

Foundational Language Models

Mounika Marreddy,

University of Bonn, Germany

mounika0559@gmail.com



Agenda

- Recap on small language models [15 mins]
- Emerging abilities of language models and why they are effective? [30 mins]
- Text-to-Text Transfer Transformer [15 mins]
- Prompting [15 mins]
- Instruction-tuning [15 mins]

Agenda

- Recap on small language models [15 mins]
- Emerging abilities of language models and why they are effective? [30 mins]
- Text-to-Text Transfer Transformer [15 mins]
- Prompting [15 mins]
- Instruction-tuning [15 mins]

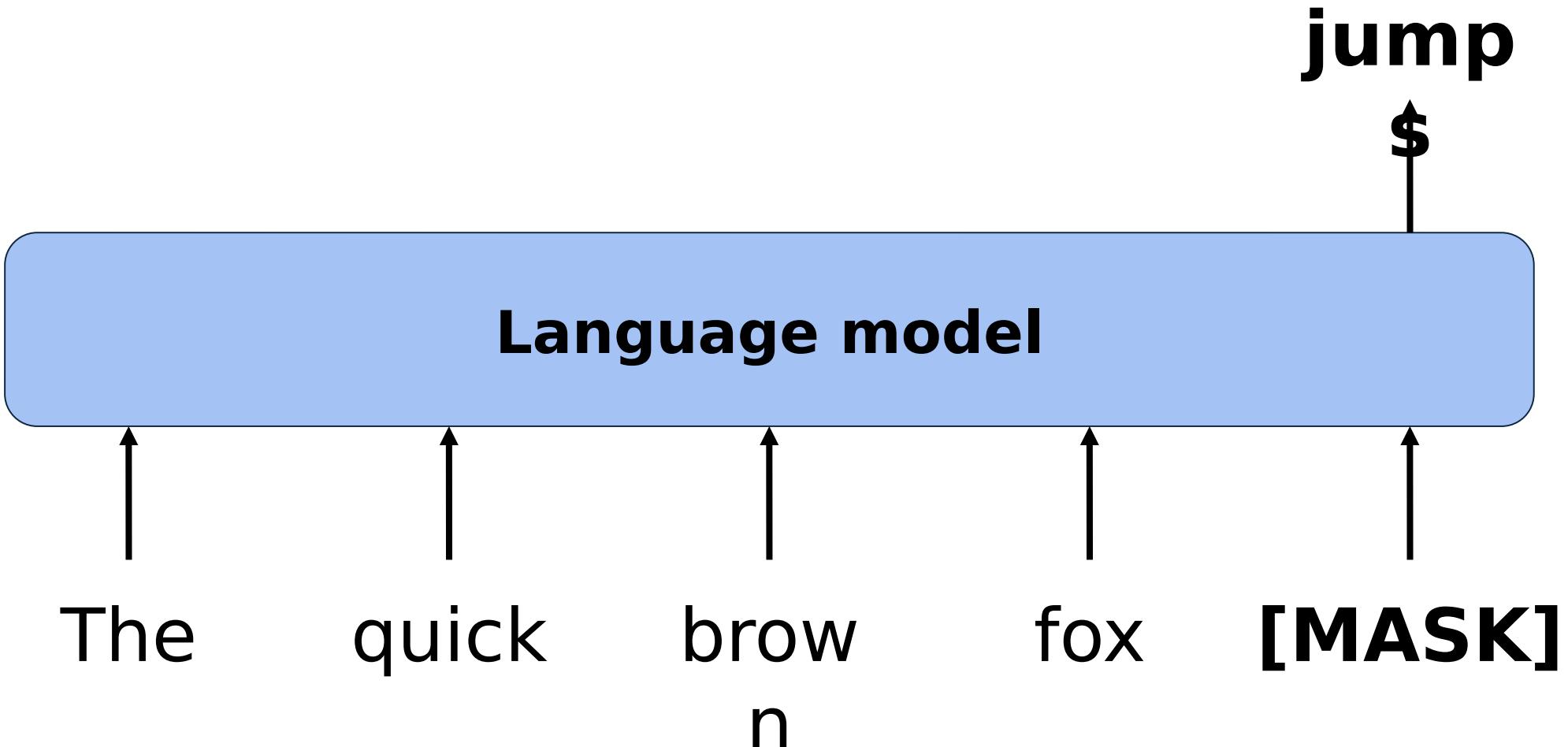
Foundation Models

Foundation models:

BERT, GPT-2 variations, T5, Flan-T5 and so on.

