

In-context learning and limits I

AntNLP 2022 Fall Seminar

Li Yuqian

In-context learning

- Given $\theta^* \in \Theta$, the prompt is a concatenation of n independent demonstrations and 1 test input x_{test} that are all conditioned on θ^*
- The goal is to predict the test output y_{test} by predicting the next token.

Demonstrations

Circulation revenue has increased by 5% in Finland. \n Positive
Panostaja did not disclose the purchase price. \n Neutral
Paying off the national debt will be extremely painful. \n Negative
The acquisition will have an immediate positive impact. \n _____

Test input



Fantastically Ordered Prompts and Where to Find Them: Overcoming Few-Shot Prompt Order Sensitivity

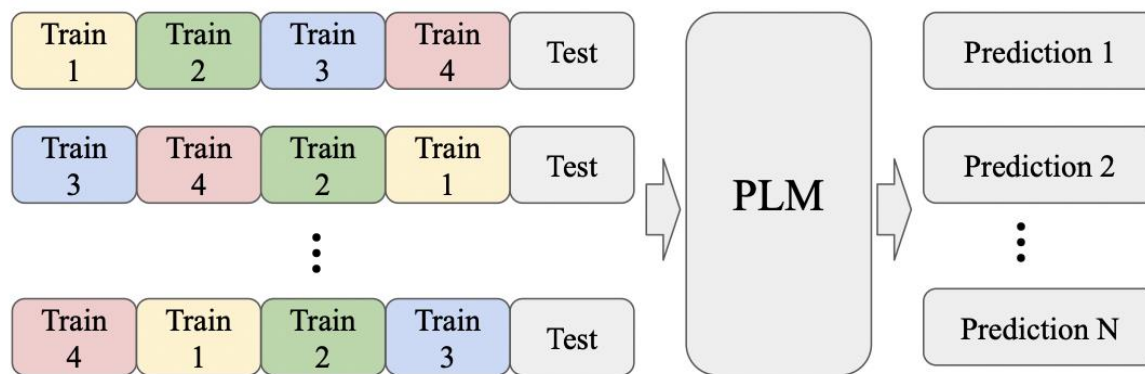
Yao Lu[†] Max Bartolo[†] Alastair Moore[‡] Sebastian Riedel[†] Pontus Stenetorp[†]

[†]University College London [‡]Mishcon de Reya LLP

`{yao.lu,m.bartolo,s.riedel,p.stenetorp}@cs.ucl.ac.uk`
`alastair.moore@mishcon.com`

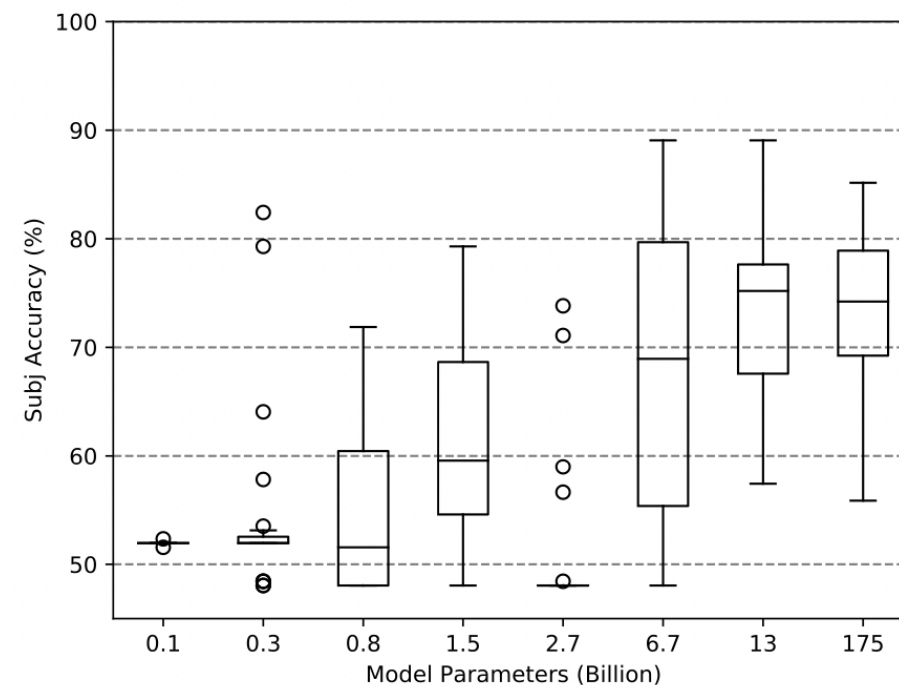
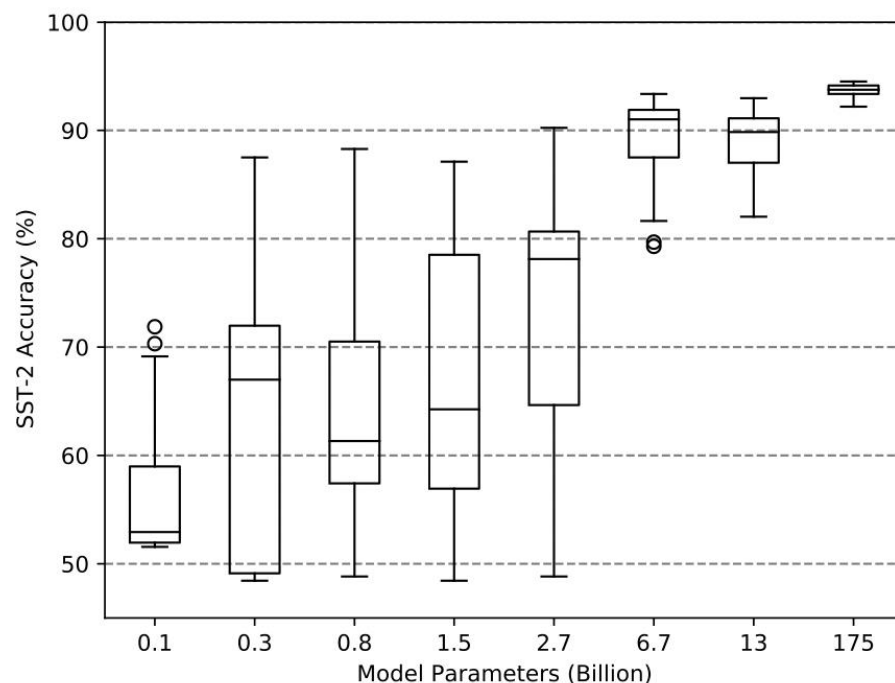
Prompt Order Sensitivity

1. Take 4 samples, create all 24 permutations, test prediction performance.



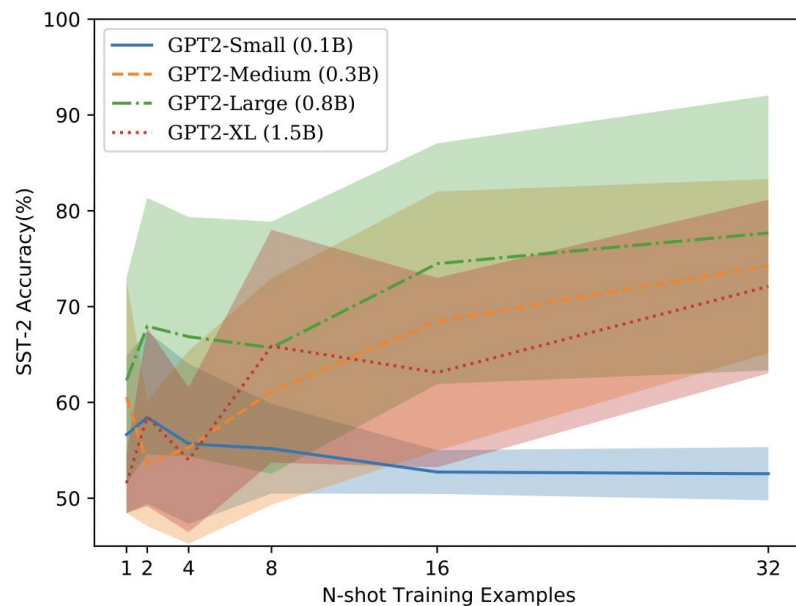
2. Test on various datasets and models (4 GPT-2, 4 GPT-3 sizes)

