NLP and Text-as-Data Speaker Series, Spring 2022

# Lipstick on a Pig:

#### Using Language Models as Few-Shot Learners



Sameer Singh sameersingh.org



#### Natural Language Processing Pipeline



### Natural Language Processing Pipeline



### Natural Language Processing Pipeline



#### What's next? Get rid of finetuning!



#### Manual Prompts: Sentiment Analysis Task Model (TM) Input: Amazing movie! Sentiment: [MASK] Task Input Pos (hu un us <u>N</u> Language model "Amazing movie!" LM

*P*("positive") > *P*("negative")

6

nlp

## In-Context Learning (Few-Shot Learning!)



nlp

### Why is in-context learning interesting?

#### Academically interesting

- What do language models learn? How do we control them?
- Practically relevant (with GPT-3)
  - effective with ~0-16 examples
  - serve one model for many tasks
  - no ML expertise needed
- Related to other ways of adapting language models
  - AutoPrompt\*: customized phrases to adapt LMs
  - Prompt/prefix tuning: continuous changes to input/weights
  - Increasingly more accurate and useful

### Today's Talk



## What are the biases introduced by this format?

## How robust are these capabilities to the pretraining data?

Input: Subpar acting. Sentiment: Negative Input: Beautiful film. Sentiment: Positive Input: Amazing. Sentiment: [MASK]



### Today's Talk

Z. Zhao, E. Wallace, S. Feng, D. Klein, S. Singh.

UCI nlp

**Calibrate Before Use: Improving Few-shot Performance of Language Models.** *International Conference on Machine Learning (ICML).* 2021

What are the biases introduced by this format?

How robust are these capabilities to the pretraining data?

 Image: Characteristic diagonalization
 Pre-Training

 Image: Characteristic diagonalization
 Unlabeled Data





### Components Of The Prompt

#### **Prompt Format**





### Components Of The Prompt

#### Training Example Selection

Input: <mark>Subpar acting.</mark> Sentiment: <mark>negative</mark>

Input: Beautiful film. Sentiment: positive

Input: Amazing. Sentiment:

Input: Good film.Sentiment: positiveInput: Don't watch.Sentiment: negativeInput: Amazing.Sentiment:



### Components Of The Prompt

Training Example Perturbation

Input: <mark>Subpar acting.</mark> Sentiment: <mark>negative</mark> Input: <mark>Beautiful film</mark>. Sentiment: <mark>positive</mark>

Input: Amazing. Sentiment:

Input: Beautiful film.Sentiment: positiveInput: Subpar acting.Sentiment: negativeInput: Amazing.Sentiment:

#### Accuracy Is Highly Sensitive To Prompt Design







...

Prompt #24



#### Example Permutation Impacts Accuracy

nlp

### Accuracy Is Highly Sensitive To Prompt Design



Example Selection Impacts Accuracy



#### Accuracy Is Highly Sensitive To Prompt Design





Example Format Impacts Accuracy



#### In-Context Learning

Input: Meh movie. Sentiment: Negative Input: Subpar acting. Sentiment: Negative Input: Beautiful film. Sentiment: Positive Input: Amazing. Sentiment: [MASK]





#### Majority Label Bias

Frequency of Positive *Test* Predictions



Frequent training answers dominate predictions

#### **Recency Bias**

**Frequency of Positive Predictions** 



Examples near end of prompt dominate predictions

#### Common Token Bias

			Token	Prob
The Model T was released by Fore	book	0.35		
Answer:			transportation	0.23
		)	school	0.11
	Language model LM	$  \longrightarrow$	village	0.03
			company	0.02

Token	Web (%)	Label (%)	Prediction (%)	
book	0.026	9	29	
transportation	0.000006	5 9	4	

Common n-grams dominate predictions

# Contextual Calibration of Language Models

Input: Subpar acting. Sentiment: Negative Input: Beautiful film. Sentiment: Positive Input: Amazing. Sentiment: \_\_\_\_\_



Input: Subpar acting. Sentiment: Negative Input: Beautiful film. Sentiment: Positive Input: N/A. Sentiment: \_\_\_\_\_

