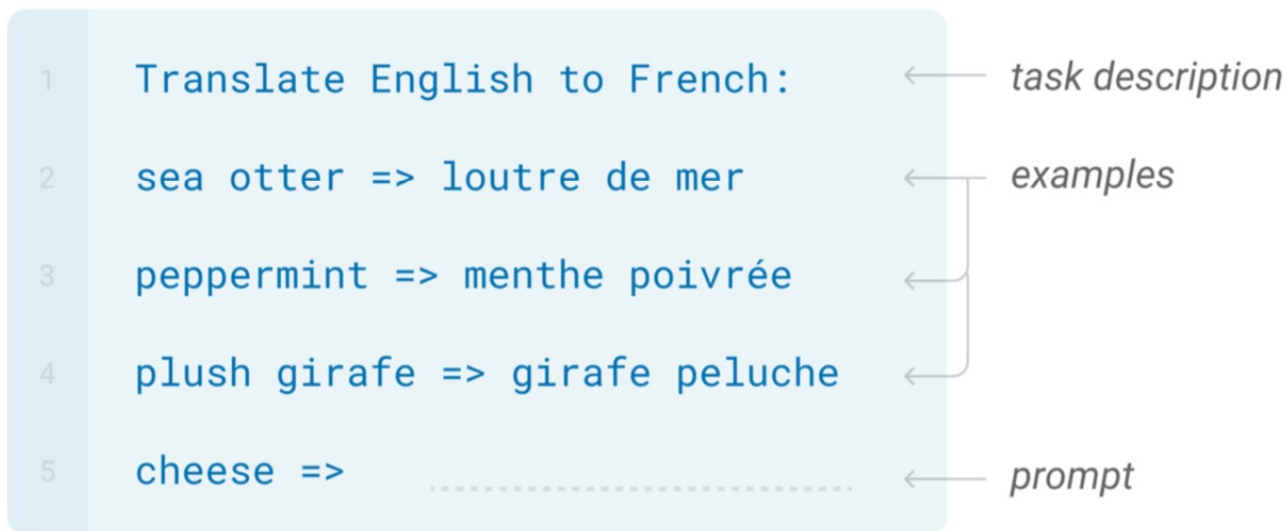


JOHNS HOPKINS  
UNIVERSITY

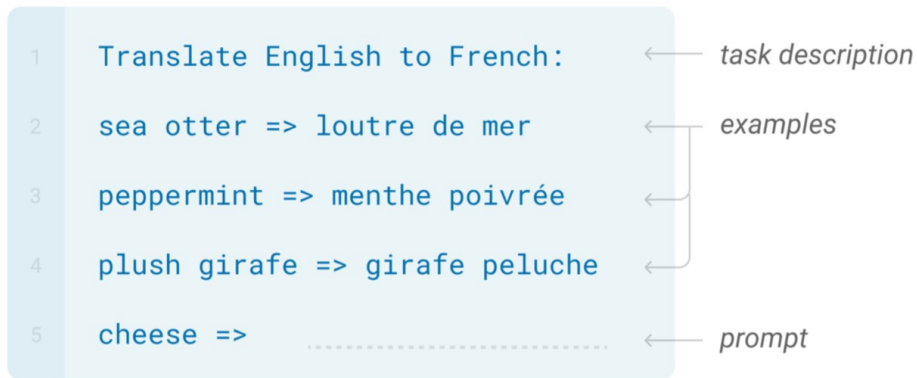
# Logistics

- HW6: The Struggles ... 😞
- HW7 will be the last homework! 😊
- Projects
  - Do you have an idea yet?
  - Do you have teammate(s) already?
  - Have you looked at the deck of ideas I shared yet?
- Project proposal deadline: Thu Mar 16
  - Should we push it back? Alternative: Tue Mar 28

# In-Context Learning



# In-Context Learning



- Learns to do a downstream task by conditioning on input-output examples!
- **No weight update** — our model is not **explicitly pre-trained** to learn from examples
  - The underlying models are quite general
- Today's focus:
  - How to use effectively in practice?
  - Fundamentally, why does it work?

# Why Do We Care About Few-Shot Learning?

Practically Useful

Intellectually Intriguing

# Practically Useful

Labeling data is costly

- May require domain expertise
  - Medical, legal, financial
- Inputs may be long/complex
  - Grammaticality
- Outputs may be complex
  - Semantic parsing

# Practically Useful

Labeling data is costly

You want to do best with what you have

- You don't want to get more data
- Emergent, time-sensitive scenarios
  - Something new happened
  - Need to react quickly!

# Practically Useful

Labeling data is costly

You want to do best with what you have

Finetuning can unstable

- Training is sensitive to hyperparams
- Not enough validation data
- We don't quite understand how finetuning works



# Practically Useful

Labeling data is costly

You want to do best with what you have

Finetuning can be unstable

Finetuning large LMs is expensive

- Expensive to train, time and memory

# Intellectually Intriguing

Potential test for “Intelligent Behavior”

- Generalization from few examples
  - Fundamental piece of intelligence
  - Often used in psychology
  - Quickly adjust to environment

# Intellectually Intriguing

Potential test for “Intelligent Behavior”

“But... Deep learning is data hungry!”

- Long-standing criticism of DL
- Understand why it doesn't work here
  - Or does it?
- What are the new limitations of DL?

# Intellectually Intriguing

Potential test for “Intelligent Behavior”

“But... Deep learning is data hungry!”

Insights into Language Modeling

- What does an LLM “know”?
- What are the biases/limitations of LLMs?
- ...

# Intellectually Intriguing

Potential test for “Intelligent Behavior”

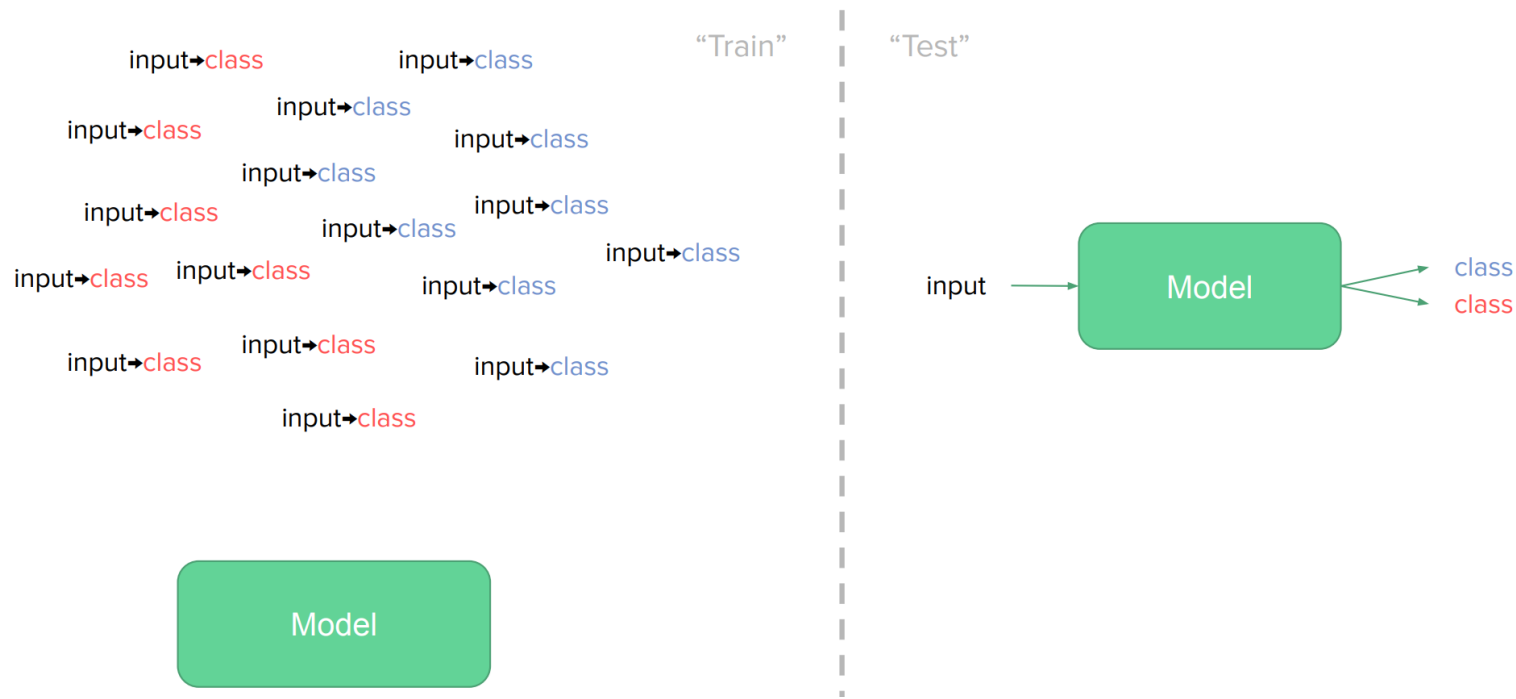
“But... Deep learning is data hungry!”