Prompt Design Study

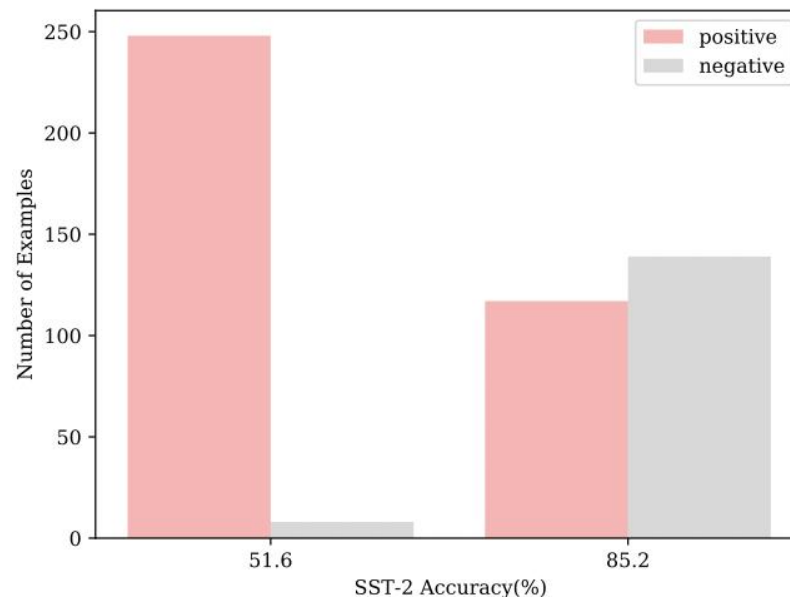


Although beneficial, increasing model size does not guarantee low variance

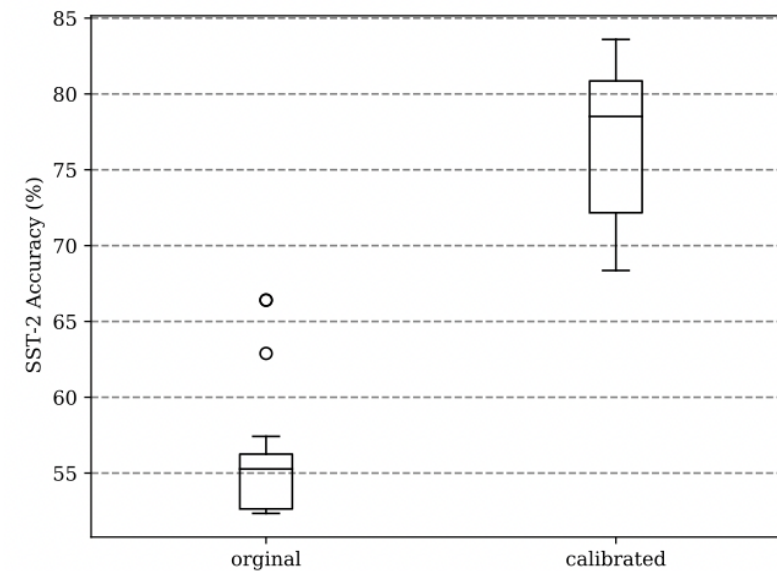
Prompt Design Study



Adding training samples does not significantly reduce variance



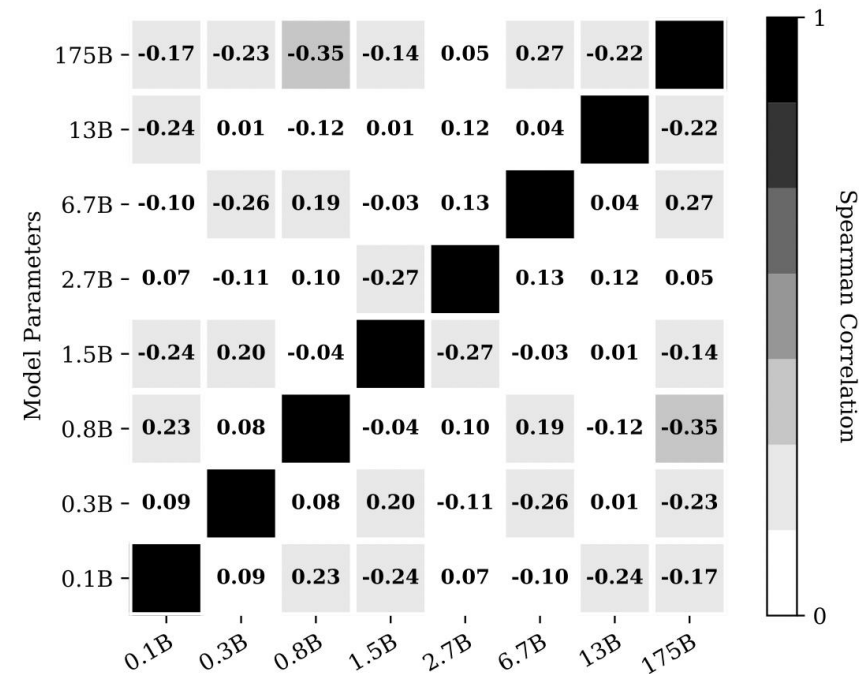
Failing prompts suffer from unbalanced label distribution



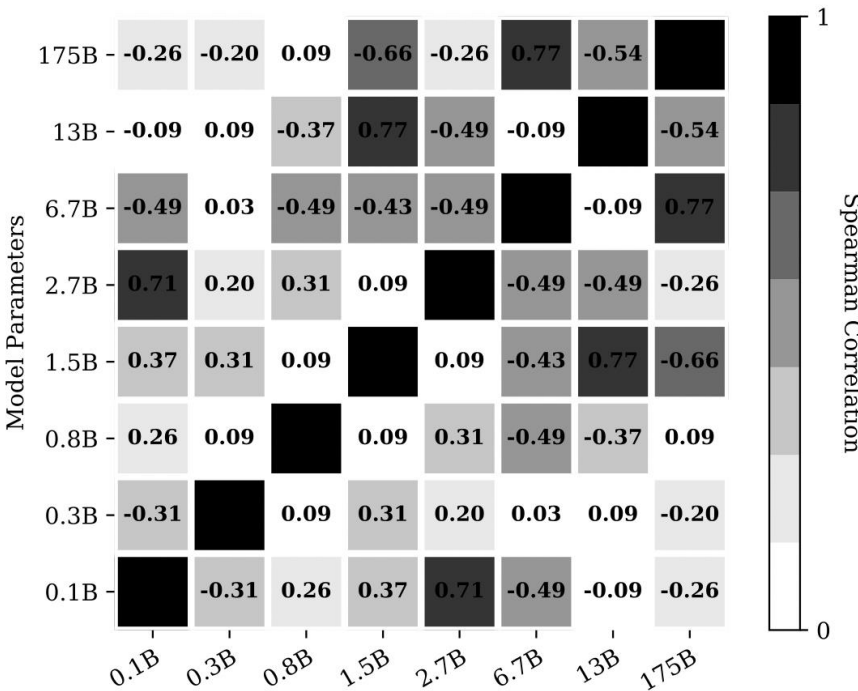
Calibration¹ improves performance but variance stays high

¹ Zhao et. al. Calibrate before use: Improving few-shot performance of language models. arXiv preprint arXiv:2102.09690.

Prompt Design Study



Performance Prompts are not transferable across models



Performant label orderings are not consistent across models

NNPP , NPNP , NPPN, PNPP , PNPN, PPNN

Spearman's rank correlation coefficient

The Spearman correlation coefficient is defined as the Pearson correlation coefficient between the rank variables.

$$r_s = \rho_{R(X), R(Y)} = \frac{\text{cov}(R(X), R(Y))}{\sigma_{R(X)} \sigma_{R(Y)}},$$

$$\rho_{X,Y} = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y} \quad \text{cov}(X, Y) = E[(X - E[X])(Y - E[Y])]$$

Only if all n ranks are *distinct integers*, it can be computed using the popular formula

$$r_s = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)},$$

where

$d_i = R(X_i) - R(Y_i)$ is the difference between the two ranks of each observation,
 n is the number of observations.

IQ, X_i	Hours of TV per week, Y_i
106	7
100	27
86	2
101	50
99	28
103	29
97	20
113	12
112	6
110	17

IQ, X_i	Hours of TV per week, Y_i	rank x_i	rank y_i	d_i	d_i^2
86	2	1	1	0	0
97	20	2	6	-4	16
99	28	3	8	-5	25
100	27	4	7	-3	9
101	50	5	10	-5	25
103	29	6	9	-3	9
106	7	7	3	4	16
110	17	8	5	3	9
112	6	9	2	7	49
113	12	10	4	6	36

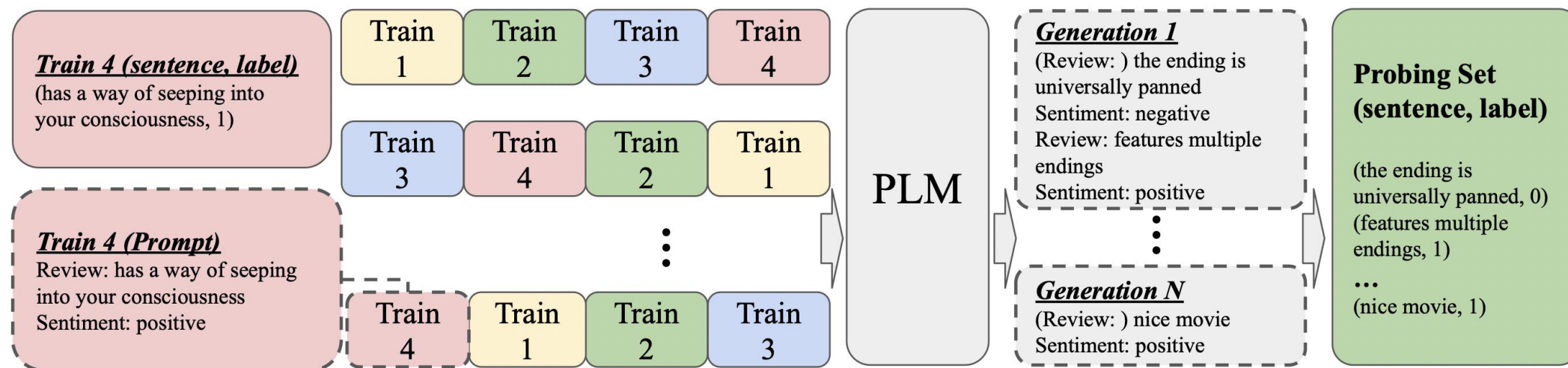
Prompt Engineering

How to find performant prompt orderings?

1. development set (need additional data)
2. automatically generate a ‘probing set’ (this paper)
 - (i) for a randomly-selected set of training samples, we use every possible ordering permutation of this set as candidates;
 - (ii) construct-ing a probing set by querying the language model using all candidate prompts as context;
 - (iii) use this probing set to identify the best ordering by ranking them using a probing metric.