"meaningless" input, but full context

#### More Accurate and Stable!

11 different datasets, 0-16 shots, GPT-2 and GPT-3 models

**Different Training Examples** 

90 80 Director Accuracy (%) AGNews Accuracy (%) 0 0 0 0 08 70 60 50 MIT GPT-3 175B 40 GPT-3 13B With Calibration With Calibration 40 01 8 16 0 4 Number of Training Examples Number of Training Examples

#### Different Prompt Formats



Improved mean and worst accuracy Reduced variance for selection and ordering

Reduced variance for formats



#### Contextual Calibration for In-context Learning + *extremely* simple fix + boosts accuracy, reduces variance

- Calibration doesn't completely solve brittleness
- Independent of the pretraining corpus

### Today's Talk



## What are the biases introduced by this format?

## How robust are these capabilities to the pretraining data?

Input: Subpar acting. Sentiment: Negative Input: Beautiful film. Sentiment: Positive Input: Amazing. Sentiment: [MASK]



## Today's Talk

#### Y. Razeghi, R. Logan, M. Gardner, S. Singh.

Impact of Pretraining Term Frequencies on Few-Shot Reasoning ArXiV. 2022

What are the biases introduced by this format?

How robust are these capabilities to the pretraining data?





LM

### Reasoning and In-context Learning

- Instead of downstream classification, let's focus on Reasoning
  - Difficult to define precisely, but it's about inference
  - Go beyond regurgitation of what it has already seen
  - Feels different from memorization of facts
- Language Models need to perform reasoning

Went for a long lunch today, it lasted \_\_\_\_\_.

Alex loves chewing bones, which is not a surprise, given that he's a \_\_\_\_\_.

I wanted it in 10 days, but it took 2 weeks, which made me \_\_\_\_\_.

- And in-context few-shot reasoning is fairly accurate!
  - But how much of this performance is robust reasoning?

### Numerical Reasoning

- One of the fundamental reasoning tasks
  - Version of common-sense reasoning
- Piece of the Neural vs Symbolic debate
  - Can LMs learn to multiply numbers?
- Good few-shot performance by big LMs
  - LMs are not explicitly trained for them

Prompt
What is 75*10?
Output:
750
What is -0.002 take away 72.75?
-72.752
Calculate -0.5 - 1039.
-1039.5
What is the difference between -1360 and 2?
1362
What is -27.95 less than -20?
7.95
Calculate -0.3 + -169.
-169.3
What is 0.7 minus 0.05?
0.65
Calculate -2 + 0.0899.
-1.9101

Example from GPT-J blog:

#### UCI nlp

### Motivating Example: Multiplication

• Good performance but not always correct

Q: What is 24 times 18? A: 432 🗸

Q: What is 23 times 18? A: 462 🗙

 $\Omega(24) \simeq 10^7$ 

 $\Omega(23) \simeq 10^6$ 

Why does the model perform differently on different instances?

Hypothesis: maybe it depends on unigram statistics in pretraining?

## Motivating Example: Multiplication

- First operand: numbers between 0-99
- Accuracy averaged over:
  - 5 choices of training instances
  - second operand: numbers in 1-50

Q: What is 24 times [x]? A: \_\_\_\_ Q: What is 23 times [x]? A: \_\_\_\_ 1.0 0.8 0.6 0.6 0.4 0.2 0.2 0.0 10<sup>7</sup> 10<sup>8</sup> Frequency

Performance of GPT-J on 2-shot multiplication

### Motivating Example: Multiplication

- First operand: numbers between 0-99
- Accuracy averaged over:
  - 5 choices of training instances
  - second operands as numbers in 1-50

Q: What is 24 times [x]? A: \_\_\_\_ Q: What is 23 times [x]? A: \_\_\_\_ Performance of GPT-J on 2-shot multiplication



### Pipeline for Evaluating this Effect





nlp



### Analysis of Language Models





#### Metric: Performance Gap

• Difference in average accuracy of the instances in the top and bottom quantiles of the distribution over term frequencies

$$\Delta(\Omega) = \operatorname{Acc}(\Omega_{>90\%}) - \operatorname{Acc}(\Omega_{<10\%})$$



### Experiment Setup

- EleutherAl GPT-models
  - GPT-J-6B
  - GPT-Neo-2.7B
  - GPT-Neo-1.3B

#### **Pretrained on Pile Dataset**

- 800GB pretraining corpus
- Publicly available!