Insights into Language Modeling

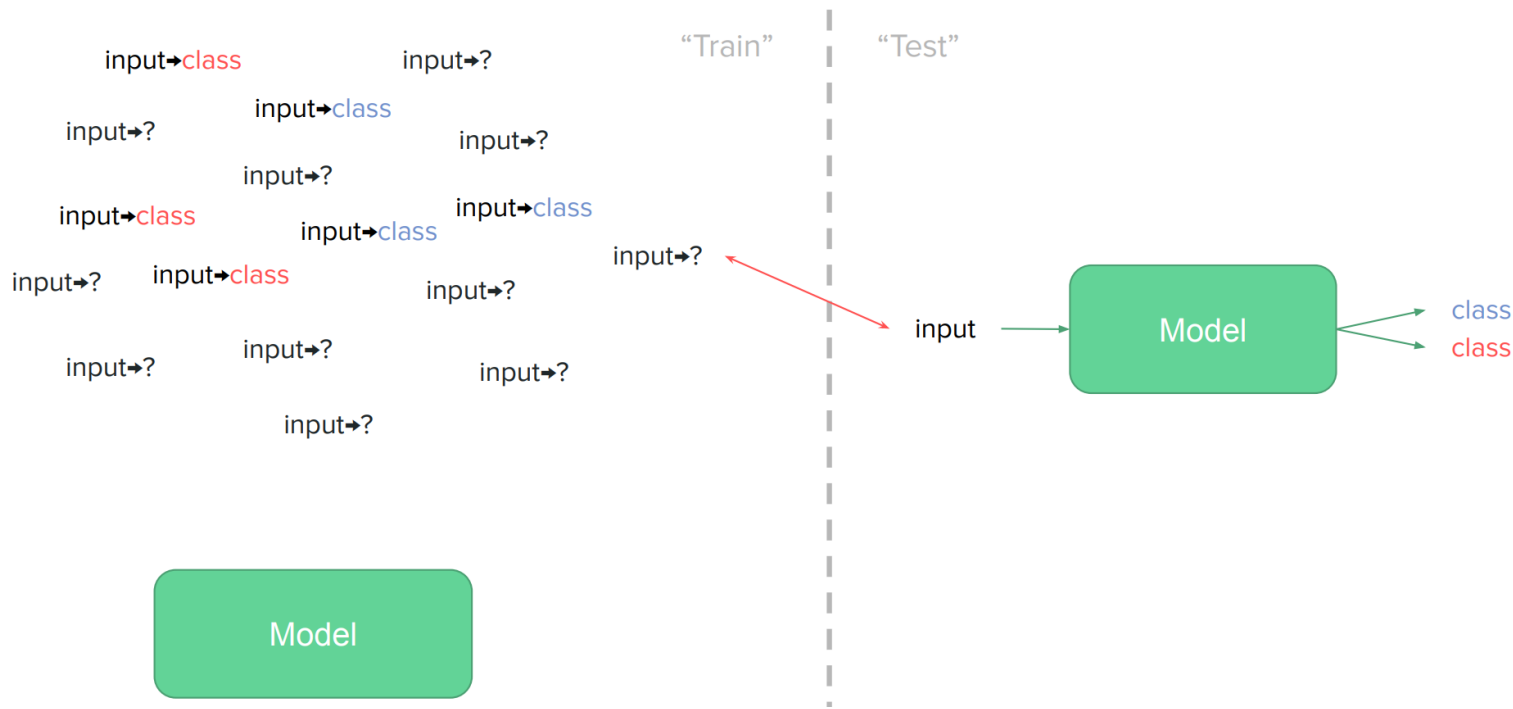
Because LARGE language models

- Training/inference/access is tough
- What else can we do? ;)

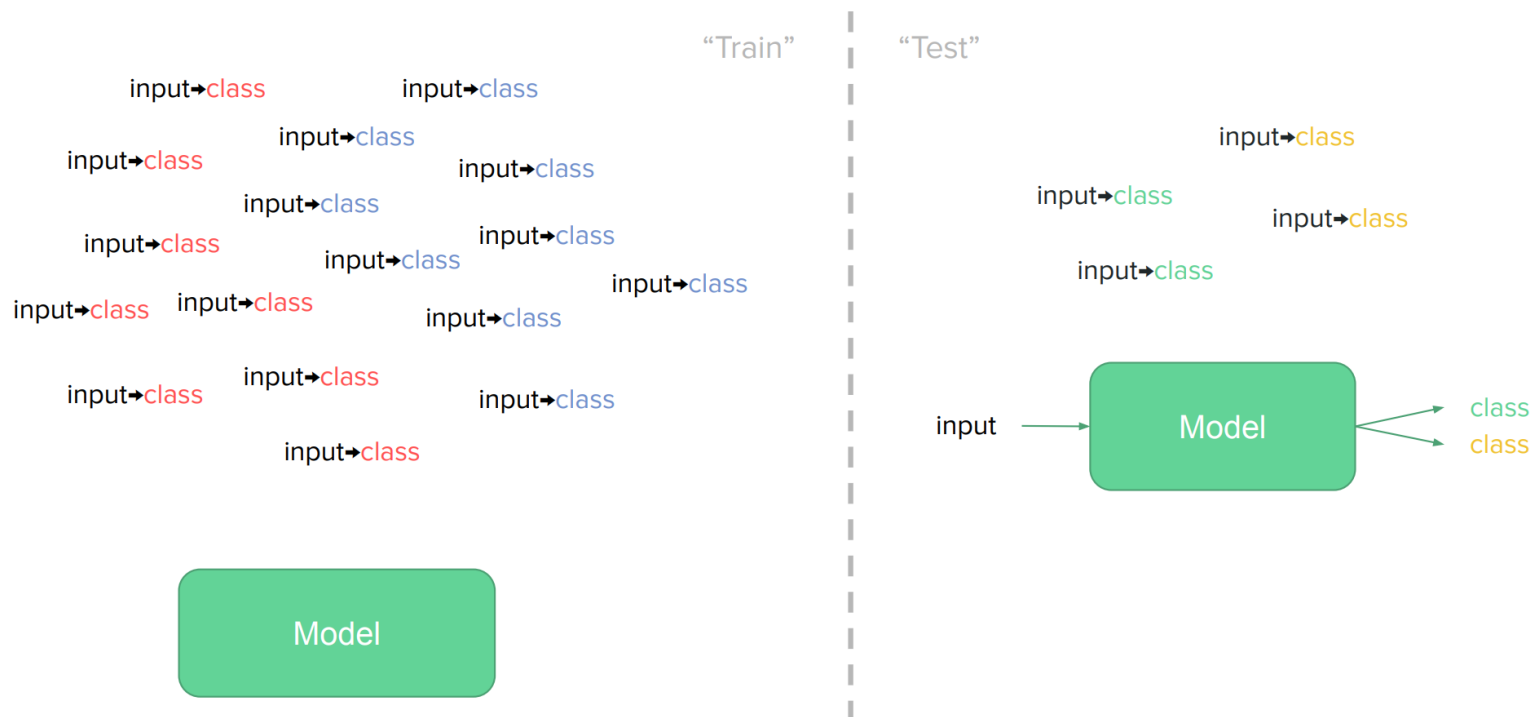
# The Broader Context: Supervised Learning