Prompt Engineering

(ii) construct-ing a probing set by querying the language model using all candidate prompts as context



Probing Metrics

$$\hat{y}_{i,m} = \operatorname{argmax}_{v \in V} P(v | c_m \oplus \mathcal{T}(x'_i); \theta)$$

$$p_m^v = \frac{\sum_i \mathbb{1}_{\{\hat{y}_{i,m}=v\}}}{|D|}$$

$$\text{GlobalE}_m = \sum_{v \in V} -p_m^v \log p_m^v$$

For prompts that avoid extremely unbalanced predictions.

$$p_{i,m}^v = P_{(x'_i, y'_i) \sim D}(v | c_m \oplus \mathcal{T}(x'_i); \theta), v \in V$$

$$\text{LocalE}_m = \frac{\sum_i \sum_{v \in V} -p_{i,m}^v \log p_{i,m}^v}{|D|}$$

To penalize overconfident predictions.

Experiments

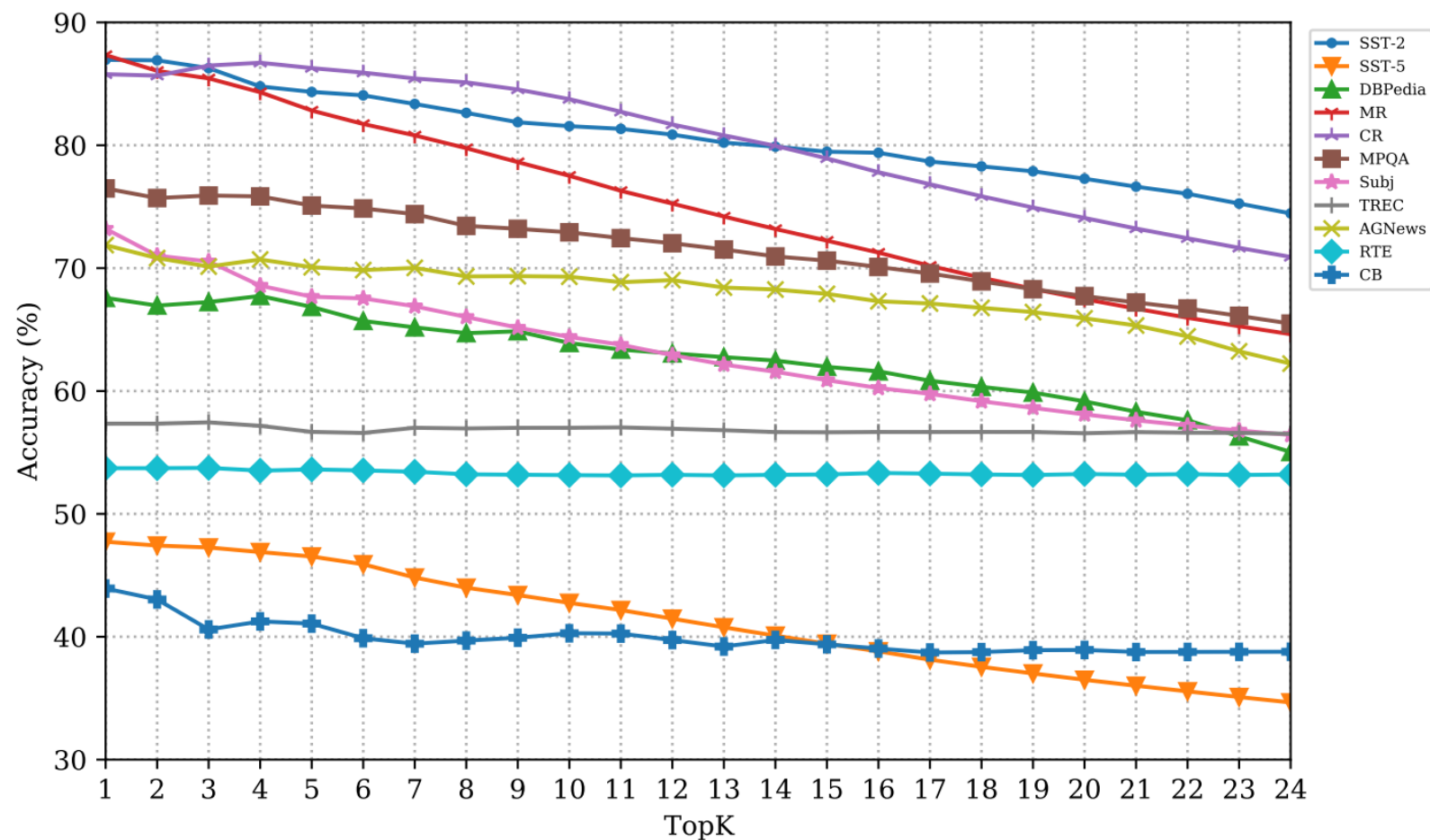
- Models: GPT-2 (with 0.1B, 0.3B, 0.8B, and 1.5B parameters), GPT-3 (with 2.7B, and 175B parameters)
- 5 different sets of randomly selected training samples; 4 training samples, 24 different permutations for each set
- select top $k=4$ samples ranked by highest entropy values
- classification datasets: SST-2, SST-5, DBPedia, MR, CR, MPQA, Subj, TREC, AGNews, RTE, CB
- Baselines:
 - Majority: predict the majority label in the dataset (lower-bound)
 - Oracle: select the top four orderings based on performance on the dev set (upper-bound)

<https://www.gwern.net/GPT-3>

	SST-2	SST-5	DBPedia	MR	CR	MPQA	Subj	TREC	AGNews	RTE	CB
Majority	50.9	23.1	9.4	50.0	50.0	50.0	50.0	18.8	25.0	52.7	51.8
Finetuning (Full)	95.0	58.7	99.3	90.8	89.4	87.8	97.0	97.4	94.7	80.9	90.5
GPT-2 0.1B	58.9 _{7.8}	29.0 _{4.9}	44.9 _{9.7}	58.6 _{7.6}	58.4 _{6.4}	68.9 _{7.1}	52.1 _{0.7}	49.2 _{4.7}	50.8 _{11.9}	49.7 _{2.7}	50.1 _{1.0}
LocalE	65.2 _{3.9}	34.4 _{3.4}	53.3 _{4.9}	66.0 _{6.3}	65.0 _{3.4}	72.5 _{6.0}	52.9 _{1.3}	48.0 _{3.9}	61.0 _{5.9}	53.0 _{3.3}	49.9 _{1.6}
GlobalE	63.8 _{5.8}	35.8 _{2.0}	56.1 _{4.3}	66.4 _{5.8}	64.8 _{2.7}	73.5 _{4.5}	53.0 _{1.3}	46.1 _{3.7}	62.1 _{5.7}	53.0 _{3.0}	50.3 _{1.6}
Oracle	73.5 _{1.7}	38.2 _{4.0}	60.5 _{4.2}	74.3 _{4.9}	70.8 _{4.4}	81.3 _{2.5}	55.2 _{1.7}	58.1 _{4.3}	70.3 _{2.8}	56.8 _{2.0}	52.1 _{1.3}
GPT-2 0.3B	61.0 _{13.2}	25.9 _{5.9}	51.7 _{7.0}	54.2 _{7.8}	56.7 _{9.4}	54.5 _{8.8}	54.4 _{7.9}	52.6 _{4.9}	47.7 _{10.6}	48.8 _{2.6}	50.2 _{5.3}
LocalE	75.3 _{4.6}	31.0 _{3.4}	47.1 _{3.7}	65.2 _{6.6}	70.9 _{6.3}	67.6 _{7.2}	66.7 _{9.3}	53.0 _{3.9}	51.2 _{7.3}	51.8 _{1.0}	47.1 _{4.2}
GlobalE	78.7 _{5.2}	31.7 _{5.2}	58.3 _{5.4}	67.0 _{5.9}	70.7 _{6.7}	68.3 _{6.9}	65.8 _{10.1}	53.3 _{4.6}	59.6 _{7.2}	51.1 _{1.9}	50.3 _{3.7}
Oracle	85.5 _{4.3}	40.5 _{6.3}	65.2 _{7.6}	74.7 _{6.1}	80.4 _{5.4}	77.3 _{2.3}	79.4 _{2.4}	63.3 _{2.9}	68.4 _{8.0}	53.9 _{1.3}	62.5 _{7.4}
GPT-2 0.8B	74.5 _{10.3}	34.7 _{8.2}	55.0 _{12.5}	64.6 _{13.1}	70.9 _{12.7}	65.5 _{8.7}	56.4 _{9.1}	56.5 _{2.7}	62.2 _{11.6}	53.2 _{2.0}	38.8 _{8.5}
LocalE	81.1 _{5.5}	40.3 _{4.7}	56.7 _{7.5}	82.6 _{4.2}	85.4 _{3.8}	73.6 _{4.8}	70.4 _{4.2}	56.2 _{1.7}	62.7 _{8.1}	53.3 _{1.6}	38.4 _{5.2}
GlobalE	84.8 _{4.1}	46.9 _{1.1}	67.7 _{3.6}	84.3 _{2.9}	86.7 _{2.5}	75.8 _{3.1}	68.6 _{6.5}	57.2 _{2.3}	70.7 _{3.6}	53.5 _{1.5}	41.2 _{4.5}
Oracle	88.9 _{1.8}	48.4 _{0.7}	72.3 _{3.3}	87.5 _{1.1}	89.9 _{0.9}	80.3 _{4.9}	76.6 _{4.1}	62.1 _{1.5}	78.1 _{1.3}	57.3 _{1.0}	53.2 _{5.3}
GPT-2 1.5B	66.8 _{10.8}	41.7 _{6.7}	82.6 _{2.5}	59.1 _{11.9}	56.9 _{9.0}	73.9 _{8.6}	59.7 _{10.4}	53.1 _{3.3}	77.6 _{7.3}	55.0 _{1.4}	53.8 _{4.7}
LocalE	76.7 _{8.2}	45.1 _{3.1}	83.8 _{1.7}	78.1 _{5.6}	71.8 _{8.0}	78.5 _{3.6}	69.7 _{5.8}	53.6 _{3.1}	79.3 _{3.7}	56.8 _{1.1}	52.6 _{3.9}
GlobalE	81.8 _{3.9}	43.5 _{4.5}	83.9 _{1.8}	77.9 _{5.7}	73.4 _{6.0}	81.4 _{2.1}	70.9 _{6.0}	55.5 _{3.0}	83.9 _{1.2}	56.3 _{1.2}	55.1 _{4.6}
Oracle	86.1 _{1.5}	50.9 _{1.0}	87.3 _{1.5}	84.0 _{2.7}	80.3 _{3.3}	85.1 _{1.4}	79.9 _{5.7}	59.0 _{2.3}	86.1 _{0.7}	58.2 _{0.6}	63.9 _{4.3}
GPT-3 2.7B	78.0 _{10.7}	35.3 _{6.9}	81.1 _{1.8}	68.0 _{12.9}	76.8 _{11.7}	66.5 _{10.3}	49.1 _{2.9}	55.3 _{4.4}	72.9 _{4.8}	48.6 _{1.9}	50.4 _{0.7}
LocalE	81.0 _{6.0}	42.3 _{4.7}	80.3 _{1.7}	75.6 _{4.1}	79.0 _{5.5}	72.5 _{5.8}	54.2 _{4.2}	54.0 _{2.6}	72.3 _{4.6}	50.4 _{1.9}	50.5 _{0.8}
GlobalE	80.2 _{4.2}	43.2 _{4.3}	81.2 _{0.9}	76.1 _{3.8}	80.3 _{3.4}	73.0 _{4.3}	54.3 _{4.0}	56.7 _{2.0}	78.1 _{1.9}	51.3 _{1.8}	51.2 _{0.8}
Oracle	89.8 _{0.7}	48.0 _{1.1}	85.4 _{1.6}	87.4 _{0.9}	90.1 _{0.7}	80.9 _{1.4}	60.3 _{10.3}	62.8 _{4.2}	81.3 _{2.9}	53.4 _{3.1}	52.5 _{1.4}
GPT-3 175B	93.9 _{0.6}	54.4 _{2.5}	95.4 _{0.9}	94.6 _{0.7}	91.0 _{1.0}	83.2 _{1.5}	71.2 _{7.3}	72.1 _{2.7}	85.1 _{1.7}	70.8 _{2.8}	75.1 _{5.1}
LocalE	93.8 _{0.5}	56.0 _{1.7}	95.5 _{0.9}	94.5 _{0.7}	91.3 _{0.5}	83.3 _{1.7}	75.0 _{4.6}	71.8 _{3.2}	85.9 _{0.7}	71.9 _{1.4}	74.6 _{4.2}
GlobalE	93.9 _{0.6}	53.2 _{2.1}	95.7 _{0.7}	94.6 _{0.2}	91.7 _{0.4}	82.0 _{0.8}	76.3 _{3.5}	73.6 _{2.5}	85.7 _{1.0}	71.8 _{1.9}	79.9 _{3.3}
Oracle	94.7 _{0.2}	58.2	96.7 _{0.2}	95.5 _{0.2}	92.6 _{0.4}	85.5 _{0.8}	81.1 _{4.9}	77.0 _{1.2}	87.7 _{0.6}	74.7 _{0.4}	83.0 _{0.9}

