

# Domain-Specific NLP

# Contents

1. Domain-Specific Models
  - a. *Don't Stop Pre-training*
  - b. Specialized Models (BioBERT, SciBERT, Galactica)
2. Unsupervised Classification Models
  - a. Out-of-the-box representations: limitations
  - b. SimCSE, E5, GTE...
  - c. Document Representation: DocBERT
3. Learning Long-Range Dependencies
  - a. Long-range attention models
  - b. State-space models: S4

# Domain-Specific Models

# Domain-Specific Models

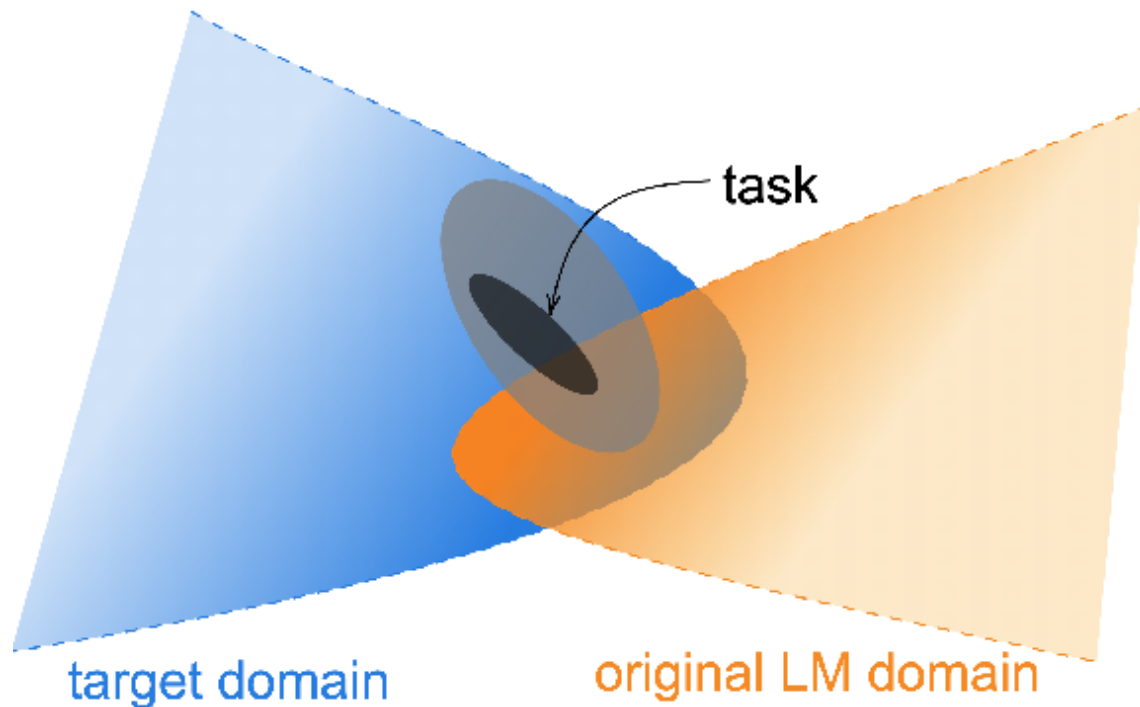
**Pretrained (Large) Language Models** are trained on content crawled over the internet, books, reports and news papers and are, hence **are open-domain**.

A **textual domain** is the **distribution over language characterizing a given topic or genre** [1].

- You are more likely to see the word "integer" in computer science than in news papers.
- An (L)LM will be more perplex to the word "integer" even though the input comes from a StackOverflow post.

# Don't Stop Pretraining

*Don't Stop Pretraining: Adapt Language Models to Domains and Tasks. [1]*



# Don't Stop Pretraining

Domain	Pretraining Corpus	# Tokens	Size	$\mathcal{L}_{\text{ROB.}}$	$\mathcal{L}_{\text{DAPT}}$
BIOMED	2.68M full-text papers from S2ORC (Lo et al., 2020)	7.55B	47GB	1.32	0.99
CS	2.22M full-text papers from S2ORC (Lo et al., 2020)	8.10B	48GB	1.63	1.34
NEWS	11.90M articles from REALNEWS (Zellers et al., 2019)	6.66B	39GB	1.08	1.16
REVIEWS	24.75M AMAZON reviews (He and McAuley, 2016)	2.11B	11GB	2.10	1.93
ROBERTA (baseline)	see Appendix §A.1	N/A	160GB	‡1.19	-

**Table 1:** List of the domain-specific unlabeled datasets. In columns 5 and 6, we report ROBERTA’s masked LM loss on 50K randomly sampled held-out documents from each domain before ( $\mathcal{L}_{\text{ROB.}}$ ) and after ( $\mathcal{L}_{\text{DAPT}}$ ) *DAPT* (lower implies a better fit on the sample). ‡ indicates that the masked LM loss is estimated on data sampled from sources similar to ROBERTA’s pretraining corpus.

# Don't Stop Pretraining

PT	100.0	54.1	34.5	27.3	19.2
News	54.1	100.0	40.0	24.9	17.3
Reviews	34.5	40.0	100.0	18.3	12.7
BioMed	27.3	24.9	18.3	100.0	21.4
CS	19.2	17.3	12.7	21.4	100.0
	PT	News	Reviews	BioMed	CS

**Figure 2:** Vocabulary overlap (%) between domains. PT denotes a sample from sources similar to ROBERTA's pretraining corpus. Vocabularies for each domain are created by considering the top 10K most frequent words (excluding stopwords) in documents sampled from each domain.