# The Broader Context: Semi-Supervised Learning

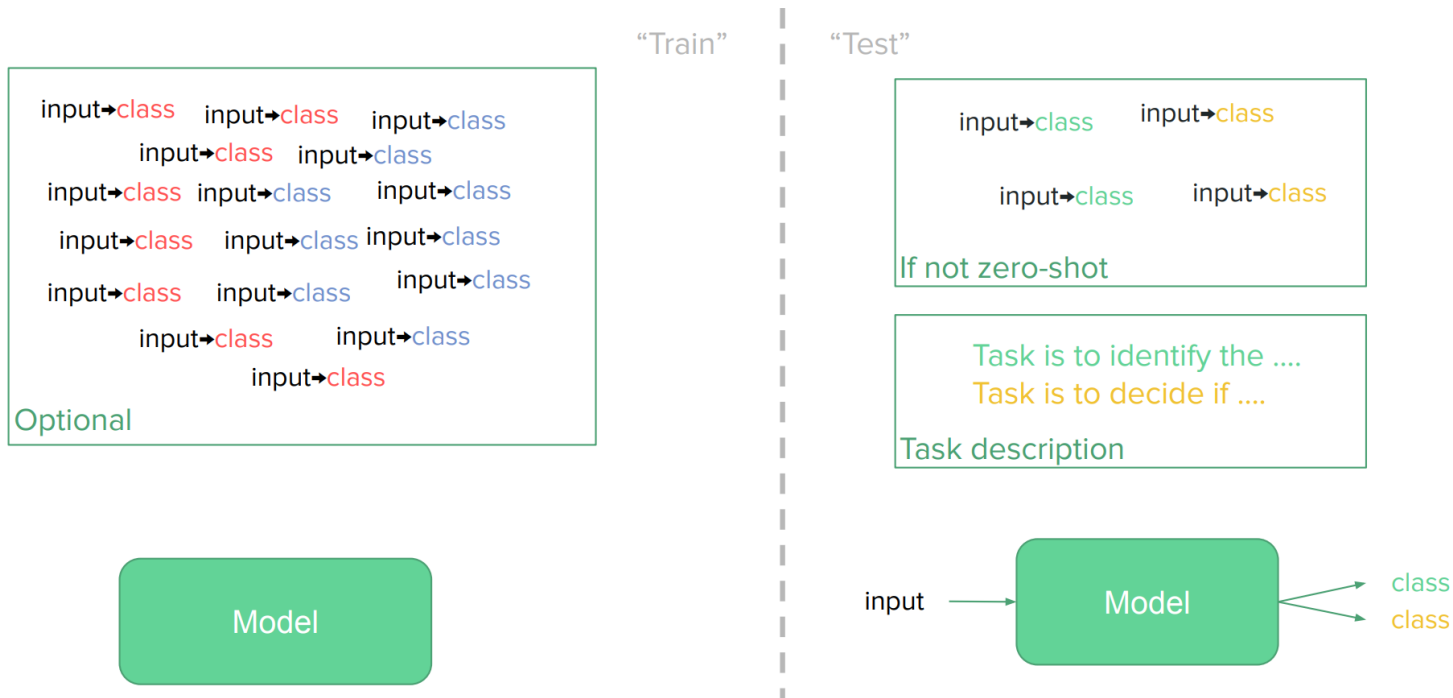


# The Broader Context: [Traditional] Few-shot Learning





# The Broader Context: [Modern] Few-shot Learning



# In-Context Prompting

## Movie review dataset

Input: An effortlessly accomplished and richly resonant work.

Label: positive

Input: A mostly tired retread of several other mob tales.

Label: negative

An effortlessly accomplished and richly resonant work. It was great! A mostly tired retread of several other mob tales. It was terrible!

A three-hour cinema master class. It was \_\_\_\_\_

Language Model

$p_1 = P(\text{It was great!} \mid \text{1st train input+output} \setminus \text{2nd train input+output} \setminus \text{A three-hour cinema master class.})$

$p_2 = P(\text{It was terrible!} \mid \text{1st train input+output} \setminus \text{2nd train input+output} \setminus \text{A three-hour cinema master class.})$

$p_1 > p_2$       "positive"

$p_1 < p_2$       "negative"

# LM Prompting: Choices of Encoding

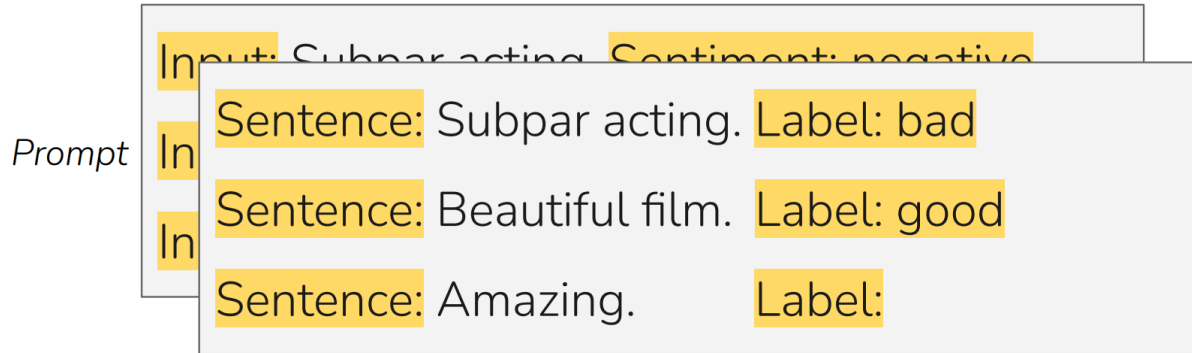
*Prompt*

Input: Subpar acting. Sentiment: negative

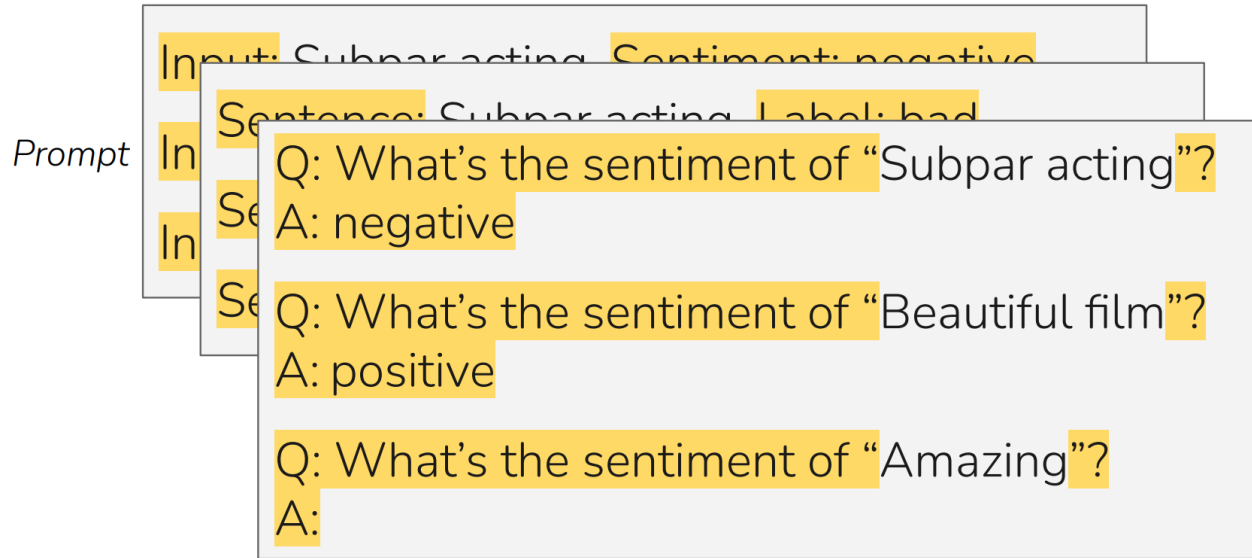
Input: Beautiful film. Sentiment: positive

Input: Amazing. Sentiment:

# LM Prompting: Choices of Encoding



# LM Prompting: Choices of Encoding



# LM Prompting: Choices of Encoding

- **Pattern:** A function that encodes the inputs.
- **Verbalizer:** A function that encodes the output.

**Pattern:**  $f(\langle x \rangle) = \text{"Input: } \langle x \rangle \text{"}$

**Verbalizer:**  $v(\langle x \rangle) = \text{"Label: } \langle x \rangle \text{"}$

**Pattern:**  $f(\langle x \rangle) = \text{"Q: What is the sentiment of } \langle x \rangle \text{"}$

**Verbalizer:**  $v(\langle x \rangle) = \text{"A: } \langle x \rangle \text{"}$

**Input:** Subpar acting. **Sentiment:** negative

**Input:** Beautiful film. **Sentiment:** positive

**Input:** Amazing. **Sentiment:**

**Q:** What's the sentiment of "Subpar acting"?

**A:** negative

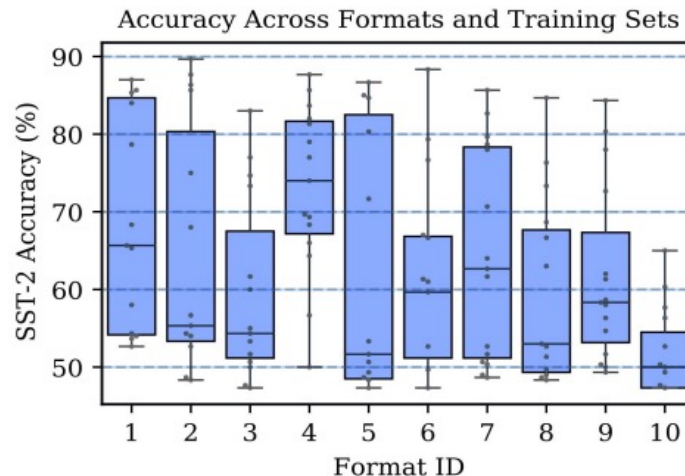
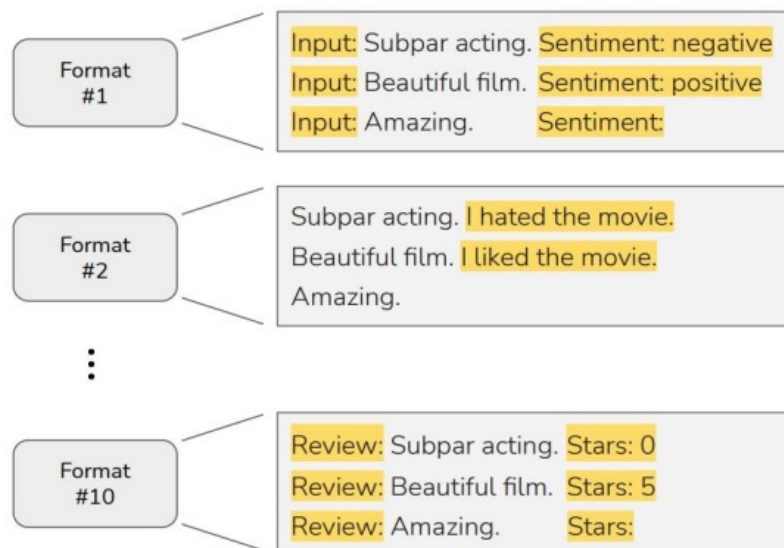
**Q:** What's the sentiment of "Beautiful film"?