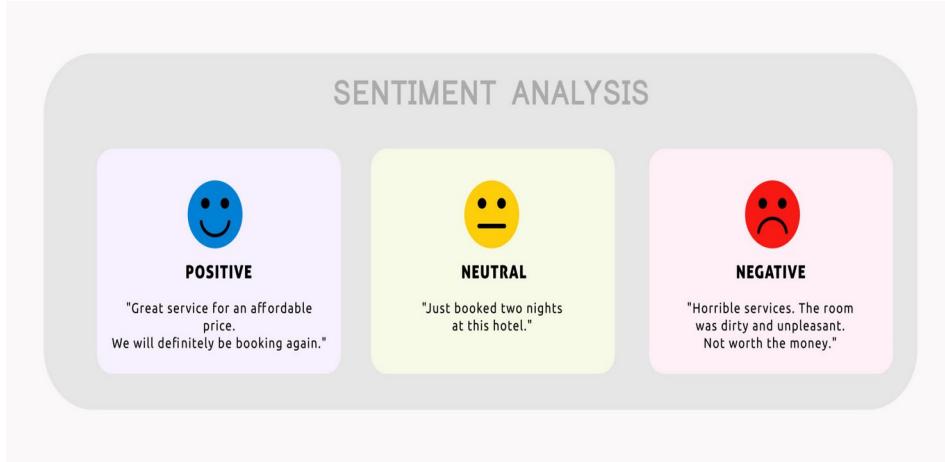
- What drives the success of these models?
 - Data
 - Hardware
 - Self-supervised learning
 - Transformer architectures

LMs are trained to predict missing words



Language models are everywhere

Sentiment
Question Answering



Passage Sentence

In meteorology, precipitation is any product of the condensation of atmospheric water vapor that falls under gravity.

Question

What causes precipitation to fall?