# Don't Stop Pretraining

Domain	Task	RoBERTa	Additional Pretraining Phases		
			DAPT	TAPT	DAPT + TAPT
BioMed	CHEMPROT	81.9 <sub>1.0</sub>	84.2 <sub>0.2</sub>	82.6 <sub>0.4</sub>	<b>84.4</b> <sub>0.4</sub>
	†RCT	87.2 <sub>0.1</sub>	87.6 <sub>0.1</sub>	87.7 <sub>0.1</sub>	<b>87.8</b> <sub>0.1</sub>
CS	ACL-ARC	63.0 <sub>5.8</sub>	75.4 <sub>2.5</sub>	67.4 <sub>1.8</sub>	<b>75.6</b> <sub>3.8</sub>
	SciERC	77.3 <sub>1.9</sub>	80.8 <sub>1.5</sub>	79.3 <sub>1.5</sub>	<b>81.3</b> <sub>1.8</sub>
NEWS	HYPERPARTISAN	86.6 <sub>0.9</sub>	88.2 <sub>5.9</sub>	<b>90.4</b> <sub>5.2</sub>	90.0 <sub>6.6</sub>
	†AGNEWS	93.9 <sub>0.2</sub>	93.9 <sub>0.2</sub>	94.5 <sub>0.1</sub>	<b>94.6</b> <sub>0.1</sub>
REVIEWS	†HELPFULNESS	65.1 <sub>3.4</sub>	66.5 <sub>1.4</sub>	68.5 <sub>1.9</sub>	<b>68.7</b> <sub>1.8</sub>
	†IMDB	95.0 <sub>0.2</sub>	95.4 <sub>0.1</sub>	95.5 <sub>0.1</sub>	<b>95.6</b> <sub>0.1</sub>

**Table 5:** Results on different phases of adaptive pretraining compared to the baseline RoBERTa (col. 1). Our approaches are *DAPT* (col. 2, §3), *TAPT* (col. 3, §4), and a combination of both (col. 4).



# Don't Stop Pretraining

"We show that **pretraining the model towards a specific task or small corpus can provide significant benefits**. Our findings suggest it may be valuable to complement work on ever-larger LMs with parallel efforts to **identify and use domain and task relevant corpora to specialize models**."

# BioBERT

"[..] the word distributions of general and biomedical corpora are quite different, which can often be a problem for biomedical text mining models." [2]

# BioBERT

: 1. List of text corpora used for BioBERT

Corpus	# of words (B)	Domain
English Wikipedia	2.5B	General
BooksCorpus	0.8B	General
PubMed Abstracts	4.5B	Biomedical
PMC Full-text articles	13.5B	Biomedical

**Table 1.** List of text corpora used for BioBERT

# BioBERT

"We showed that **pre-training BERT on biomedical corpora is crucial in applying it to the biomedical domain**. Requiring minimal task-specific architectural modification, **BioBERT outperforms previous models on biomedical text mining tasks** such as NER, RE and QA."

# SciBERT

"[...] while both BERT and ELMo have released pretrained models, they are still trained on general domain corpora such as news articles and Wikipedia." [3]

# SciBERT

Field	Task	Dataset	SOTA	BERT-Base		SciBERT	
				Frozen	Finetune	Frozen	Finetune
Bio	NER	BC5CDR (Li et al., 2016)	88.85 <sup>7</sup>	85.08	86.72	88.73	<b>90.01</b>
		JNLPBA (Collier and Kim, 2004)	<b>78.58</b>	74.05	76.09	75.77	77.28
		NCBI-disease (Dogan et al., 2014)	<b>89.36</b>	84.06	86.88	86.39	88.57
	PICO	EBM-NLP (Nye et al., 2018)	66.30	61.44	71.53	68.30	<b>72.28</b>
	DEP	GENIA (Kim et al., 2003) - LAS	<b>91.92</b>	90.22	90.33	90.36	90.43
		GENIA (Kim et al., 2003) - UAS	<b>92.84</b>	91.84	91.89	92.00	91.99
	REL	ChemProt (Kringelum et al., 2016)	76.68	68.21	79.14	75.03	<b>83.64</b>
CS	NER	SciERC (Luan et al., 2018)	64.20	63.58	65.24	65.77	<b>67.57</b>
	REL	SciERC (Luan et al., 2018)	n/a	72.74	78.71	75.25	<b>79.97</b>
	CLS	ACL-ARC (Jurgens et al., 2018)	67.9	62.04	63.91	60.74	<b>70.98</b>
Multi	CLS	Paper Field	n/a	63.64	65.37	64.38	<b>65.71</b>
		SciCite (Cohan et al., 2019)	84.0	84.31	84.85	<b>85.42</b>	<b>85.49</b>
Average				73.58	77.16	76.01	79.27

**Table 1:** Test performances of all BERT variants on all tasks and datasets. [...]

# SciBERT

Task	Dataset	BIOBERT	SCIERT
NER	BC5CDR	88.85	90.01
	JNLPBA	77.59	77.28
	NCBI-disease	89.36	88.57
REL	ChemProt	76.68	83.64

**Table 2:** Comparing SciBERT with the reported BioBERT results on biomedical datasets.

# SciBERT

NB: SciBERT was trained on curated textual data ; not trained on code or script for example --at least not trained directly and purposefully on this kind of data

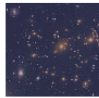
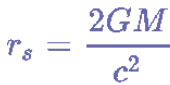
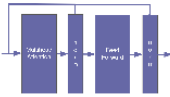
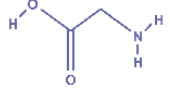

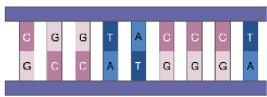


# Galactica

"Computing has indeed revolutionized how research is conducted, but information overload remains an overwhelming problem [...]. In this paper, we argue for a better way through large language models. Unlike search engines, language models can potentially store, combine and reason about scientific knowledge." [4]

- Galactica was trained on a rather small highly curated dataset.
- All the data was standardized as markdown text.

# Galactica

Modality	Entity	Sequence	
Text	Abell 370	Abell 370 is a cluster...	
LaTeX	Schwarzschild radius	$r_s = \frac{2GM}{c^2}$	
Code	Transformer	<code>class Transformer(nn.Module)</code>	
SMILES	Glycine	<chem>C(C(=O)O)N</chem>	
AA Sequence	Collagen $\alpha$ -1(II) chain	MIRLGAPQTL..	
DNA Sequence	Human genome	CGGTACCCTC..	