Results

- Ranking using Entropy-based probing is robust



Results

- Entropy-based probing is effective across templates

	Template 1	Template 2	Template 3	Template 4
GPT-2 0.1B	58.9 _{7.8}	57.5 _{6.8}	58.1 _{7.4}	56.6 _{6.6}
LocalE	65.2 _{3.9}	60.7 _{4.6}	65.4 _{4.8}	61.0 _{4.7}
GlobalE	63.8 _{5.8}	59.0 _{2.9}	64.3 _{4.8}	63.5 _{4.8}
GPT-2 0.3B	61.0 _{13.2}	63.9 _{11.3}	68.3 _{11.8}	59.2 _{6.4}
LocalE	75.3 _{4.6}	70.0 _{7.2}	80.2 _{4.2}	62.2 _{3.4}
GlobalE	78.7 _{5.2}	73.3 _{4.5}	81.3 _{4.1}	62.8 _{4.3}
GPT-2 0.8B	74.5 _{10.3}	66.6 _{10.6}	70.3 _{10.5}	63.7 _{8.9}
LocalE	81.1 _{5.5}	80.0 _{5.6}	73.7 _{6.2}	71.3 _{4.5}
GlobalE	84.8 _{4.1}	80.9 _{3.6}	79.8 _{3.9}	70.7 _{5.3}
GPT-2 1.5B	66.8 _{10.8}	80.4 _{7.6}	54.5 _{7.9}	69.1 _{10.5}
LocalE	76.7 _{8.2}	83.1 _{3.6}	66.9 _{7.5}	72.7 _{5.5}
GlobalE	81.8 _{3.9}	83.4 _{3.2}	67.2 _{6.1}	74.2 _{5.3}

ID	Template	Label Mapping
1	Review: {Sentence} Sentiment: {Label}	positive/negative
2	Input: {Sentence} Prediction: {Label}	positive/negative
3	Review: {Sentence} Sentiment: {Label}	good/bad
4	{Sentence} It was {Label}	good/bad

Table 4: Different Templates for SST-2

Prompt selection performance of different templates on SST-2

Future directions

- Linguistic perspective
 - Are there any linguistic commonalities in these good orders?
 - How do these good orders arise?
 - Does it correlate with some linguistic distributions in the pre-trained corpus?
- Mathematical perspective
 - Does the uncertainty issue really come from biased/over-confident predictions?
 - Where does the uncertainty come from? (Error of estimated distribution towards ground-truth)
 - PAC-bayes or something?

Cons - Vicky

1. Unexplored theoretical grounding | Lack of transferability
 - a. Prompt ordering affects performance greatly but is not transferable
 - i. Why does ordering matter? Why is it not transferable? Is this similar to brute-force
 - b. Probing metrics: Each motivation explained, but does not explain why only these two / how these two compare, and reason about their differing performances
2. Ablations not fully covered
 - a. Argument on template invariance: Singled out sentiment analysis that inherently has limited template formats
 - b. Lack of coverage on the 11 tasks evaluated: Pointed out sentence-pair tasks, but what about others? Complete breakdown beneficial
 - c. Argument on probing to be better than train-devel split: Is it really better than original data, or is the split unfair? (Train set cut to half, expected drop)
3. General comments
 - a. Figure captions can be improved
 - i. Fig 1: Lack of description on variation within single sample run
 - ii. Fig 3: Insufficient description on variance shade
 - iii. Fig 4/5: Insufficient description on correlation value (small = worse)
 - b. Introduce some context earlier for better grounding
 - i. Reason for choosing 4-shot (limited by window size)
 - ii. Each sample run is averaged across 5 subsets, each with 24 permutations

On the Effect of Pretraining Corpora on In-context Learning by a Large-scale Language Model

**Seongjin Shin^{*,1} Sang-Woo Lee^{*,1,2} Hwijeen Ahn¹ Sungdong Kim²
HyoungSeok Kim¹ Boseop Kim¹ Kyunghyun Cho³ Gichang Lee¹
Woomyoung Park¹ Jung-Woo Ha^{1,2} Nako Sung¹**

NAVER CLOVA¹ NAVER AI Lab² NYU³

Pretraining corpus: HyperCLOVA

Name	Description	Tokens
Blog	Blog corpus	273.6B
Cafe	Online community corpus	83.3B
News	News corpus	73.8B
Comments	Crawled comments corpus	41.1B
KiN	Korean QnA website corpus	27.3B
Modu	Collection of five datasets	6.0B
Ency	Encyclopedia corpus	1.7B
Others	Other corpus	55.0B
Total		561.8B

max training set size: 150B tokens

validation set size: 10000 examples

Table 1: Descriptions of HyperCLOVA corpus ([Kim et al., 2021](#)).

Downstream corpus: HyperCLOVA

Name	Description
NSMC	a binary sentiment classification dataset on movie review
KorQuAD	a machine reading comprehension dataset similar to SQuAD 1.0
AI Hub translation	Korean-English parallel sentences from news, government websites, legal documents, etc
YNAT	a topic classification problem with seven classes

Relationship between pretraining corpus and downstream task

KorQuAD—Ency

YMAT—News

AI Hub—News

AI Hub—KiN

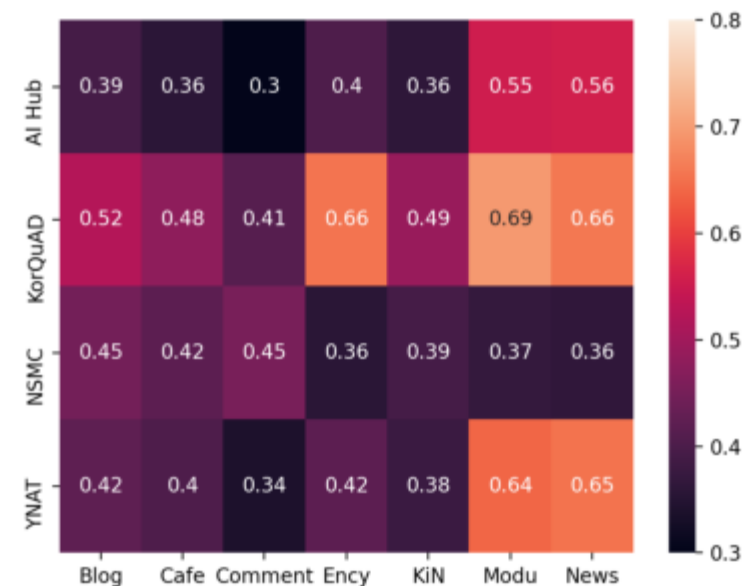


Figure 1: Vocabulary overlap ratio between pretraining corpus and downstream task. Top 1,000 nouns are used to calculate the ratio. Nouns are extracted using our in-house part-of-speech tagger.