**A:** positive

**Q:** What's the sentiment of "Amazing"?

**A:**

# In-Context Learning: Sensitivity to Encoding



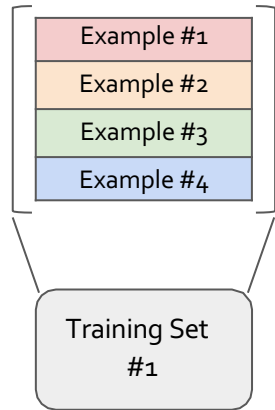
In-context learning is highly sensitive to prompt format (training sets and patterns/verbalizers)

# In-Context Learning: Sensitivity to Demo. Permutations

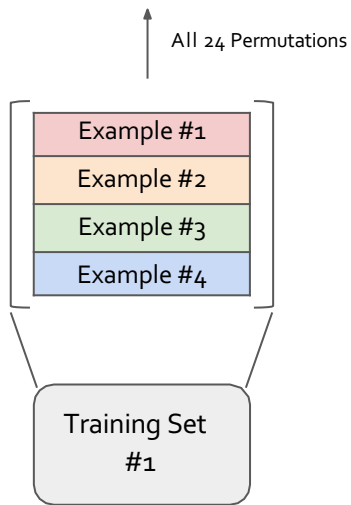
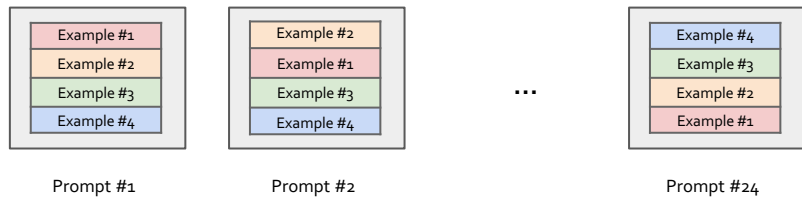
Training Set  
#1



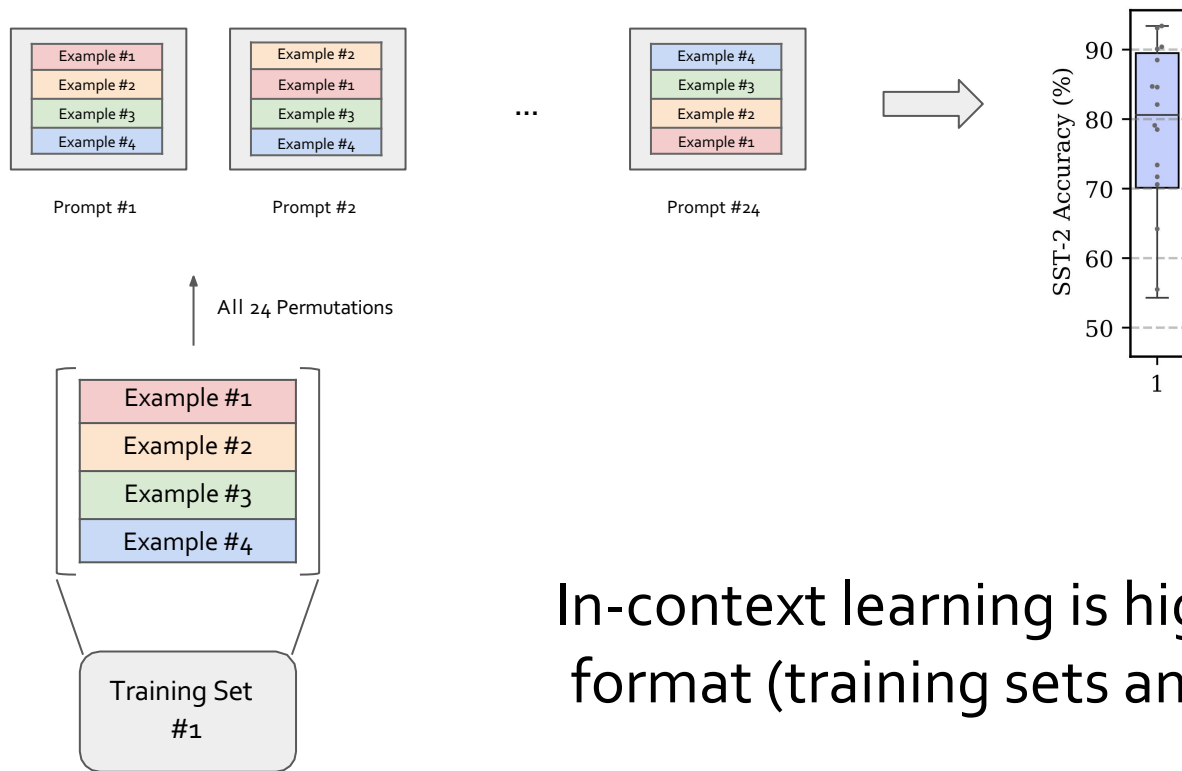
# In-Context Learning: Sensitivity to Demo. Permutations



# In-Context Learning: Sensitivity to Demo. Permutations

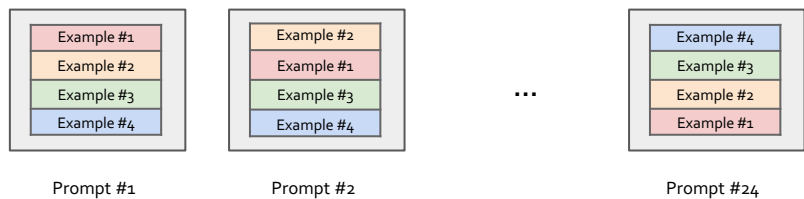


# In-Context Learning: Sensitivity to Demo. Permutations

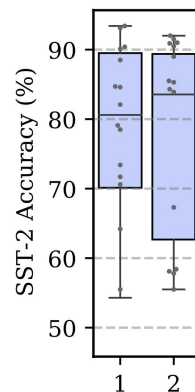
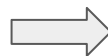
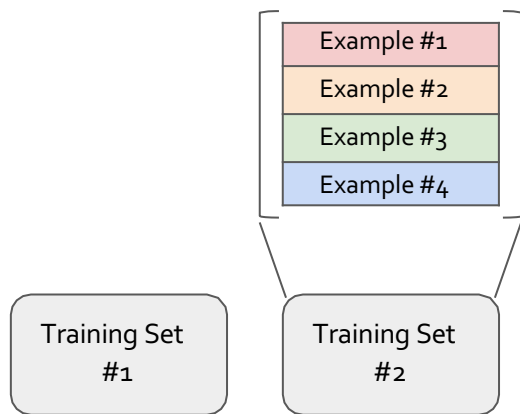


In-context learning is highly sensitive to prompt format (training sets and patterns/verbalizers)