**Table 1: Tokenizing Nature.** Galactica trains on text sequences that represent scientific phenomena.

**Table 1:** Tokenizing Nature. Galactica trains on text sequences that represent scientific phenomena.

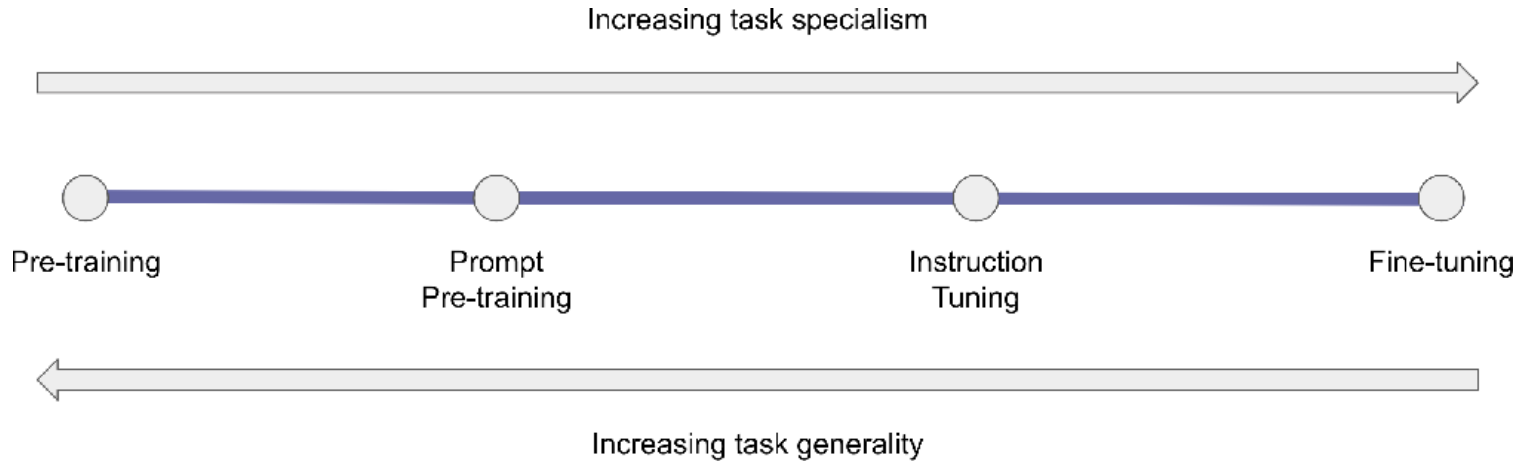
# Galactica

1. **Citations:** we wrap citations with special reference tokens [START\_REF] and [END\_REF].
2. **Step-by-Step Reasoning:** we wrap step-by-step reasoning with a working memory token , mimicking an internal working memory context.
3. **Mathematics:** for mathematical content, with or without LaTeX, we split ASCII operations into individual characters. Parentheses are treated like digits. The rest of the operations allow for unsplit repetitions. Operation characters are !"#\$%&'\*+,-./:;<=>?^\_`| and parentheses are ()[]{}.

4. **Numbers:** we split digits into individual tokens. For example 737612.62 -> 7,3,7,6,1,2,,6,2.
5. **SMILES formula:** we wrap sequences with [START\_SMILES] and [END\_SMILES] and apply characterbased tokenization. Similarly we use [START\_I\_SMILES] and [END\_I\_SMILES] where isomeric SMILES is denoted. For example, C(C(=O)O)N → C,(,C,(,=,O,),O,),N.
6. **Amino acid sequences:** we wrap sequences with [START\_AMINO] and [END\_AMINO] and apply character-based tokenization, treating each amino acid character as a single token. For example, MIRLGAPQTL -> M,I,R,L,G,A,P,Q,T,L.

7. **DNA sequences:** we also apply a character-based tokenization, treating each nucleotide base as a token, where the start tokens are [START\_DNA] and [END\_DNA]. For example, CGGTACCCTC -> C, G, G, T, A, C, C, C, T, C.

# Galactica



**Figure 5:** Prompt Pre-training. Pre-training weighs all tokens equally as part of the self-supervised loss. This leads to a weak relative signal for tasks of interest, meaning model scale has to be large to work. Instruction tuning boosts performance post hoc, and can generalize to unseen tasks of interest, but it risks performance in tasks that are distant from instruction set tasks.

Prompt pre-training has a weaker task of interest bias than instruction tuning but less risk of

# Galactica

- **GeLU Activation** - GeLU activations for all model sizes.
- **Context Window** - a 2048 length context window.
- **No Biases** - following PaLM, we do not use biases in any of the dense kernels or layer norms.
- **Learned Positional Embeddings** - learned positional embeddings for the model.
- **Vocabulary** - vocabulary of 50k tokens using BPE. The vocabulary was generated from a randomly selected 2% subset of the training data.

# Galactica

*Gaussian Error Linear Units function (GeLu)*

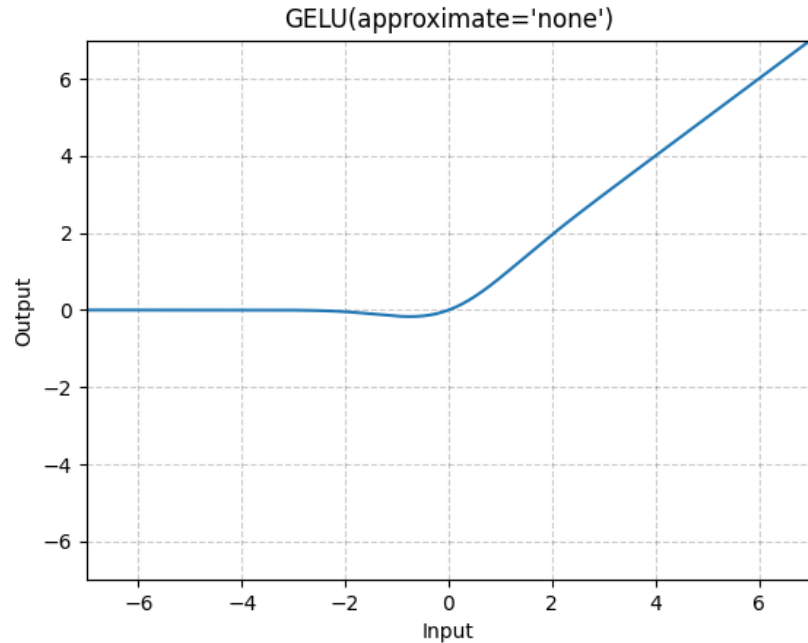
$$GELU(x) = x * \Phi(x)$$

Where  $\Phi(x)$  is the Gaussian function.

$$GELU(x) \approx x * \frac{1}{2} (1 + \tanh(\frac{2}{\pi} * (x + 0.044715 * x^3)))$$