Answer Candidate

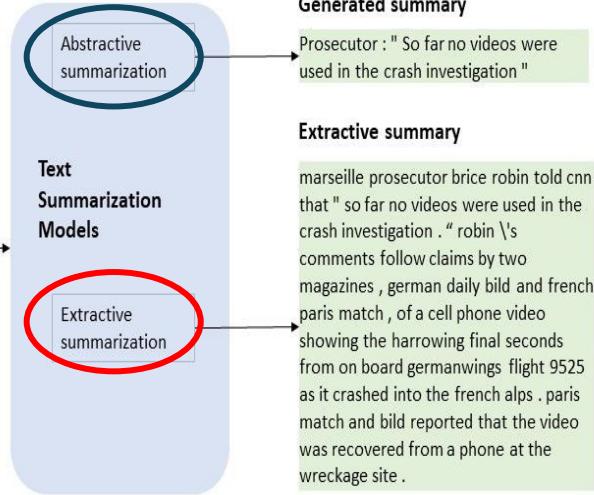
gravity

Summarization

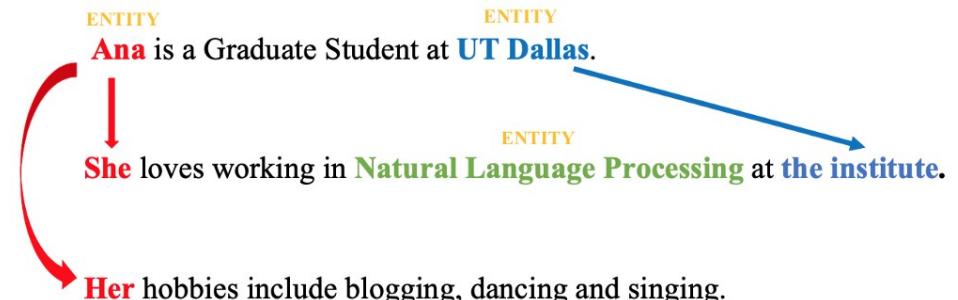
Input Article

Marseille, France (CNN) The French prosecutor leading an investigation into the crash of Germanwings Flight 9525 insisted Wednesday that he was not aware of any video footage from on board the plane. Marseille prosecutor Brice Robin told CNN that " so far no videos were used in the crash investigation ." He added, " A person who has such a video needs to immediately give it to the investigators ." Robin's comments follow claims by two magazines, German daily Bild and French Paris Match, of a cell phone video showing the harrowing final seconds from on board Germanwings Flight 9525 as it crashed into the French Alps . All 150 on board were killed. Paris Match and Bild reported that the video was recovered from a phone at the wreckage site. ...

Text



Coreference Resolution



Language models are everywhere



Small language
models



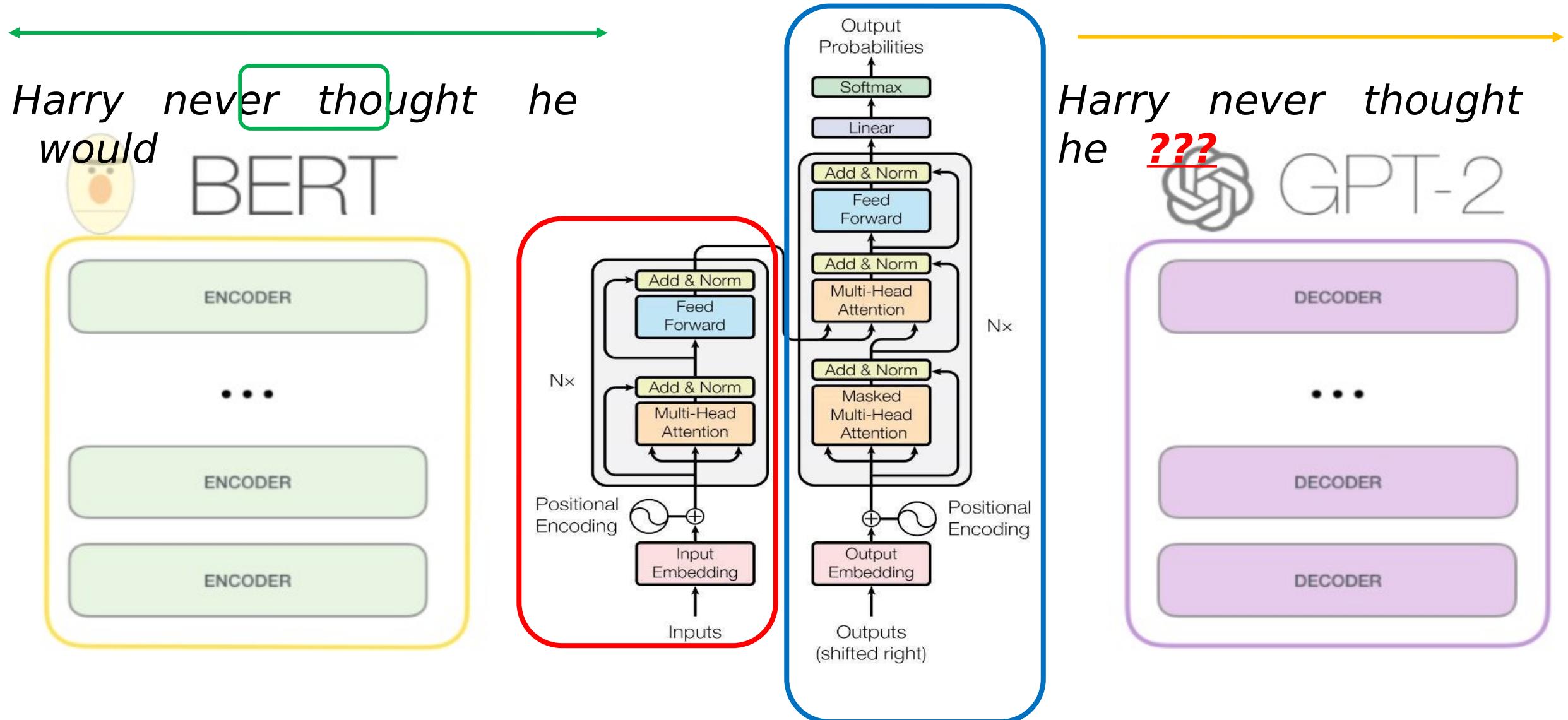
MMLU (Massive Multitask
Language Understanding)

BIG-bench A green oval surrounds the text "BIG-bench" and a small brown chair icon.