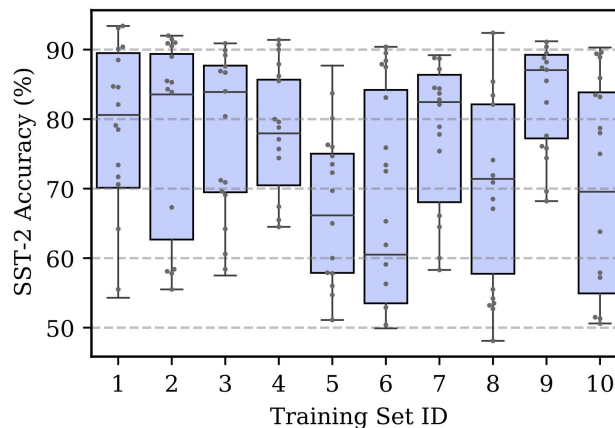
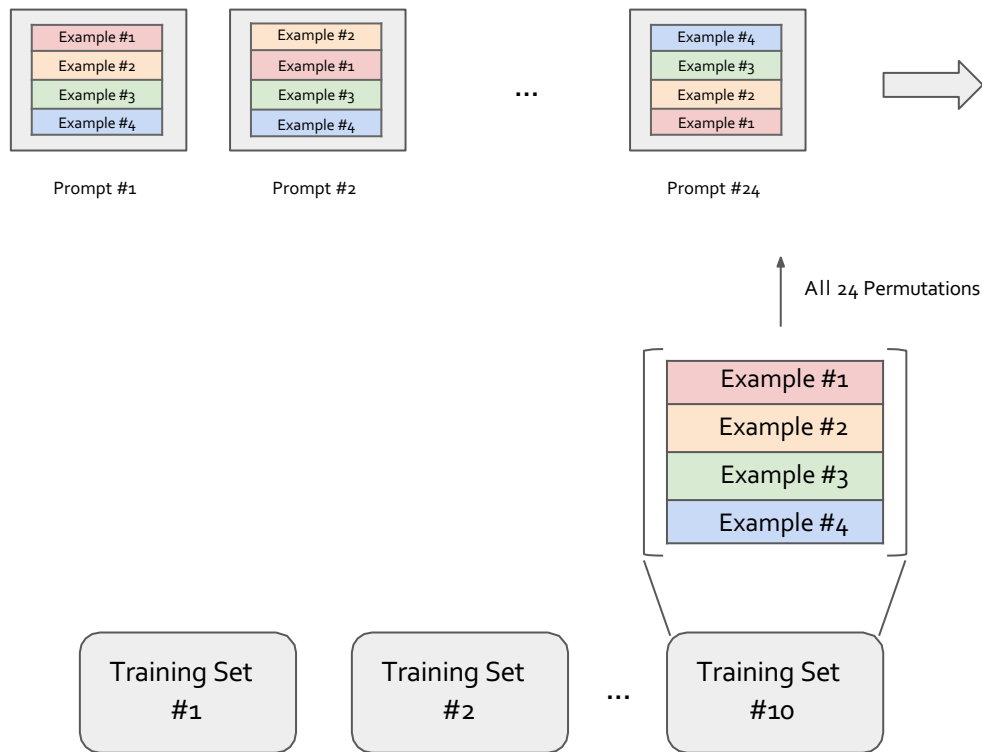
# In-Context Learning: Sensitivity to Demo. Permutations



All 24 Permutations



# In-Context Learning: Sensitivity to Demo. Permutations



The choice of demonstrations and their order is quite important.

# Sensitivity to Wording (Framing) of Prompts

- Framing of prompts matters a lot.

Craft a question that requires commonsense to be answered. Based on the given context, craft a common-sense question, especially those that are LONG, INTERESTING, and COMPLEX. The goal is to write questions that are easy for humans and hard for AI machines! To create such questions, here are some suggestions: A. What may (or may not) be the plausible reason for an event? B. What may (or may not) happen before (or after, or during) an event? ...



Generate questions such that you use

- 'what may happen',
- 'will ...?',
- 'why might',
- 'what may have caused',
- 'what may be true about',
- 'what is probably true about',
- 'what must'

and similar phrases in your question based on the input context.

# Sensitivity to Wording (Framing) of Prompts

- Prompts can often be phrased in a language that are easier to be understood by language models.
- Generally, it is easier for LMs to follow shorter, crisp, itemized prompts.

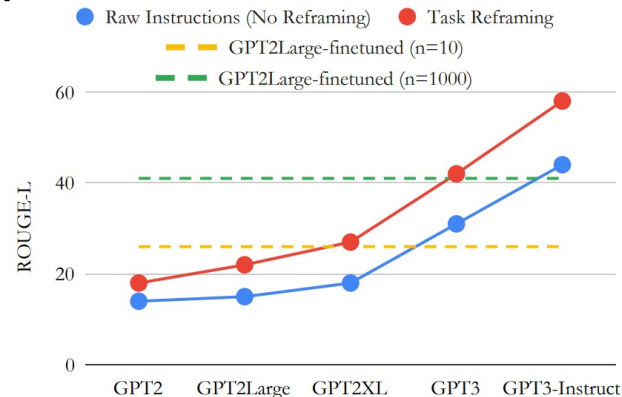


Figure 2: Across a variety of model sizes, **reframed prompts** consistently show considerable performance gain over **raw task instructions (no reframing)** in a few-shot learning setup. Since fine-tuning GPT3 is prohibitively expensive, we show the performance of fine-tuning smaller models (**horizontal lines**). This results indicates that *evaluating* reframed prompts on a large model like GPT3-instruct (red line) might be more effective than *fine-tuning* a smaller model like GPT2Large (green line) with 200× more data. Details of the experiments in §4.

# Summary Thus Far

- There are many possible ways to encode in-context examples of a **fixed task**
  - Many possible patterns/verbalizers
  - The choice of demonstrations
  - Ordering of the examples
  - ....
- It turns out there is a **huge variance** in performance depending on the encoding.
  - You can treat them as hyper-parameters
  - You should **not** choose these encodings based on the test data.
- Generally, it is better to use encoding that **makes the sequence closer to language modeling** — closer to what is observed during pretraining.



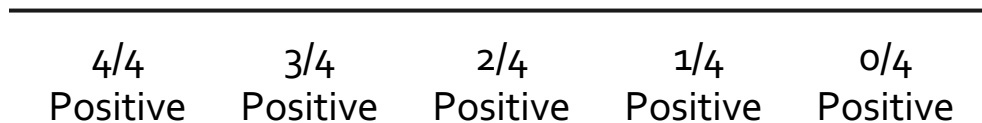
# What Causes These Variances?

- Here we will provide several justifications ...