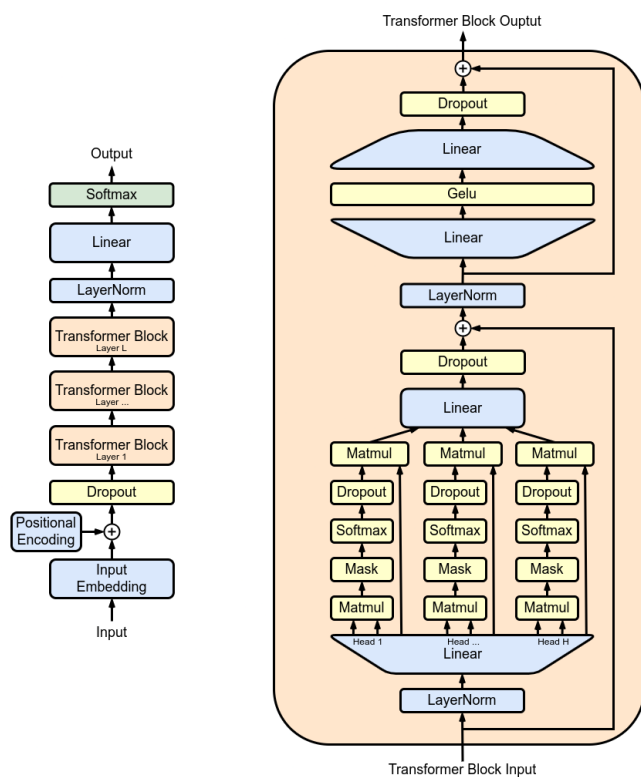
# Galactica



- Allows small negative values when  $x < 0$ .
- Avoids the dying ReLU problem.

# Galactica

## Why no biases?



# Unsupervised Classification Models

# Out-of-the-box representations

pass

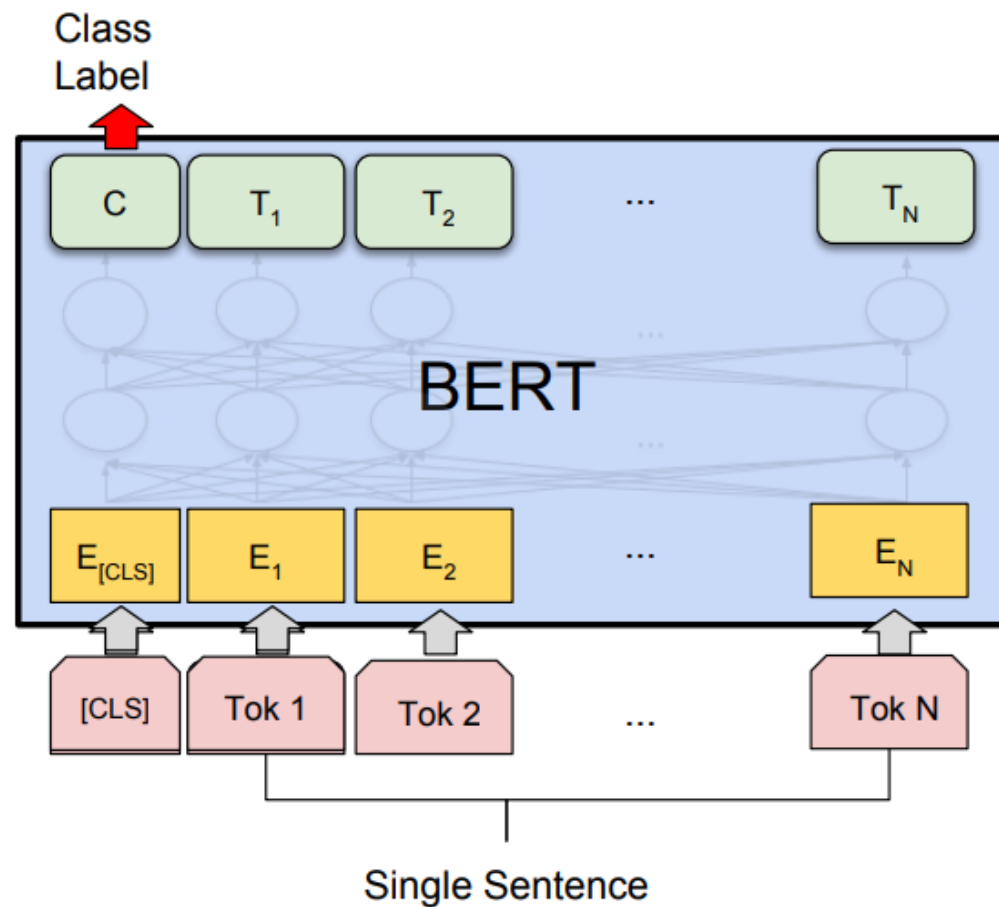
# Sentiment Analysis

# Sentiment Analysis

**Sentiment analysis** is a sentence classification task aiming at **automatically mapping data to their sentiment**.

It can be **binary** classification (e.g., positive or negative) or **multiclass** (e.g., enthusiasm, anger, etc)

# Sentiment Analysis



# Sentiment Analysis

The loss can be the likes of cross-entropy (CE), binary cross-entropy (BCE) or KL-Divergence (KL).

$$\mathcal{L}_{CE} = -\frac{1}{N} \sum_{n'=1}^N y^{(n)} \cdot \log(f(\mathbf{x}, \theta)^{(n)})$$

$$\mathcal{L}_{BCE} = -y^{(n)} \cdot \log(f(\mathbf{x}, \theta)^{(n)}) + (1 - y^{(n)}) \cdot (1 - f(\mathbf{x}, \theta)^{(n)})$$

$$\mathcal{L}_{KL} = -\frac{1}{N} \sum_{n'=1}^N y^{(n)} \cdot \log\left(\frac{y^{(n)}}{f(\mathbf{x}, \theta)^{(n)}}\right)$$