<https://gluebenchmark.com/>
<https://supergluebenchmark.com/>

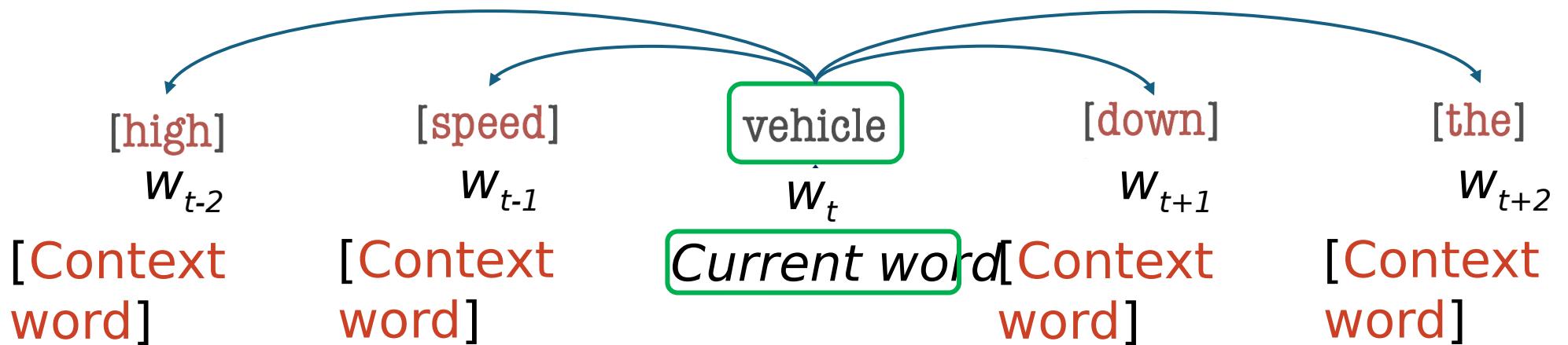
<https://crfm.stanford.edu/helm/lite/latest/>
<https://paperswithcode.com/dataset/mmlu>
<https://paperswithcode.com/dataset/big-bench>