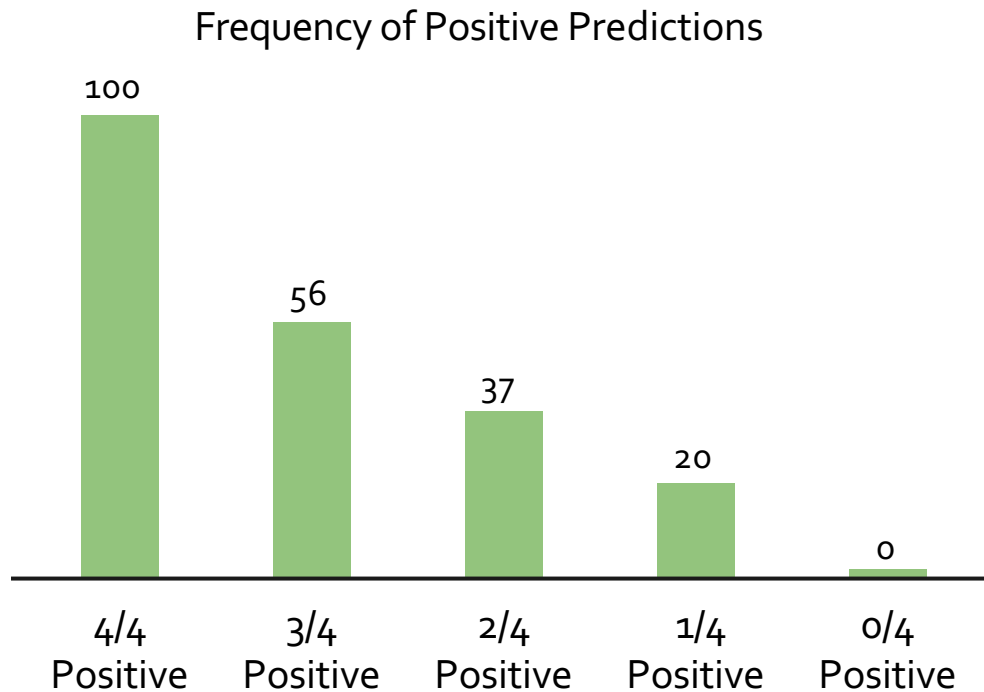
# Majority Label Bias

# Majority Label Bias

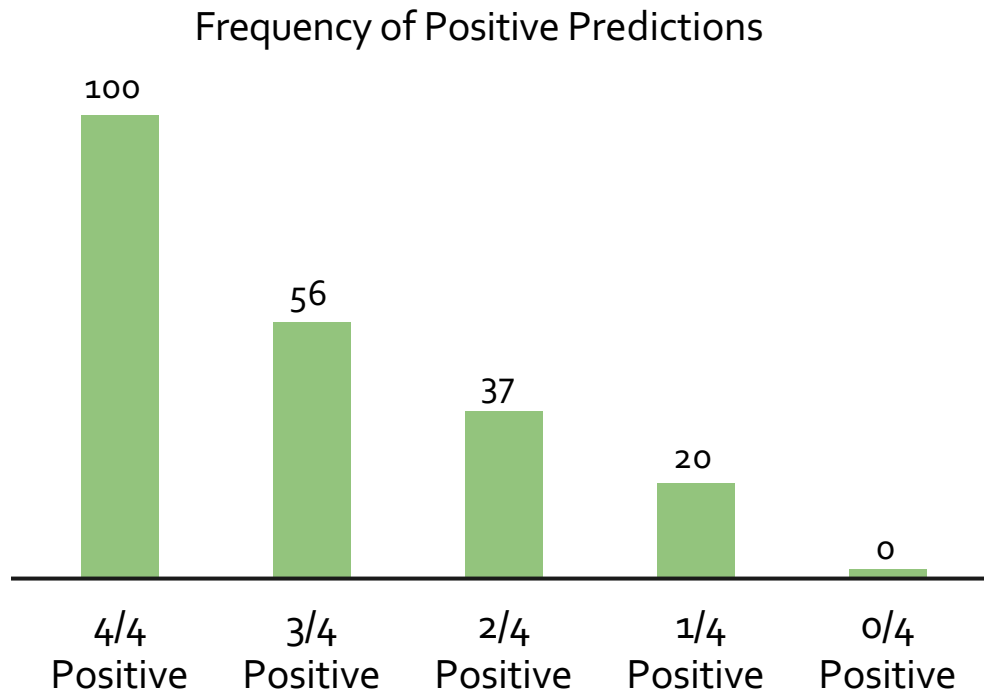
Frequency of Positive Predictions



# Majority Label Bias



# Majority Label Bias



Majority label bias: frequent training answers dominate predictions  
Explain some of the variances across example selections

# Recency Bias

# Recency Bias

Frequency of Positive Predictions

---

NPPP

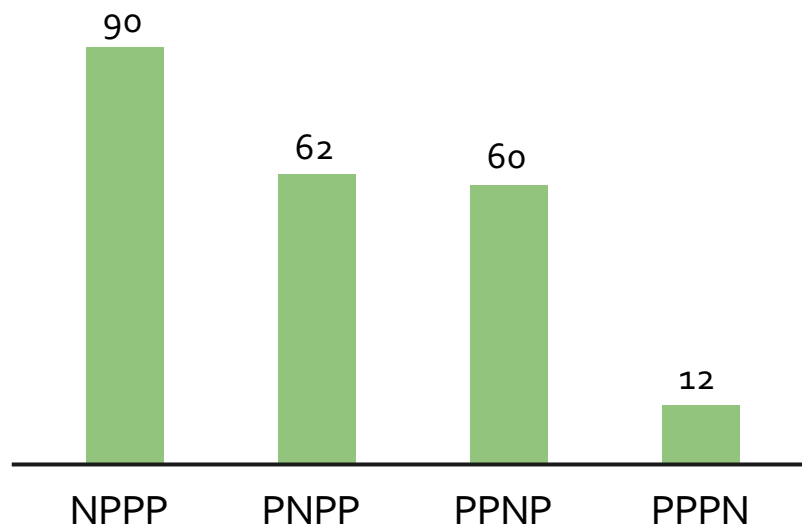
PNPP

PPNP

PPPN

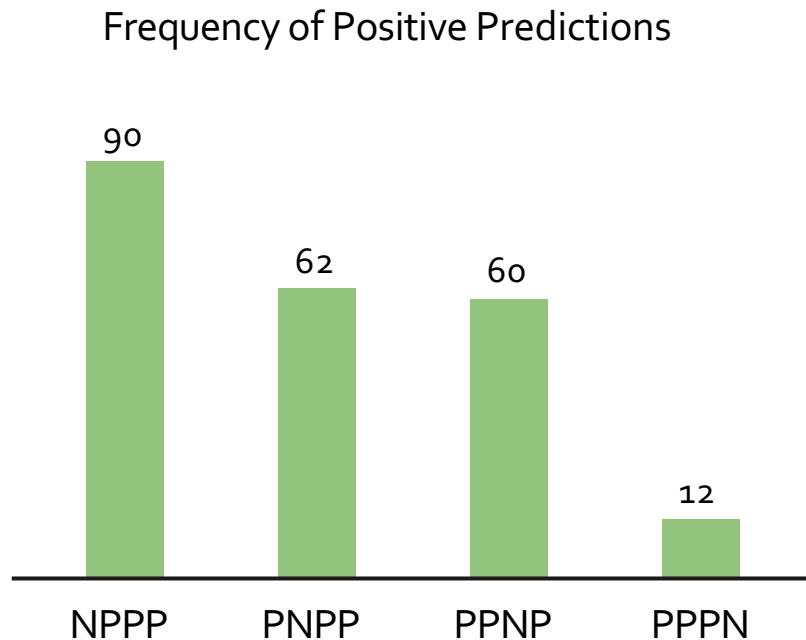
# Recency Bias

Frequency of Positive Predictions





# Recency Bias

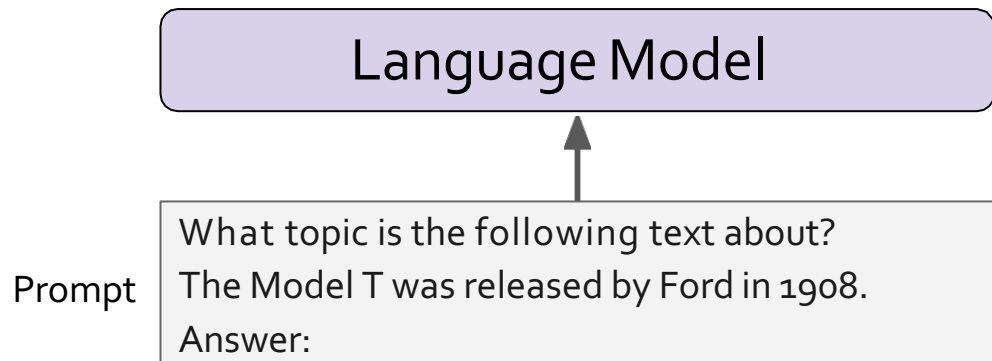


Recency bias: examples near end of prompt dominate predictions

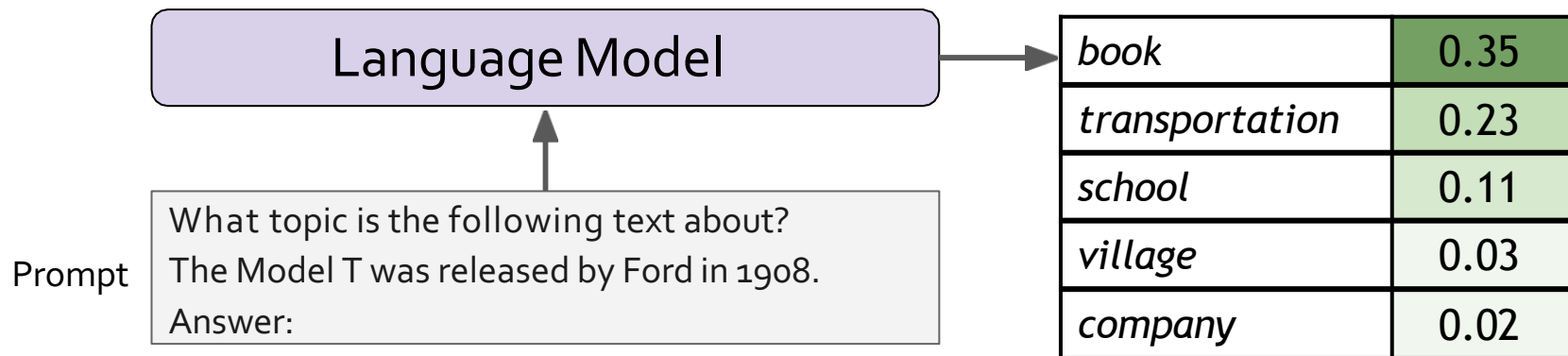
- Explains variance across example permutations