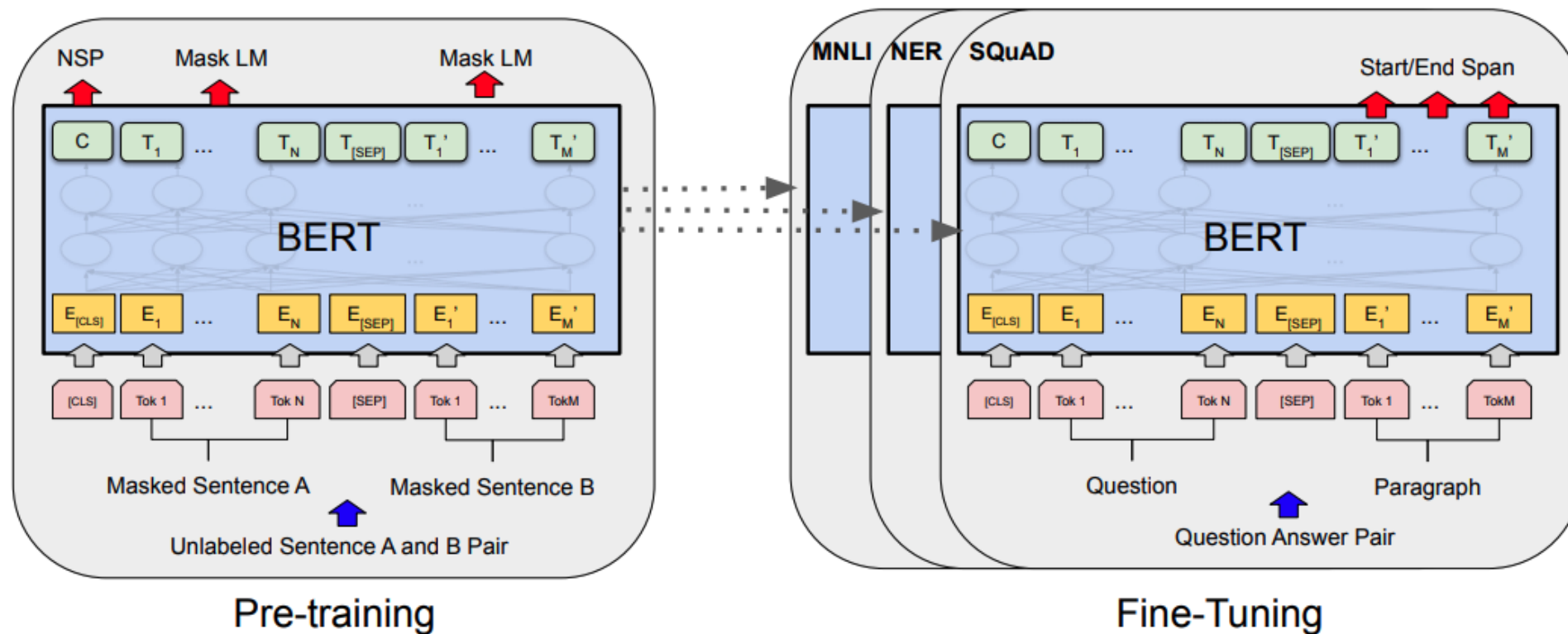
# Question Answering (QA)

# Question Answering (QA)

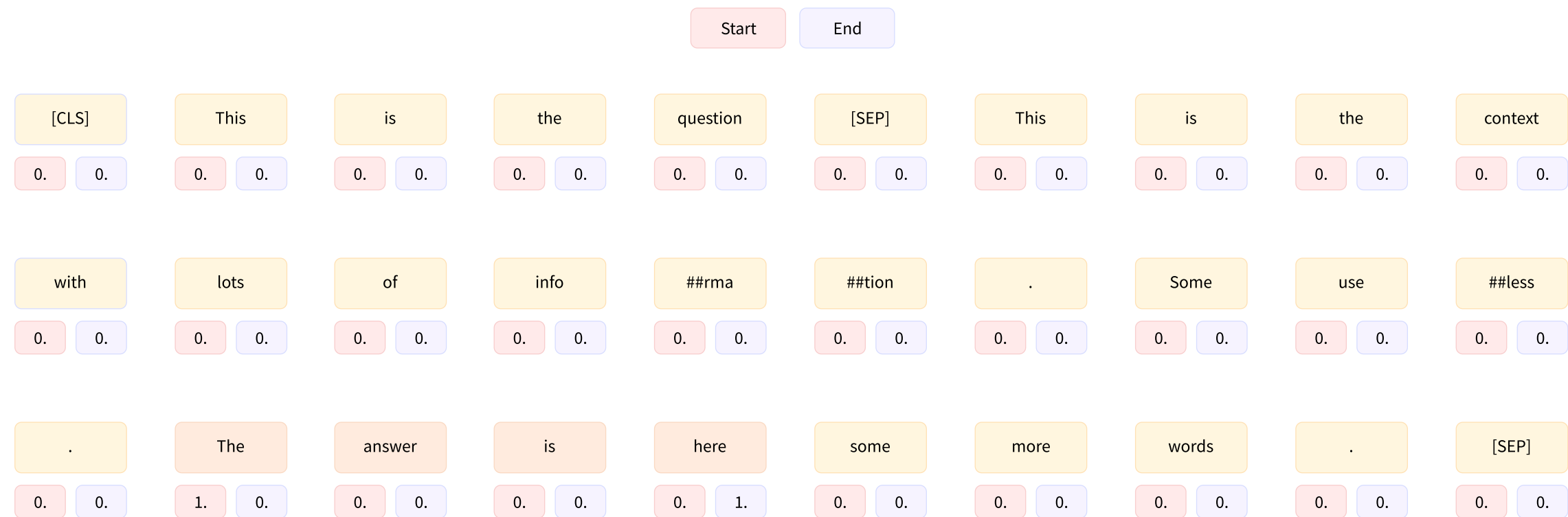
**QA** is the task of **retrieving a span of text from a context** that is best suited to answer a question.