Experiment

- 1. How large do the source and the size of pretraining corpora have the effects on emerging in-context learning ability?
- 2. What is the effect of combining various corpora?
- 3. How large does domain relevance of corpus influence on model performances of the downstream task?
- 4. How strong is the correlation between validation perplexity and in-context learning of language models?

Effect of Corpus Source

Model	Corpus Train	PPL	NSMC (Acc)	KorQuAD (EM)	KorQuAD (F1)	AI Hub (BLEU)		YNAT (F1)
						Ko→En	En→Ko	
Majority	-	-	50.35	0.0	0.0	0.0	0.0	8.26
ALL	150B	119.99	84.59	56.17	73.47	6.15	23.36	59.57
ALL w/o Others	150B	119.66	84.59	56.49	74.20	6.14	23.21	50.76
Blog	150B	152.40	83.50	50.74	69.34	3.82	20.11	60.68
Cafe	82.5B	170.85	<u>57.77</u>	<u>3.12</u>	<u>14.26</u>	<u>2.83</u>	16.53	<u>11.04</u>
News	73.1B	234.78	<u>50.72</u>	<u>0.14</u>	<u>9.96</u>	<u>1.10</u>	15.88	<u>14.36</u>
Comments	40.7B	225.39	79.78	<u>14.69</u>	<u>33.33</u>	<u>0.79</u>	<u>5.06</u>	36.17
KiN	27.0B	187.80	<u>54.73</u>	<u>4.85</u>	<u>18.99</u>	6.81	18.16	<u>9.23</u>
Modu	5.9B	226.01	69.91	30.20	49.29	<u>1.21</u>	<u>6.13</u>	43.27
Ency	1.7B	549.40	<u>53.81</u>	<u>0.71</u>	<u>11.88</u>	<u>0.58</u>	<u>0.69</u>	<u>27.99</u>
Blog 54B	54.0B	155.69	83.06	49.13	68.10	3.93	21.12	57.97
Blog 27B	27.0B	165.60	80.27	<u>10.91</u>	<u>23.41</u>	5.35	12.32	48.19
Cafe 27B	27.0B	169.81	<u>49.91</u>	<u>1.37</u>	<u>13.98</u>	4.25	20.74	<u>8.60</u>
News 27B	27.0B	239.79	<u>50.64</u>	<u>0.80</u>	<u>8.02</u>	<u>2.42</u>	15.78	<u>27.20</u>
Comments 27B	27.0B	229.65	80.50	<u>13.02</u>	<u>31.53</u>	<u>1.70</u>	<u>3.28</u>	<u>25.79</u>

In-context few-shot learning performance with different pretraining corpus

Effect of Corpus Size

# of tokens	NSMC (Acc)	KorQuAD (EM)	AI Hub (BLEU)		YNAT (F1)
			Ko→En	En→Ko	
150B	84.59	56.17	6.15	23.36	59.57
56B	84.35	55.13	5.47	22.98	51.89
6B	74.70	36.72	3.97	17.81	30.24

Table 6: In-context few-shot learning performance of ALL with different size of the pretraining data. The dataset is randomly sampled from the original corpus.

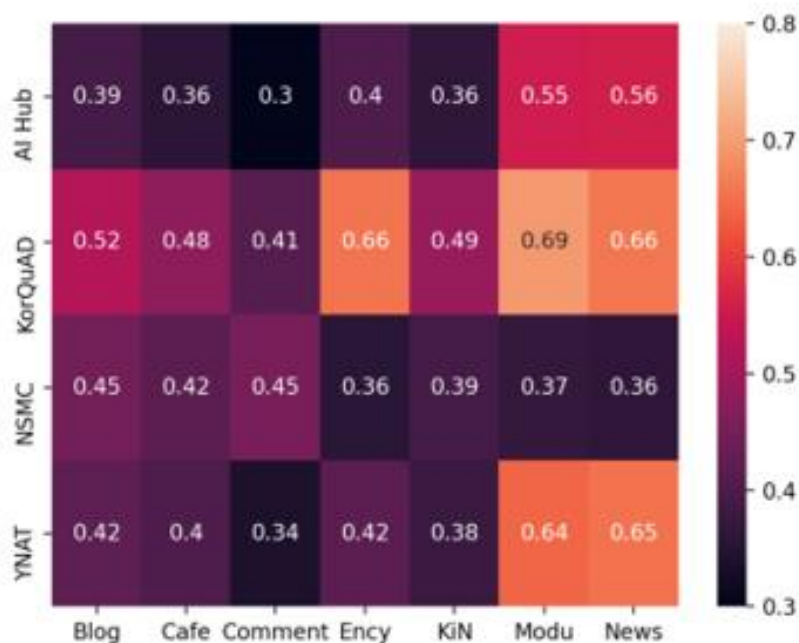
Effect of Combining Corpora

Corpus Type	Corpus Train	PPL	NSMC (Acc)	KorQuAD (EM)	KorQuAD (F1)	AI Hub (BLEU)		YNAT (F1)
ALL	150B	119.99	84.59	56.17	73.47	Ko→En	En→Ko	59.57
The Case where In-context few-shot learning Emerges by Combining Two Poor Corpora								
KiN+Ency	28.7B	164.69	59.17	42.09	61.00	8.99	23.12	42.84
Cafe+KiN	109.5B	141.92	76.42	38.45	59.00	8.41	23.41	56.96
The Case where In-context few-shot learning Does Not Emerge by Combining Two Poor Corpora								
Cafe+News	150B	154.20	54.15	8.95	22.72	4.45	17.77	8.19
The Case of Combining In-context few-shot Emerging Corpora								
Blog+Comments+Modu	150B	144.67	82.82	54.94	72.27	4.09	21.17	65.01
The Case of Adding News into KiN+Ency to Try to Enhance the Performance of YNAT								
News+KiN+Ency	101.8B	142.13	75.96	35.42	55.60	8.70	23.38	27.54

Table 4: In-context few-shot learning performance with different corpus combination.

Effect of Domain Relevance: few-shot

The close relationship between a pretraining corpus and a downstream task does not always guarantee in-context few-shot learning ability on the downstream task.



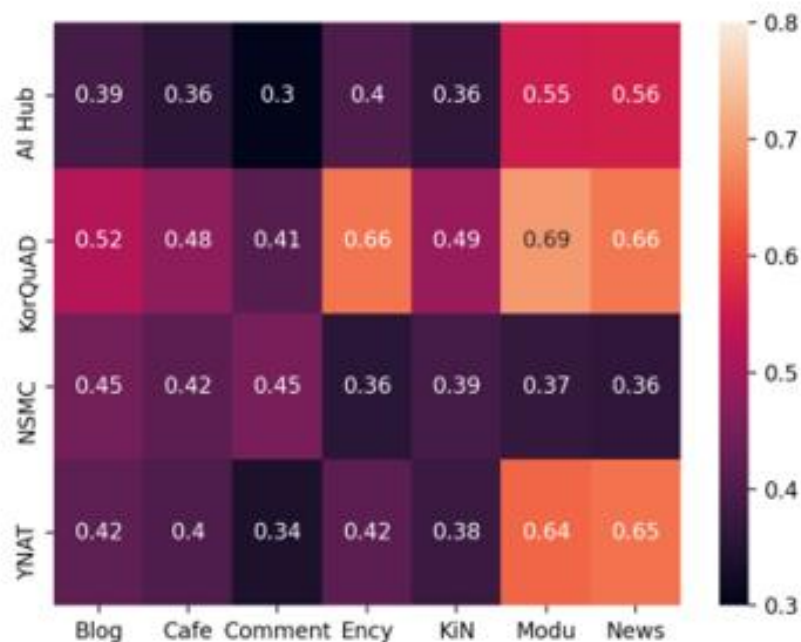
Model	Corpus Train	PPL	NSMC (Acc)	KorQuAD (EM)	KorQuAD (F1)	AI Hub (BLEU) Ko→En	AI Hub (BLEU) En→Ko	YNAT (F1)
Majority	-	-	50.35	0.0	0.0	0.0	0.0	8.26
ALL	150B	119.99	84.59	56.17	73.47	6.15	23.36	59.57
ALL w/o Others	150B	119.66	84.59	56.49	74.20	6.14	23.21	50.76
Blog	150B	152.40	83.50	50.74	69.34	3.82	20.11	60.68
Cafe	82.5B	170.85	57.77	3.12	14.26	2.83	16.53	11.04
News	73.1B	234.78	50.72	0.14	9.96	1.10	15.88	14.36
Comments	40.7B	225.39	79.78	14.69	33.33	0.79	5.06	36.17
KiN	27.0B	187.80	54.73	4.85	18.99	6.81	18.16	9.23
Modu	5.9B	226.01	69.91	30.20	49.29	1.21	6.13	43.27
Ency	1.7B	549.40	53.81	0.71	11.88	0.58	0.69	27.99
Blog 54B	54.0B	155.69	83.06	49.13	68.10	3.93	21.12	57.97
Blog 27B	27.0B	165.60	80.27	10.91	23.41	5.35	12.32	48.19
Cafe 27B	27.0B	169.81	49.91	1.37	13.98	4.25	20.74	8.60
News 27B	27.0B	239.79	50.64	0.80	8.02	2.42	15.78	27.20
Comments 27B	27.0B	229.65	80.50	13.02	31.53	1.70	3.28	25.79

Corpus Type	Corpus Train	PPL	NSMC (Acc)	KorQuAD (EM)	KorQuAD (F1)	AI Hub (BLEU) Ko→En	AI Hub (BLEU) En→Ko	YNAT (F1)
ALL	150B	119.99	84.59	56.17	73.47	6.15	23.36	59.57
The Case where In-context few-shot learning Emerges by Combining Two Poor Corpora								
KiN+Ency	28.7B	164.69	59.17	42.09	61.00	8.99	23.12	42.84
Cafe+KiN	109.5B	141.92	76.42	38.45	59.00	8.41	23.41	56.96
The Case where In-context few-shot learning Does Not Emerge by Combining Two Poor Corpora								
Cafe+News	150B	154.20	54.15	8.95	22.72	4.45	17.77	8.19
The Case of Combining In-context few-shot Emerging Corpora								
Blog+Comments+Modu	150B	144.67	82.82	54.94	72.27	4.09	21.17	65.01
The Case of Adding News into KiN+Ency to Try to Enhance the Performance of YNAT								
News+KiN+Ency	101.8B	142.13	75.96	35.42	55.60	8.70	23.38	27.54

Table 4: In-context few-shot learning performance with different corpus combination.