# Common Token Bias

# Common Token Bias



# Common Token Bias

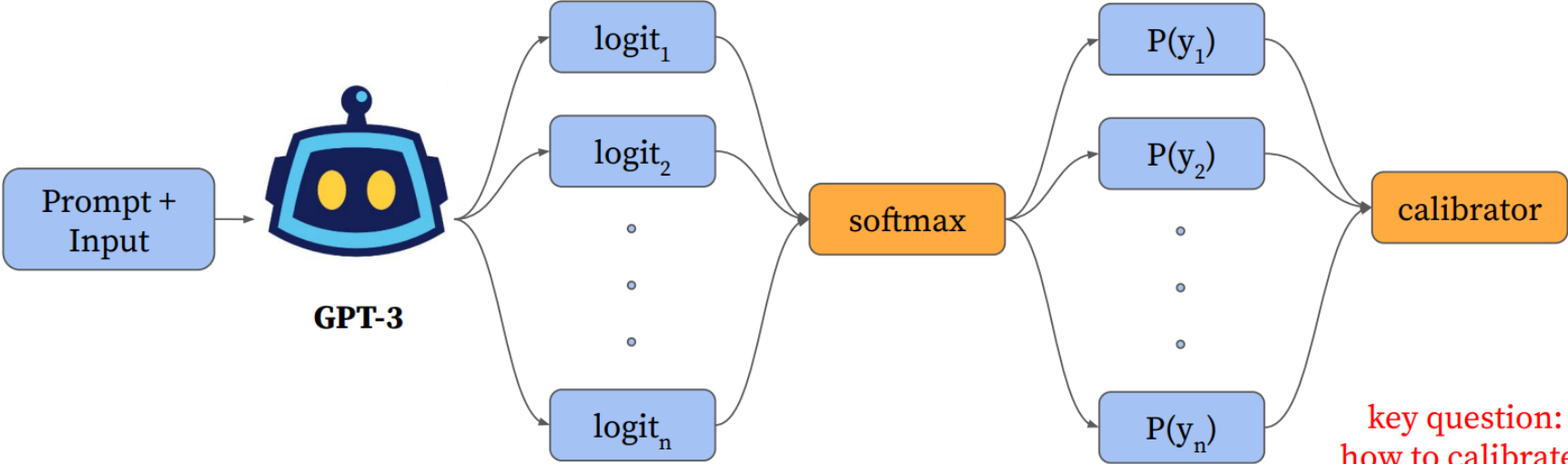


Model is biased towards predicting the incorrect frequent token "book" even when both "book" and "transportation" are equally likely labels in the dataset

	Token	Web (%)	Label (%)	Prediction (%)
✗	<i>book</i>	0.026	9	<u>29</u>
✓	<i>transportation</i>	0.0000006	9	<u>4</u>

- Common token bias: common n-grams dominate predictions
  - helps explain variance across prompt formats

# Calibrating LM Probabilities



[Slide credit: Howard Yen, Vishvak Murahari]

# Calibrating LM Probabilities

Step 1: Estimate the bias

# Calibrating LM Probabilities

## Step 1: Estimate the bias

Insert "content-free" test input into prompt

Input: Subpar acting. Sentiment: negative

Input: Beautiful film. Sentiment: positive

Input: \_\_\_\_\_ Sentiment:

# Calibrating LM Probabilities

## Step 1: Estimate the bias

Insert "content-free" test input into prompt

Input: Subpar acting. Sentiment: negative

Input: Beautiful film. Sentiment: positive

Input: \_\_\_\_\_ Sentiment:

Get model's prediction

<i>positive</i>	0.65
<i>negative</i>	0.35



# Calibrating LM Probabilities

## Step 1: Estimate the bias

Insert "content-free" test input into prompt

Input: Subpar acting. Sentiment: negative

Input: Beautiful film. Sentiment: positive

Input: \_\_\_\_\_ Sentiment:

Get model's prediction

<i>positive</i>	0.65
<i>negative</i>	0.35

## Step 2: Counter the bias

# Calibrating LM Probabilities

## Step 1: Estimate the bias

Insert “content-free” test input into prompt

Input: Subpar acting. Sentiment: negative  
Input: Beautiful film. Sentiment: positive  
Input: \_\_\_\_\_ Sentiment: \_\_\_\_\_

Get model’s prediction

<i>positive</i>	0.65
<i>negative</i>	0.35

## Step 2: Counter the bias (“Platt Scaling”)

“Calibrate” predictions with affine transformation

$$\hat{\mathbf{q}} = \text{softmax}(\mathbf{W}\hat{\mathbf{p}} + \mathbf{b})$$

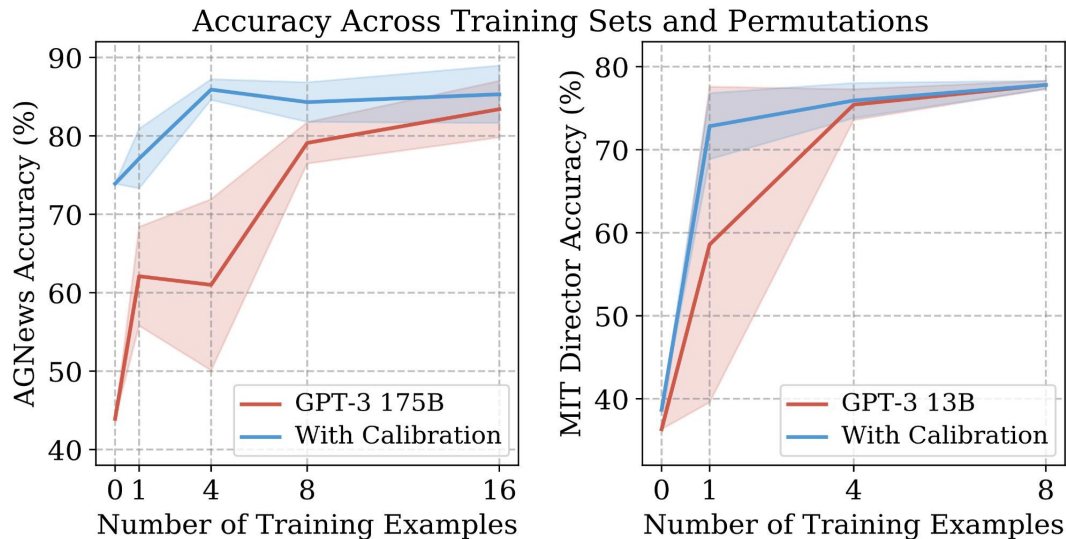






# Effect of Calibration

- Improves mean and worst-case accuracy
- Reduces variance across training sets and permutations



# Surface Form Competition

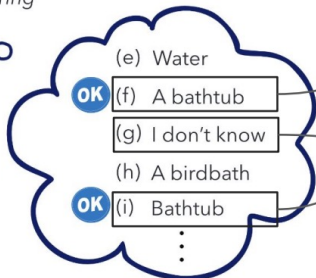
A human wants to submerge himself in water, what should he use?

Humans *select* options



- ✗ (a) Coffee cup
- ✓ (b) Whirlpool bath
- ✗ (c) Cup
- ✗ (d) Puddle

Language Models assign probability to every possible string



OK = right concept, wrong surface form

$$P(\text{Bathtub} \mid x) = 0.8$$

$$P(\text{Whirlpool bath} \mid x) \leq 0.2$$

**Competes for probability mass**



**Generic output always assigned high probability**

Every correct string is assigned lower scores than expected

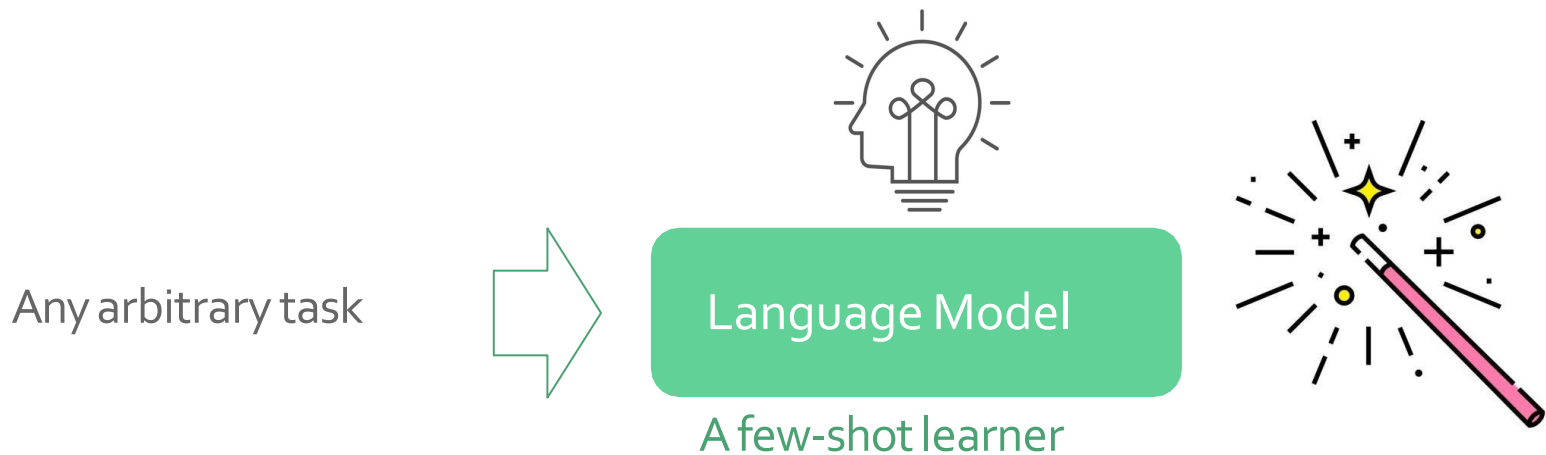
Surface forms are competing for probability mass — skew the probabilities.  
There are some ideas on how to calibrate for these issues (see [Holtzman et al 2021](#)).