This task is extractive -> **information retrieval**

# Question Answering (QA)



# Question Answering (QA)



# Question Answering (QA)

The loss is the cross entropy over the output of the starting token and the ending one:

$$\mathcal{L}_{CE_{QA}} = \mathcal{L}_{CE_{start}} + \mathcal{L}_{CE_{end}}$$

# Natural Language Inference (NLI)

# Natural Language Inference (NLI)

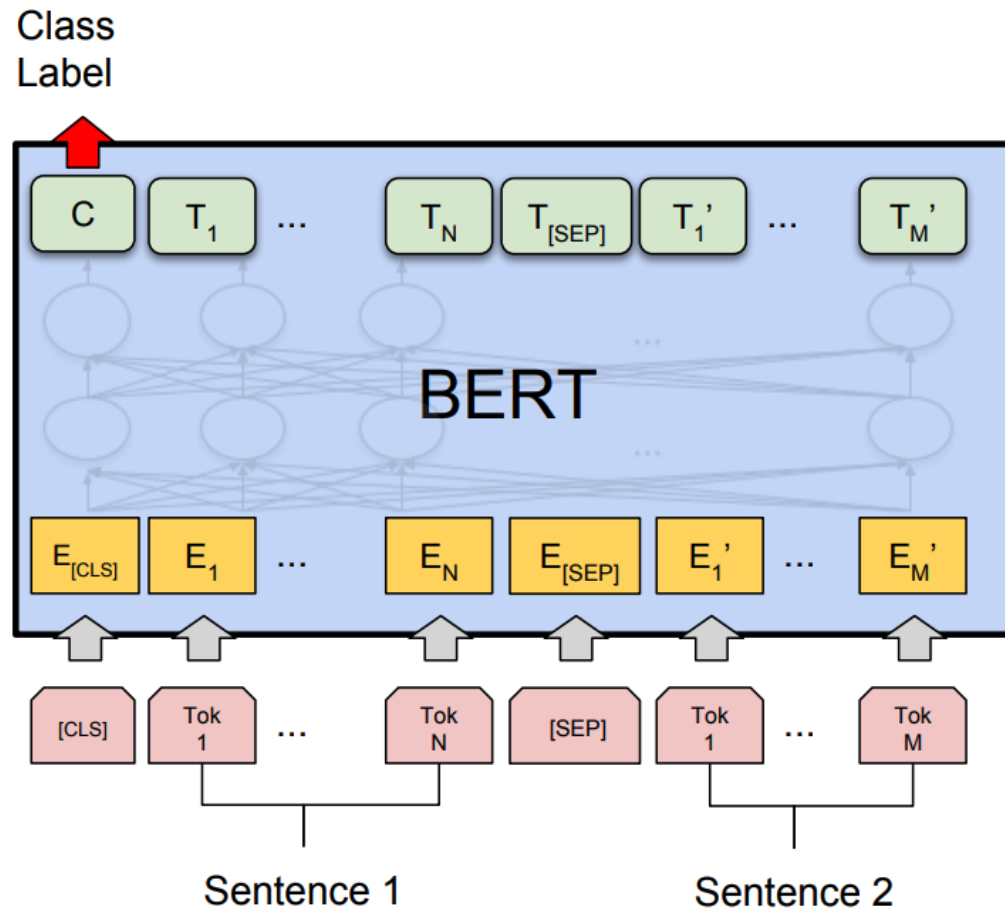
**NLI** is the task of **determining whether a "hypothesis" is true (entailment), false (contradiction), or undetermined (neutral)** given a "premise". [1]

# Natural Language Inference (NLI)

Premise	Label	Hypothesis
A man inspects the uniform of a figure in some East Asian country.	contradiction	The man is sleeping.
An older and younger man smiling.	neutral	Two men are smiling and laughing at the cats playing on the floor.
A soccer game with multiple males playing.	entailment	Some men are playing a sport.



# Natural Language Inference (NLI)



# Natural Language Inference (NLI)

The loss is simply the cross entropy or the divergence over the output of the **CLS** token and the true label.

$$\mathcal{L}_{NLI} = \mathcal{L}_{CE_{CLS}}$$

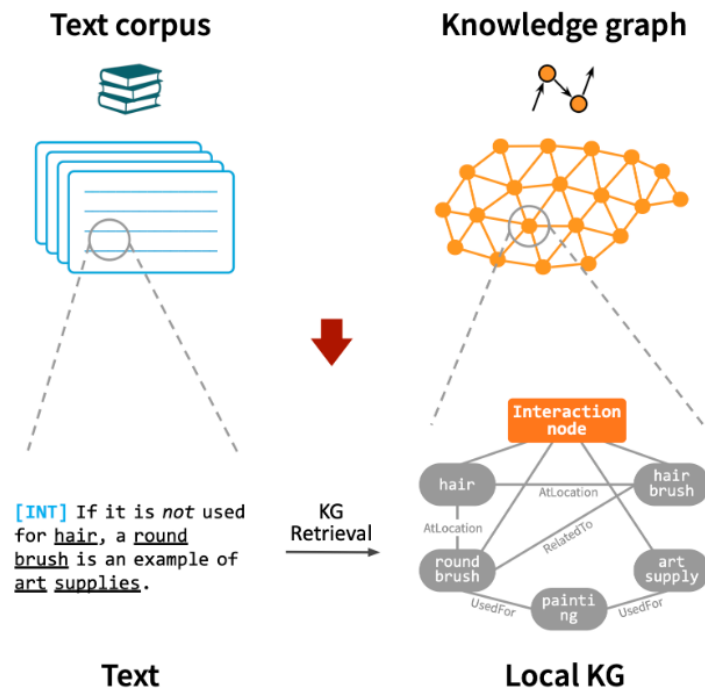
We are trying to compress the information about both sentence in one **CLS** token via attention and decide about their relationship.

Is it possible to help the model inferring more information with less text data?