Effect of Domain Relevance: zero-shot

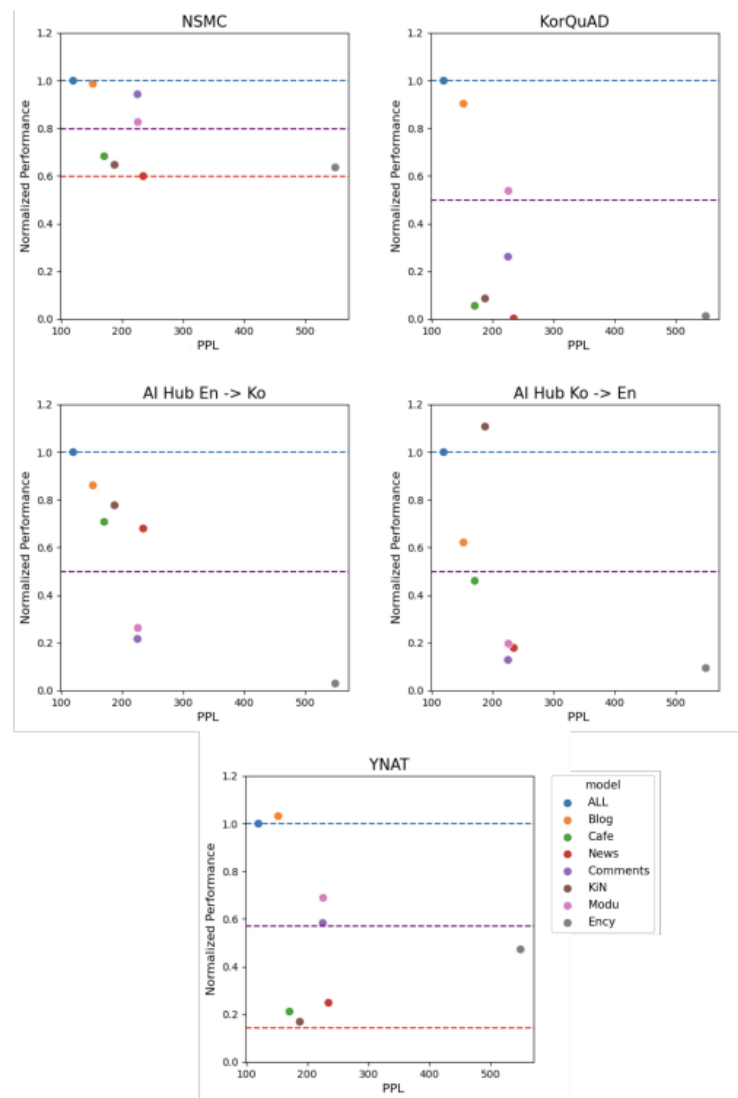
Domain relevance seems to affect more positively



Model	Corpus Train	PPL	NSMC (Acc)	KorQuAD (EM)	KorQuAD (F1)	AI Hub (BLEU)		YNAT (F1)
						Ko→En	En→Ko	
Majority	-	-	50.35	0.0	0.0	0.0	0.0	8.26
ALL	150B	119.99	<u>61.68</u>	56.17	73.47	7.43	24.81	42.79
Blog	150B	152.40	75.28	50.74	69.34	5.44	22.88	49.34
Cafe	82.5B	170.85	69.38	<u>3.12</u>	<u>14.26</u>	4.34	16.44	38.12
News	73.1B	234.78	<u>54.96</u>	<u>0.14</u>	<u>9.96</u>	<u>1.28</u>	<u>10.21</u>	48.03
Comments	40.7B	225.39	<u>57.69</u>	<u>14.69</u>	<u>33.33</u>	<u>1.98</u>	<u>3.94</u>	<u>32.48</u>
KiN	27.0B	187.80	<u>65.43</u>	<u>4.85</u>	<u>18.99</u>	4.64	<u>10.42</u>	36.06
Modu	5.9B	226.01	72.50	30.22	49.30	<u>2.39</u>	<u>7.55</u>	35.28
Ency	1.7B	549.40	<u>42.96</u>	<u>14.01</u>	<u>31.51</u>	<u>0.80</u>	<u>0.77</u>	<u>30.22</u>

Corpus Type	Corpus Train	PPL	NSMC (Acc)	KorQuAD (EM)	KorQuAD (F1)	AI Hub (BLEU)		YNAT (F1)
						Ko→En	En→Ko	
ALL	150B	119.99	<u>61.88</u>	56.17	73.47	7.43	24.81	42.79
KiN+Ency	28.7B	164.69	<u>56.78</u>	42.09	61.00	11.51	24.93	37.71
Cafe+KiN	109.5B	141.92	<u>59.27</u>	38.45	59.00	10.12	24.95	45.44
Cafe+News	150B	154.20	<u>66.92</u>	<u>8.95</u>	<u>22.85</u>	3.49	15.77	47.34
Blog+Comments+Modu	150B	144.67	69.15	54.94	72.27	6.06	22.03	48.25
News+KiN+Ency	101.8B	142.13	<u>61.49</u>	35.42	55.60	10.18	24.13	51.89

Perplexity and Downstream Task



In-context few-shot learning performance of various corpus models and their PPL.

Conclusion

- Corpus Source: In-context learning performance depends heavily on corpus sources, and with some sources, in-context learning does not work effectively.
- Corpus Combination: In-context learning ability can emerge by fusing two corpora, even when each on its own does not result in in-context learning.
- Domain Relevance: Pretraining with a corpus related to a downstream task seems to help in-context zero-shot learning performance, but is not indicative of the competitive in-context few-shot learning performance.
- Perplexity: Although perplexity and in-context learning accuracies correlate well when training a single model, perplexity alone does not reflect the difference in in-context learning accuracies across different language models.

Rethinking the Role of Demonstrations: What Makes In-Context Learning Work?

**Sewon Min^{1,2} Xinxu Lyu¹ Ari Holtzman¹ Mikel Artetxe²
Mike Lewis² Hannaneh Hajishirzi^{1,3} Luke Zettlemoyer^{1,2}**

¹University of Washington ²Meta AI ³Allen Institute for AI
`{sewon, alrope, ahai, hannaneh, lsz}@cs.washington.edu`
`{artetxe, mikelewis}@meta.com`

- 1. The input-label mapping, i.e., whether each input x_i is paired with a correct label y_i .
- 2. The distribution of the input text, i.e., the underlying distribution that $x_1 \dots x_k$ are from.
- 3. The label space, i.e., the space covered by $y_1 \dots y_k$.
- 4. The format—specifically, the use of input-label pairing as the format

Model

Model	# Params	Public	Meta-trained
GPT-2 Large	774M	✓	✗
MetaICL	774M	✓	✓
GPT-J	6B	✓	✗
fairseq 6.7B [†]	6.7B	✓	✗
fairseq 13B [†]	13B	✓	✗
GPT-3	175B [‡]	✗	✗

MetalCL

	Meta-training	Inference
Task	C meta-training tasks	An unseen <i>target</i> task
Data given	Training examples $\mathcal{T}_i = \{(x_j^i, y_j^i)\}_{j=1}^{N_i}, \forall i \in [1, C]$ ($N_i \gg k$)	Training examples $(x_1, y_1), \dots, (x_k, y_k)$, Test input x
Objective	For each iteration, 1. Sample task $i \in [1, C]$ 2. Sample $k + 1$ examples from \mathcal{T}_i : $(x_1, y_1), \dots, (x_{k+1}, y_{k+1})$ 3. Maximize $P(y_{k+1} x_1, y_1, \dots, x_k, y_k, x_{k+1})$	$\operatorname{argmax}_{c \in \mathcal{C}} P(c x_1, y_1, \dots, x_k, y_k, x)$

Direct vs Channel Models

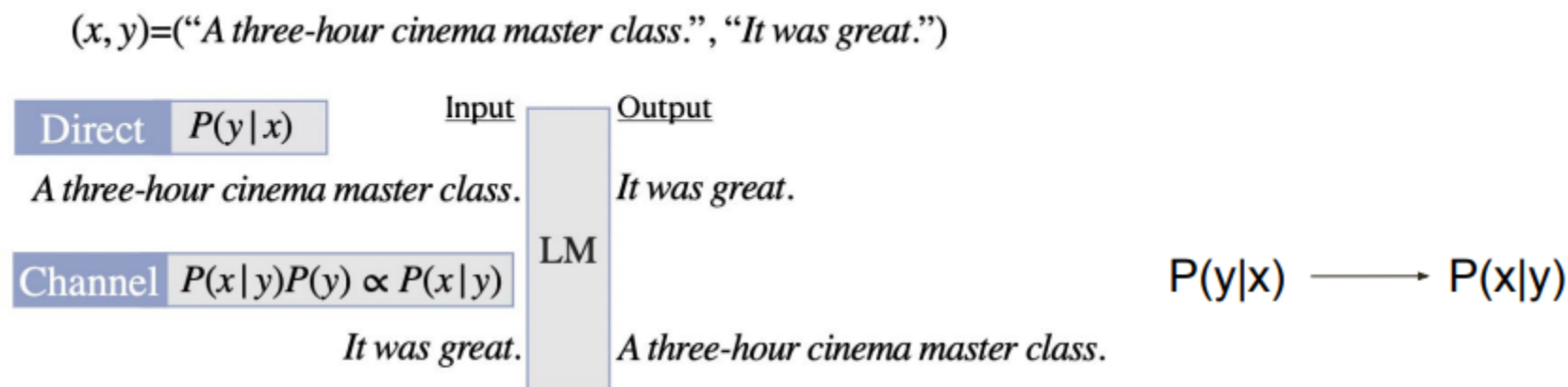


Figure 1: An illustration of the direct model and the channel model for language model prompting in the sentiment analysis task.