#### Training examples in the prompt:

- Randomly choose *k* examples
- 5 choice of random seeds



#### Arithmetic Reasoning



Q: What is 24 plus [x]? A: \_\_\_\_ Q: What is 24 times [x]? A: \_\_\_\_

k	Multipl	ication	Addition		
	Acc. $\Delta_1$		Acc.	$\Delta_1$	
0	5.4	18.0	1.6	8.4	
2	35.9	77.6	88.2	16.8	
4	39.2	70.8	91.4	15.0	
8	42.9	74.6	89.6	16.3	
16	40.9	73.3	88.6	16.4	

#### **Operation Inference**

1.0



Q: What is 24 # [x]?	A:
Q: What is 24 # [x]?	A:

k	Multiplie	cation (#)	Addition (#)		
	Acc.	$\Delta_1$	Acc.	$\Delta_1$	
0	-	-	-	-	
2	3.1	14.1	7.8	18.1	
4	5.7	20.9	9.8	24.8	
8	9.4	31.3	19.8	31.0	
16	11.0	39.6	26.2	38.5	

#### 37

#### Time Unit Conversion

- Minute to Seconds
- Hour to Minutes
- Day to Hour
- Week to Day
- Month to Week
- Year to Month
- Decade to Year

- Q: What is 24 minutes in seconds? A: \_\_\_\_
  - Q: What is 24 hours in minutes? A: \_\_\_\_
  - Q: What is 24 days in hours? A: \_\_\_\_
  - Q: What is 24 weeks in days? A: \_\_\_\_
  - Q: What is 24 months in weeks? A: \_\_\_\_
  - Q: What is 24 years in months? A: \_\_\_\_
  - Q: What is 24 decades in years? A: \_\_\_\_





#### Time Unit Conversion



#### Time Unit Conversion

k	Min→Sec		Hour→Min		Day→Hour		Week→Day	
	Acc.	$\Delta_{1,2}$	Acc.	$\Delta_{1,2}$	Acc.	$\Delta_{1,2}$	Acc.	$\Delta_{1,2}$
0	1.3	0.0	1.0	0.0	1.0	0.0	1.0	0.0
2	25.5	62.5	19.4	58.0	12.1	28.9	13.1	43.5
4	35.5	60.0	29.1	76.4	22.7	46.4	19.2	40.9
8	49.9	72.1	36.3	74.6	31.0	59.1	28.6	70.6
16	58.4	82.7	42.8	80.1	43.3	62.8	28.0	22.1

nlp

#### Effect of Model Size



#### multiplication

- As we increase size of model
  - Models get more accurate
  - But, more impacted by pretraining
- Number of shots is inconsistent
  - more training doesn't lead to robust reasoning by itself
- Difficult to detangle accuracy
  - By scale itself is not a solution

Effect of Pretraining on Reasoning + *high* impact on reasoning performance + raises questions about how to design, and evaluate, LMs

we are not making a causal statement about reasoning
only evaluated on numerical reasoning

### Today's Talk



## What are the biases introduced by this format?

## How robust are these capabilities to the pretraining data?

Input: Subpar acting. Sentiment: Negative Input: Beautiful film. Sentiment: Positive Input: Amazing. Sentiment: [MASK]



### What Can We Do?

- More diverse data is better!
  - Will suffer from Zipf's Law
  - Future is more unique than the past

- Augmentation during pretraining?
  - Add data to address specific reasoning
  - Good for fixing the issues we have observed
  - Doesn't feel like the end goal



 $10^6$ 

10

 $10^4$   $10^4$   $10^3$ 

 $10^{2}$  $10^{1}$ 

10

10<sup>6</sup> 10<sup>5</sup>

 $10^{2}$ 

 $10^{1}$ 

#### UCI nlp

### What Can We Do?

- Maybe scaling further will help?
  - Ultimately, they will just generalize perfectly?



- Neuro-symbolic language modeling?
  - Give LMS access to KGs, calculators, etc.
  - Barack's Wife Hillary ... [ACL 2019] \*
- Other losses for pretraining?
  - Should words really compete with each other?



\* https://arxiv.org/abs/1906.07241







# Thank you!

@sameer\_ sameer@uci.edu sameersingh.org