# Going Further: LM as Knowledge Graphs

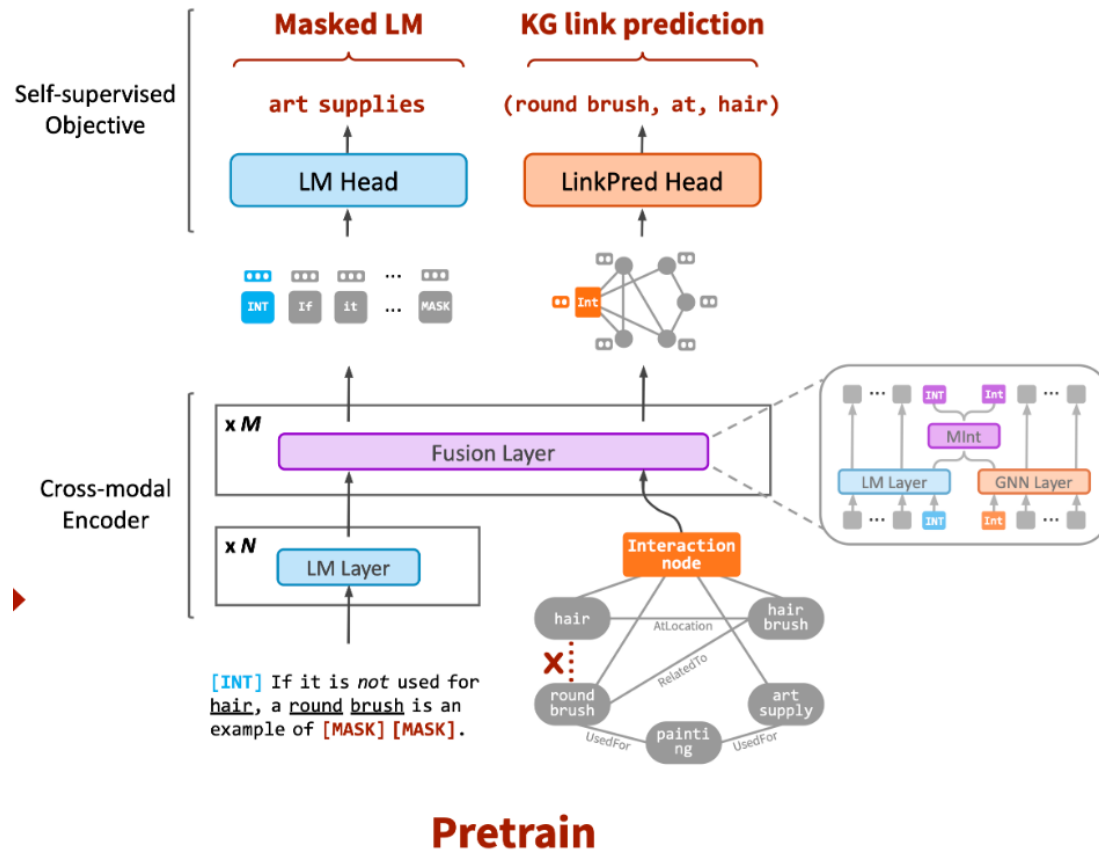
Yasunaga, M., Bosselut, A., Ren, H., Zhang, X., Manning, C. D., Liang, P. S., & Leskovec, J. (2022). [Deep bidirectional language-knowledge graph pretraining](#). Advances in Neural Information Processing Systems, 35, 37309-37323.

# Going Further: LM as Knowledge Graphs



**Raw data**

# Going Further: LM as Knowledge Graphs



# Going Further: LM as Knowledge Graphs

This architecture *involves a KG ready to use beforehand and pre-training from scratch*. How can we better **perform NLP task without having to retrain or fine-tune** a model?

# **Exploit LLMs capacities: Chain-of-thoughts & In context Learning**

# Exploit LLMs capacities

ICL enables LLMs to learn new tasks using natural language prompts without explicit retraining or fine-tuning.

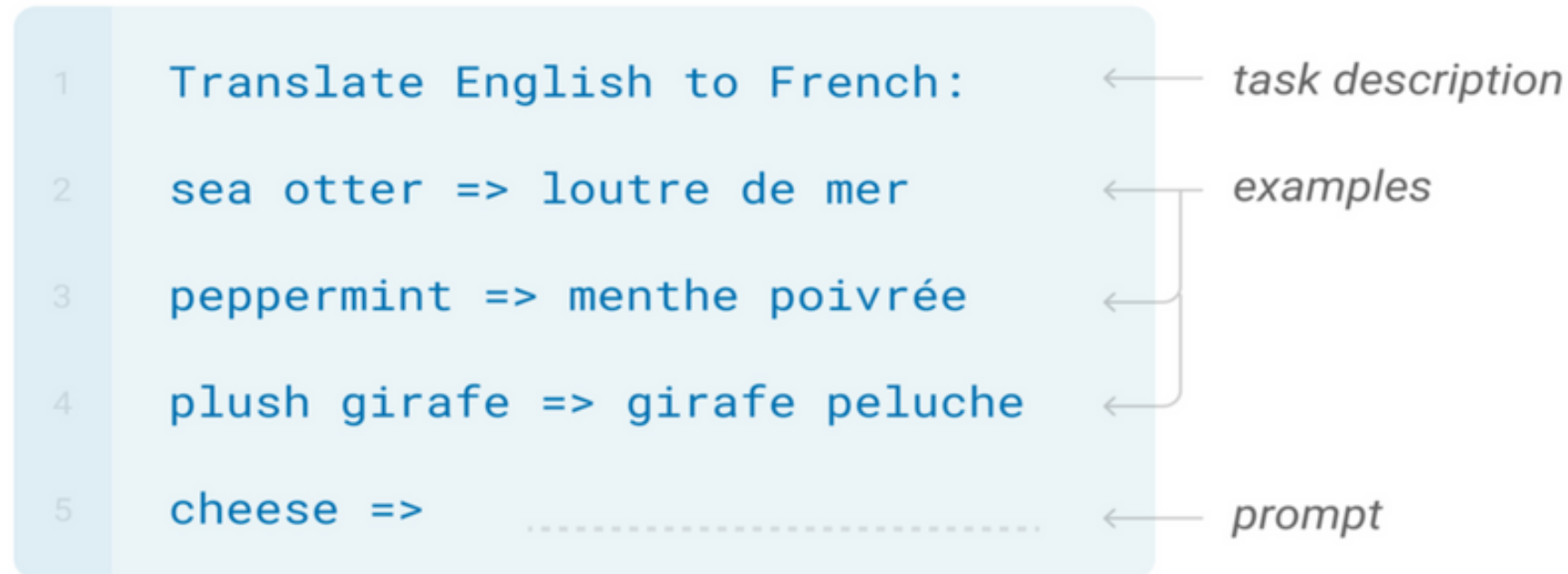
The efficacy of ICL is closely tied to the model's scale, training data quality, and domain specificity.