Evaluation Data

Classification Tasks

Dataset	# Train	# Eval
<i>Task category: Sentiment analysis</i>		
financial_phrasebank	1,811	453
poem_sentiment	892	105
<i>Task category: Paraphrase detection</i>		
medical_questions_pairs	2,438	610
glue-mrpc	3,668	408
<i>Task category: Natural language inference</i>		
glue-wnli	635	71
climate_fever	1,228	307
glue-rte	2,490	277
superglue-cb	250	56
sick	4,439	495
<i>Task category: Hate speech detection</i>		
hate_speech18	8,562	2,141
ethos-national_origin	346	87
ethos-race	346	87
ethos-religion	346	87
tweet_eval-hate	8,993	999
tweet_eval-stance_atheism	461	52
tweet_eval-stance_feminist	597	67
<i>Task category: Question answering</i>		
quarel	1,941	278
openbookqa	4,957	500
qasc	8,134	926
commonsense_qa	9,741	1,221
ai2_arc	1,119	299
<i>Task category: Sentence completion</i>		
codah	1665	556
superglue-copa	400	100
dream	6116	2040
quartz-with_knowledge	2696	384
quartz-no_knowledge	2696	384

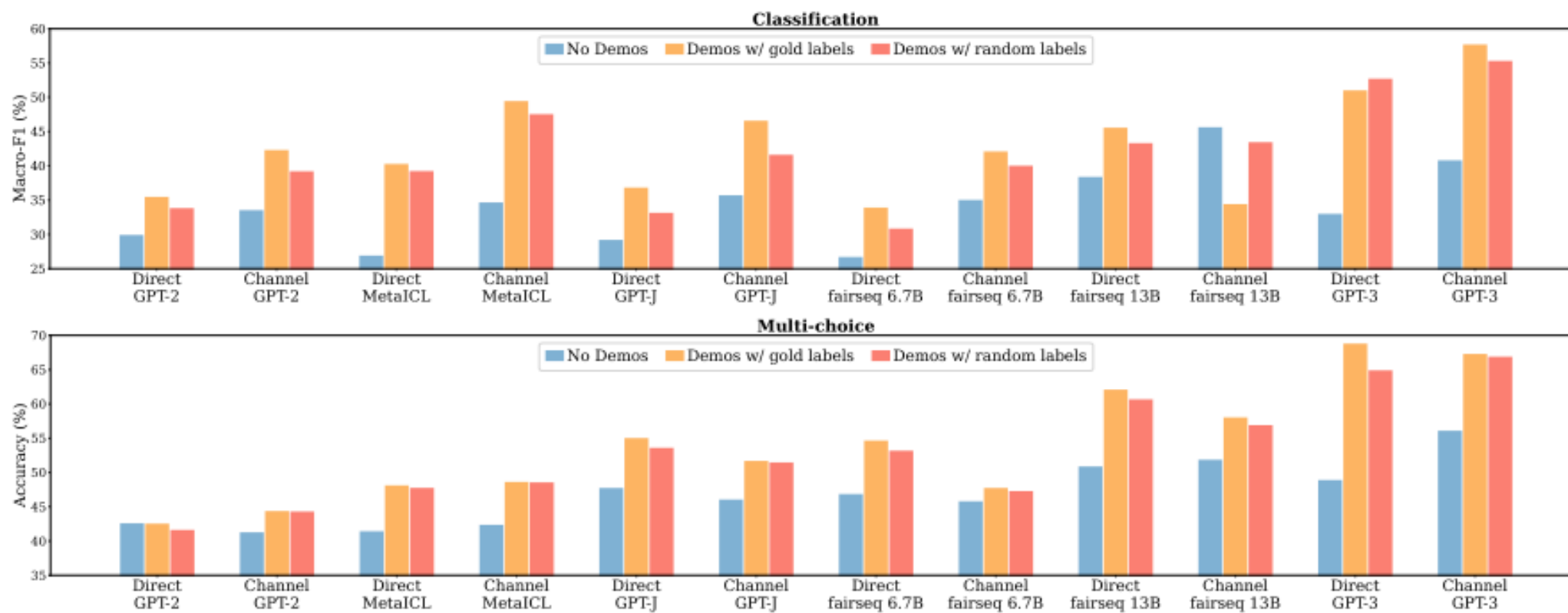
Multiple Choice Tasks

Some Details

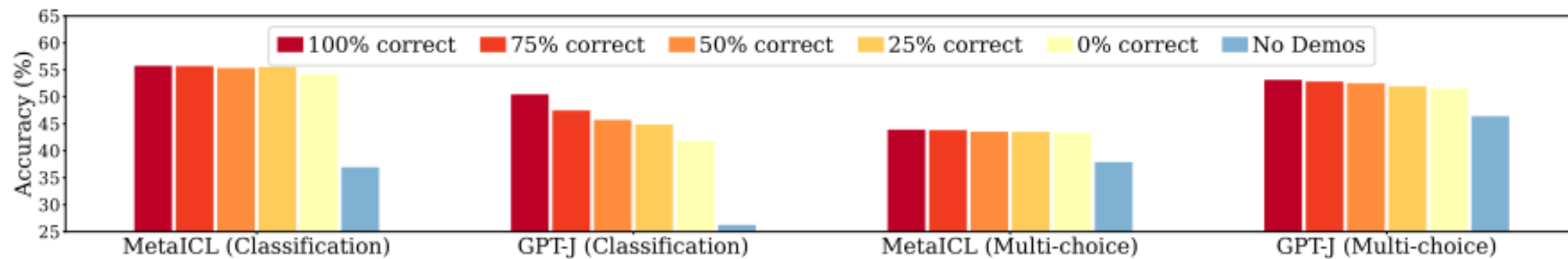
- $k = 16$ examples as demonstrations
- report Macro-F1 for classification tasks and Accuracy for multi-choice tasks
- minimal templates

Dataset	Type	Example
MRPC	Minimal	sentence 1: Cisco pared spending to compensate for sluggish sales . [SEP] sentence 2: In response to sluggish sales , Cisco pared spending . \n {equivalent not_equivalent}
	Manual	Cisco pared spending to compensate for sluggish sales . \n The question is: In response to sluggish sales , Cisco pared spending . True or False? \n The answer is:{True False}
RTE	Minimal	sentence 1: The girl was found in Drummondville. [SEP] sentence 2: Drummondville contains the girl. \n {entailment not_entailment}
	Manual	The girl was found in Drummondville. \n The question is: Drummondville contains the girl. True or False? \n The answer is:{True False}
Tweet_eval-hate	Minimal	The Truth about #Immigration \n {hate non-hate}
	Manual	Tweet: The Truth about #Immigration \n Sentiment: {against favor}
SICK	Minimal	sentence 1: A man is screaming. [SEP] sentence 2: A man is scared. \n {contradiction entailment neutral}
	Manual	A man is screaming. \n The question is: A man is scared. True or False? \n The answer is: {False True Not sure}
poem-sentiment	Minimal	willis sneered: \n {negative no_impact positive}
	Manual	willis sneered: \n The sentiment is: {negative no_impact positive}
OpenbookQA	Minimal	What creates a valley? \n {feet rock water sand}
	Manual	The question is: What creates a valley? \n The answer is: {feet rock water sand}
CommonsenseQA	Minimal	What blocks sunshine? \n {summer park desktop sea moon}
	Manual	The question is: What blocks sunshine? \n The answer is: {summer park desktop sea moon}
COPA	Minimal	Effect: I coughed. \n {Cause: I inhaled smoke. Cause: I lowered my voice.}
	Manual	I coughed because {I inhaled smoke. I lowered my voice.}
ARC	Minimal	Which biome has the most vegetation? \n {desert forest grassland tundra}
	Manual	The question is: Which biome has the most vegetation? \n The answer is: {desert forest grassland tundra}