# Summary Thus Far

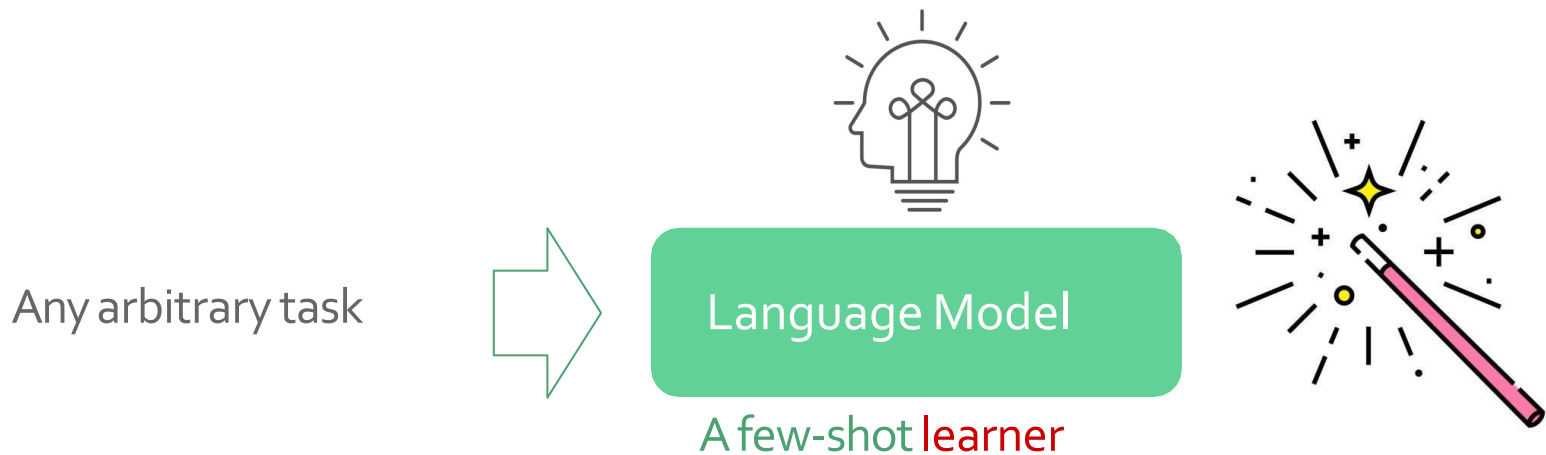
- LM prompting & In-context learning show promising results, but their performance is highly unstable/brittle.
- Better scoring: Calibration
- Other factors:
  - Better **formation** of demonstrations
  - Better **choice** of demonstrative examples
  - Better **ordering** of demonstrative examples



# How/Why does In-context Learning Work?

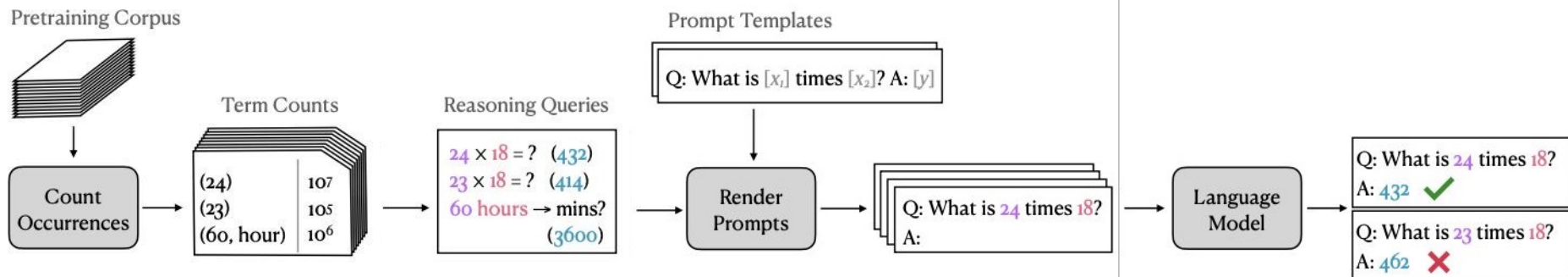


# How/Why does In-context Learning Work?



# Impact of Pretraining Term Frequencies

- For each task, identify relevant terms from each instance—numbers and units
- Count co-occurrences of these terms in the pretraining data (term pairs or triples within a fixed window)



## Impact of Pretraining Term Frequencies

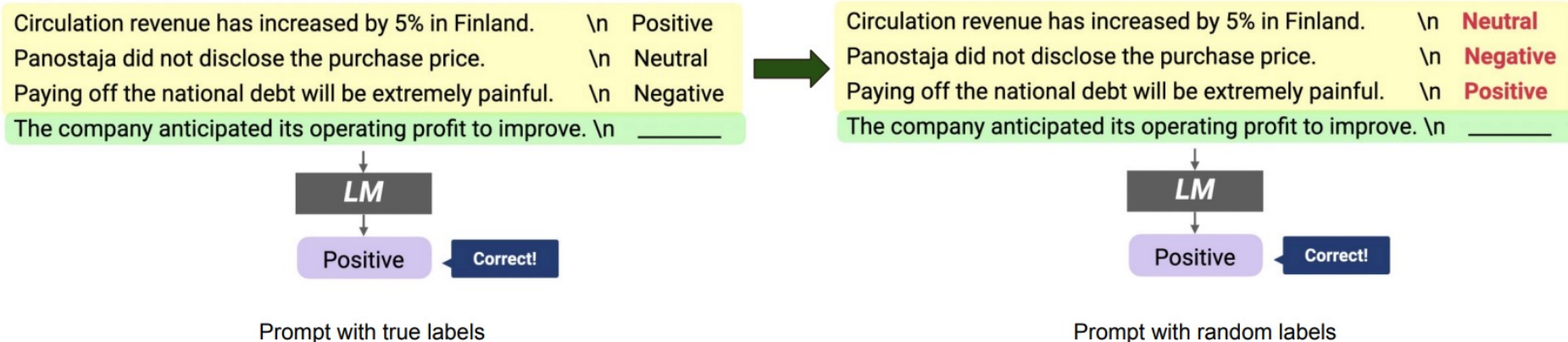
This may also indicate that, demonstrations do not teach a new task; instead, it is about locating an already-learned task during pretraining (Reynolds & McDonell, 2021)

But that brings up the question of how much LMs **actually** reason when solving these tasks. 🤔 Overlooking the impact of pretraining data can be misleading in evaluation!

In-context learning performance is highly correlated with term frequencies during pretraining

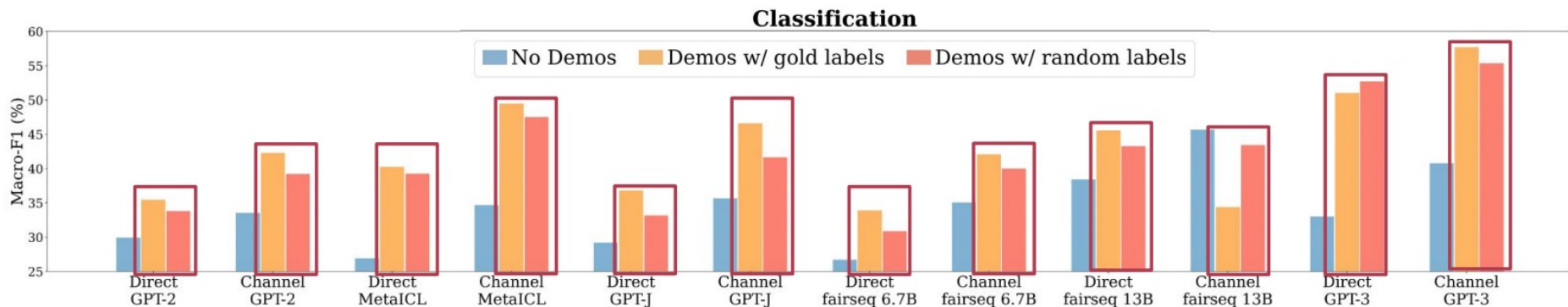
# Impact of Input-Output Mapping

- Study the effect of randomizing labels in demonstrations.
  - Randomly sample a label from the correct label space



# Impact of Input-Output Mapping

- Models see a small performance drop (0–5% absolute) with random labels



Comparisons between no-examples (blue), examples with ground truth outputs (yellow) and examples with random outputs (red)

- Takeaway:** ground truth input-label mapping in the prompt is not as important as we thought

# Impact of Input-Output Mapping

- Vary number of demonstrations
- **Takeaway:**
  - Performance drop from using gold labels to using random labels is consistently small across varying  $k$ , ranging from 0.8–1.6%
  - Using small number of examples with random labels is better than no examples