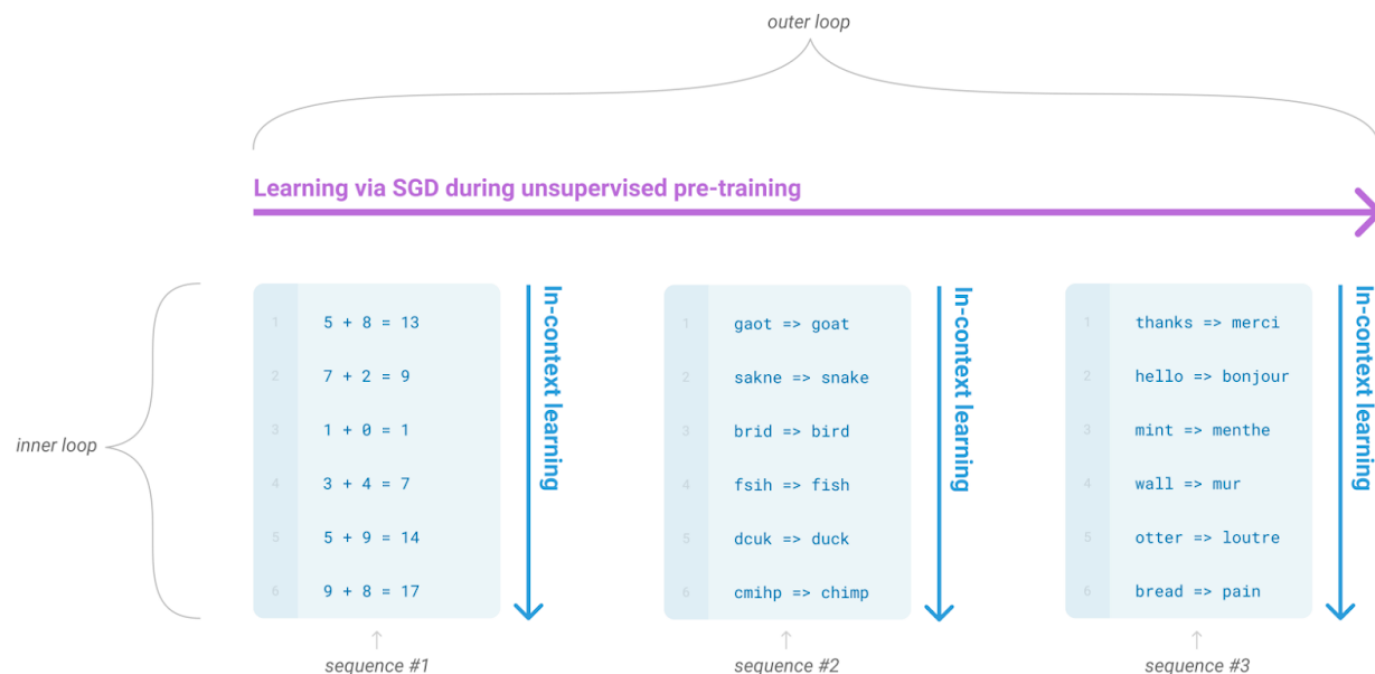
# Exploit LLMs capacities

## Few-shot

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



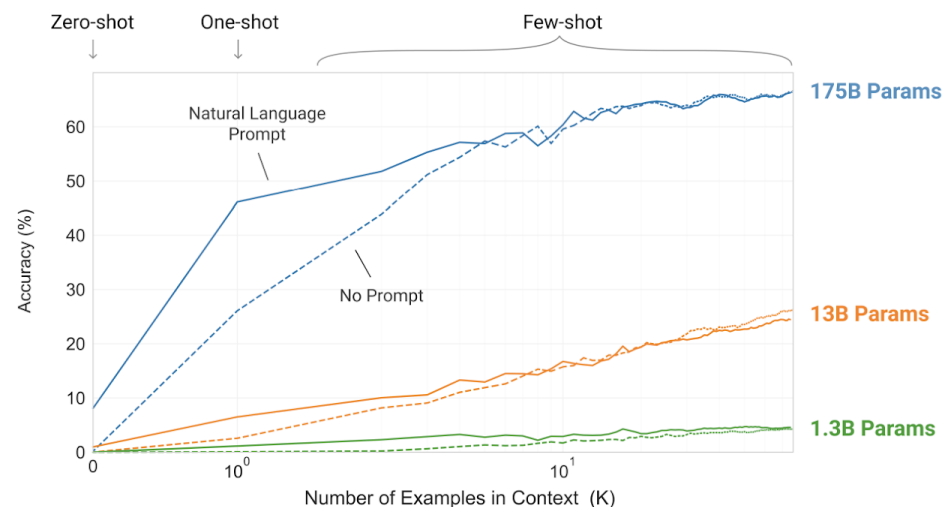
# Exploit LLMs capacities



**Figure 1.1: Language model meta-learning.** During unsupervised pre-training, a language model develops a broad set of skills and pattern recognition abilities. It then uses these abilities at inference time to rapidly adapt to or recognize the desired task. We use the term “in-context learning” to describe the inner loop of this process, which occurs within the forward-pass upon each sequence. The sequences in this diagram are not intended to be representative of the data a model would see during pre-training, but are intended to show that there are sometimes repeated sub-tasks embedded within a single sequence.

# Exploit LLMs capacities

$$p(\text{output}|\text{prompt}) = \int_{\text{concept}} p(\text{output}|\text{concept}, \text{prompt})p(\text{concept}|\text{prompt})d(\text{concept}). \quad (1)$$



**Figure 1.2: Larger models make increasingly efficient use of in-context information.** We show in-context learning performance on a simple task requiring the model to remove random symbols from a word, both with and without a natural language task description (see Sec. 3.9.2). The steeper “in-context learning curves” for large models demonstrate improved ability to learn a task from contextual information. We see qualitatively similar behavior across a wide range of tasks.

# Exploit LLMs capacities

(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. ✗

(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 ✗

(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. ✓

# Questions?

# References

[1] Gururangan, Suchin, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A. Smith. “[Don’t Stop Pretraining: Adapt Language Models to Domains and Tasks.](#)” In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, edited by Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault, 8342–60. Online: Association for Computational Linguistics, 2020.

[2] Lee, Jinhyuk, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. “[BioBERT: A Pre-Trained Biomedical Language Representation Model for Biomedical Text Mining](#).” *Bioinformatics* 36, no. 4 (February 15, 2020): 1234–40.

[3] Beltagy, Iz, Kyle Lo, and Arman Cohan. “[SciBERT: A Pretrained Language Model for Scientific Text](#).” In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, edited by Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan, 3615–20. Hong Kong, China: Association for Computational Linguistics, 2019.

[4] Taylor, Ross, Marcin Kardas, Guillem Cucurull, Thomas Scialom, Anthony Hartshorn, Elvis Saravia, Andrew Poulton, Viktor Kerkez, and Robert Stojnic. “[Galactica: A Large Language Model for Science.](#)” arXiv, November 16, 2022.