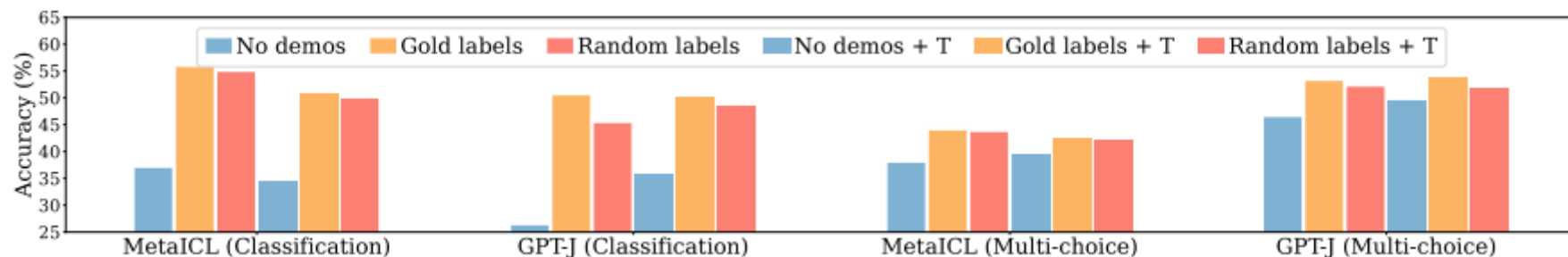
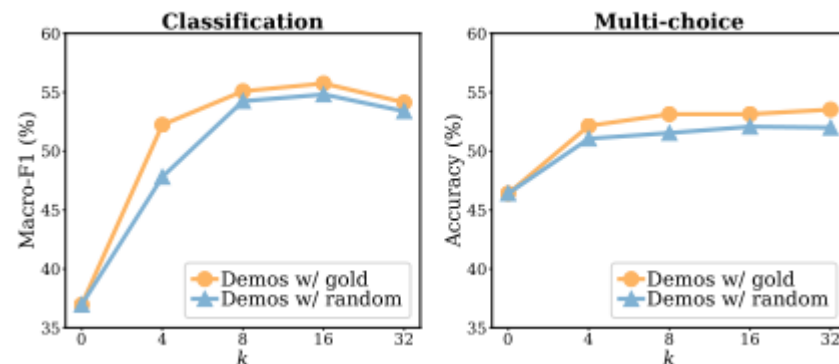
Experiment: Input-label Mapping



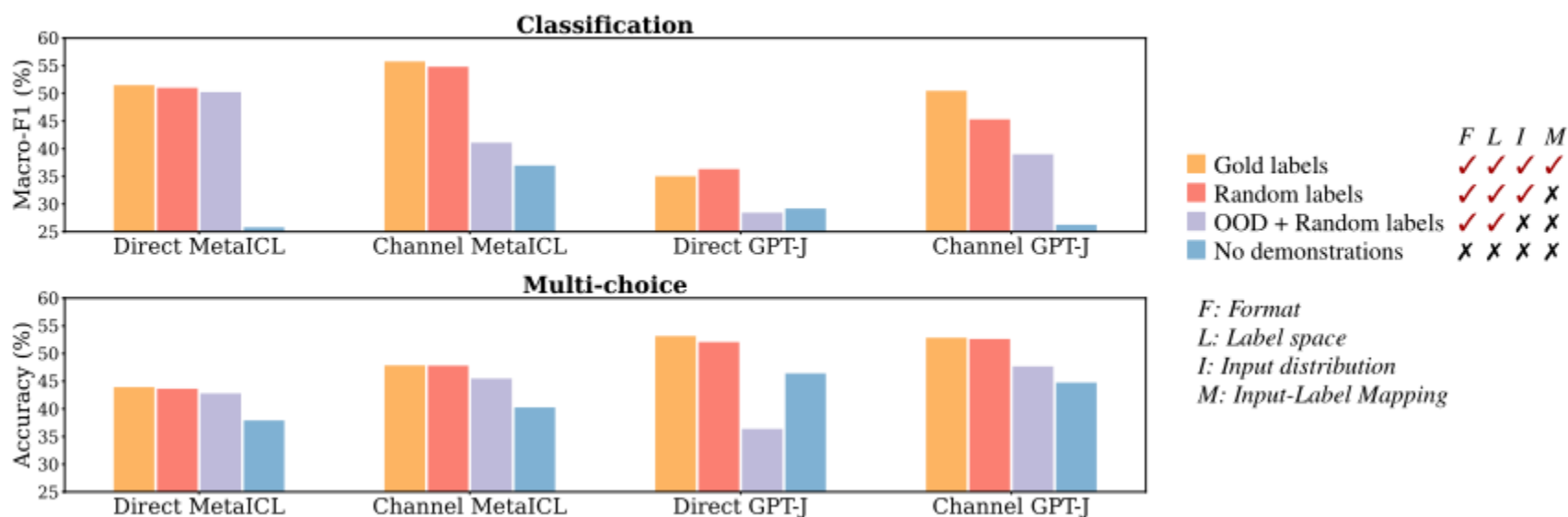
Experiment: Input-label Mapping



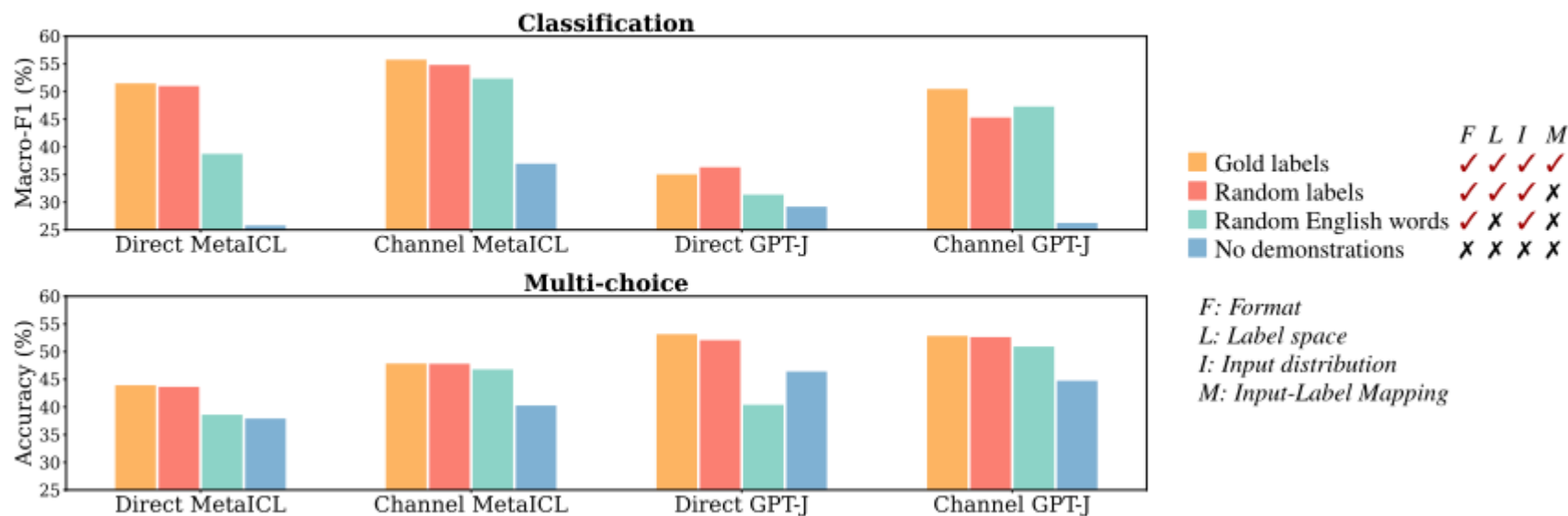
Experiment: Input-label Mapping



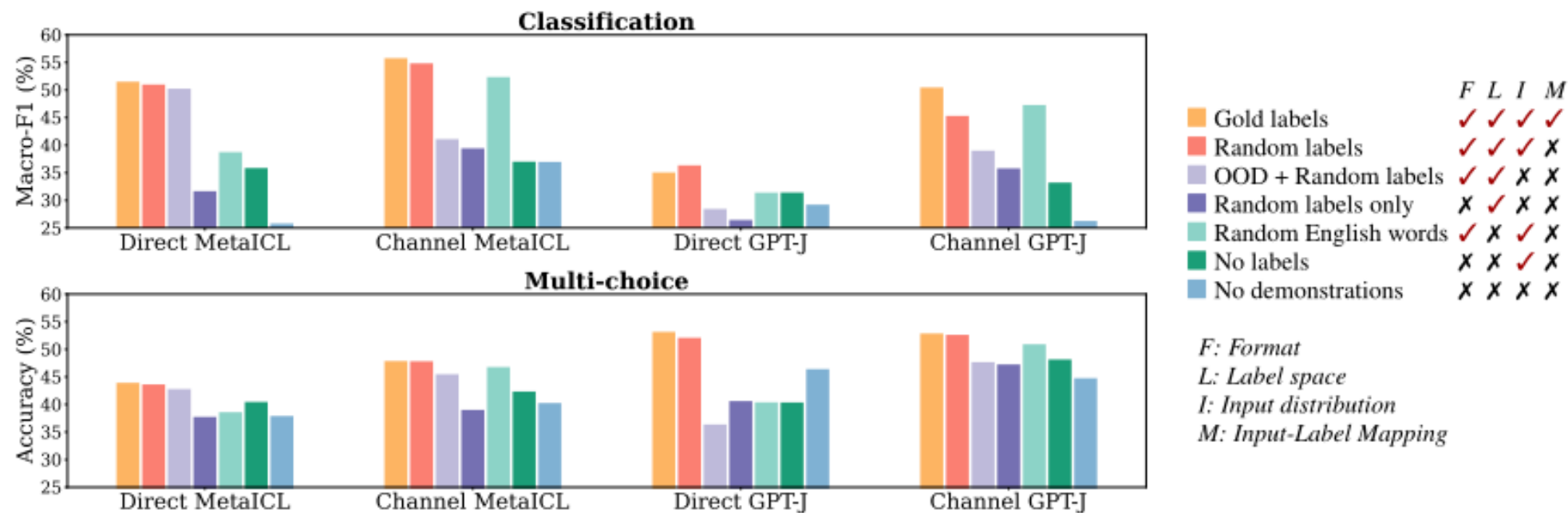
Experiment: Impact of the distribution of the input text



Experiment: Impact of the label space



Experiment: Impact of the format



Experiment: meta-training

the ground truth input-label mapping matters even less

keeping the format of the demonstrations matters even more

nearly zero influence of the input-label mapping and the input distribution in Direct MetalCL

nearly zero influence of the input-label mapping and the output space in Channel MetalCL

(1) the input-label mapping is likely harder to exploit,

(2) the format is likely easier to exploit, and

(3) the space of the text that the model is trained to generate is likely easier to exploit than the space of the text that the model conditions on.