Transformer

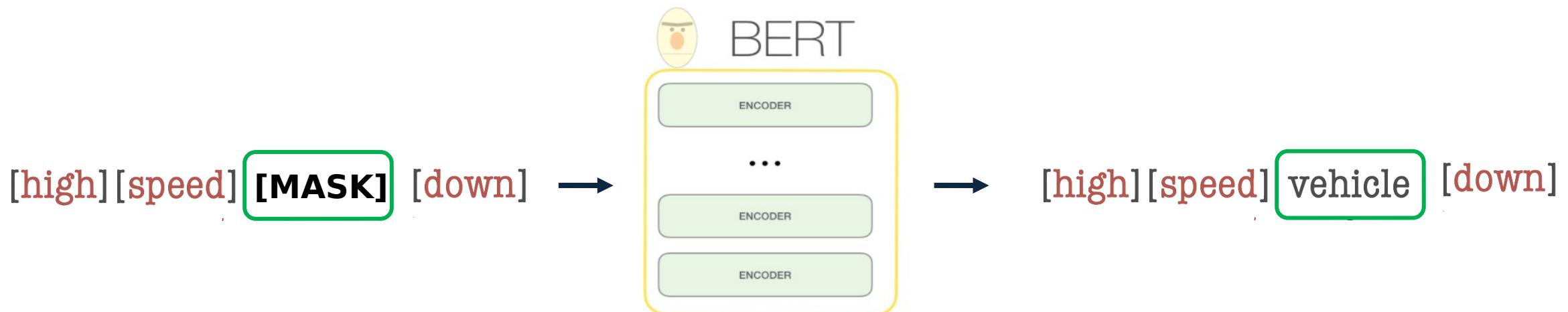


BERT: Workflow

1. Self-attention with Bi-directional context

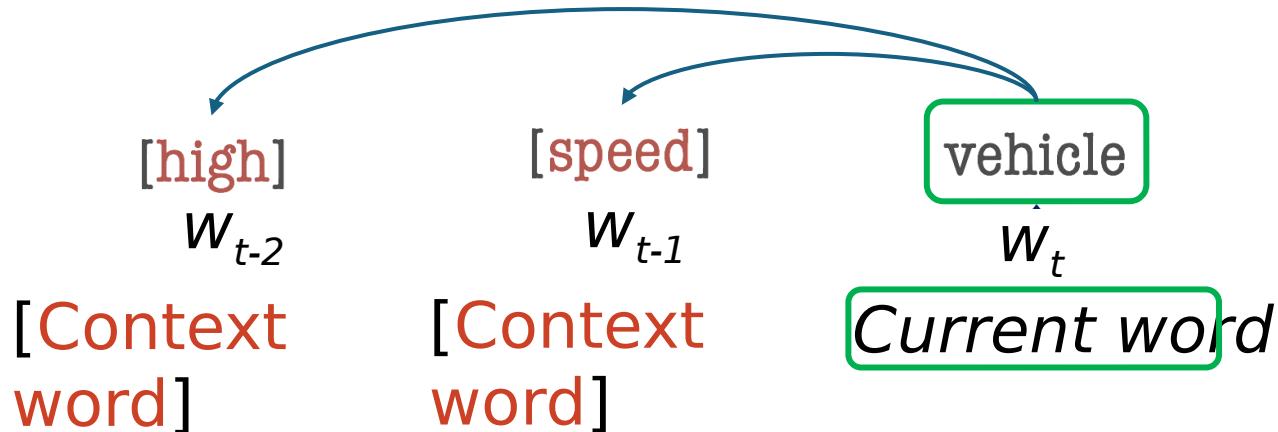


2. Masked language modeling (MLM)



GPT2: Workflow

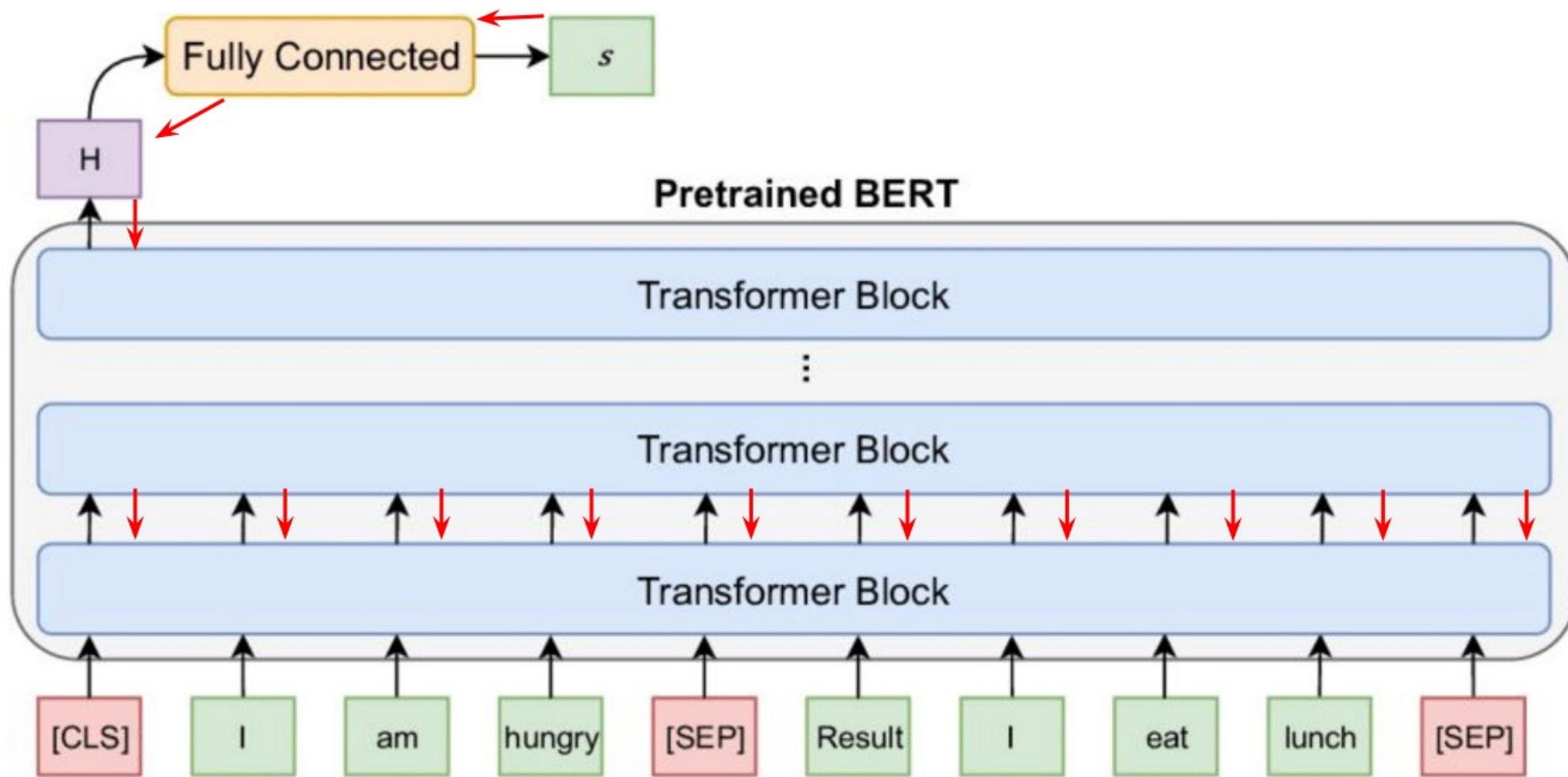
1. Self-attention with Uni-directional context



2. Causal language modeling (CLM)



Fine tuning: tune pretrain language model on a task



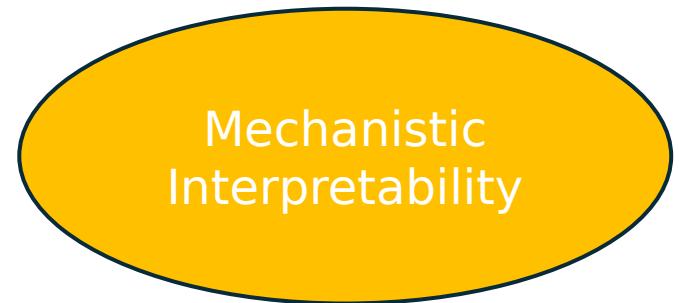
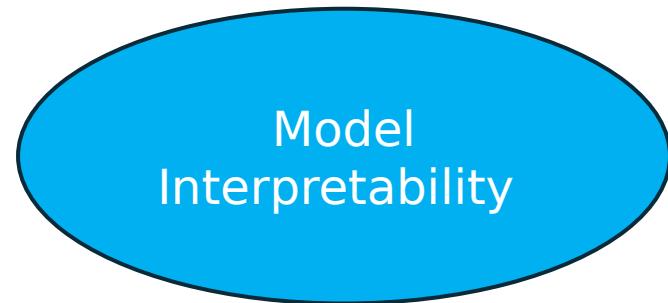
Agenda

- Recap on small language models [5 min]
- Emerging abilities of language models and why they are effective? [25 min]
- Text-to-Text Transfer Transformer [15 mins]
- Prompting [15 mins]
- Instruction-tuning [15 mins]

Agenda

- Recap on small language models [5 min]
- Emerging abilities of language models and why they are effective? [25 min]
 - Analyzing and Interpreting language models
- Text-to-Text Transfer Transformer [15 mins]
- Prompting [15 mins]
- Instruction-tuning [15 mins]

Analyzing and Interpreting LMs

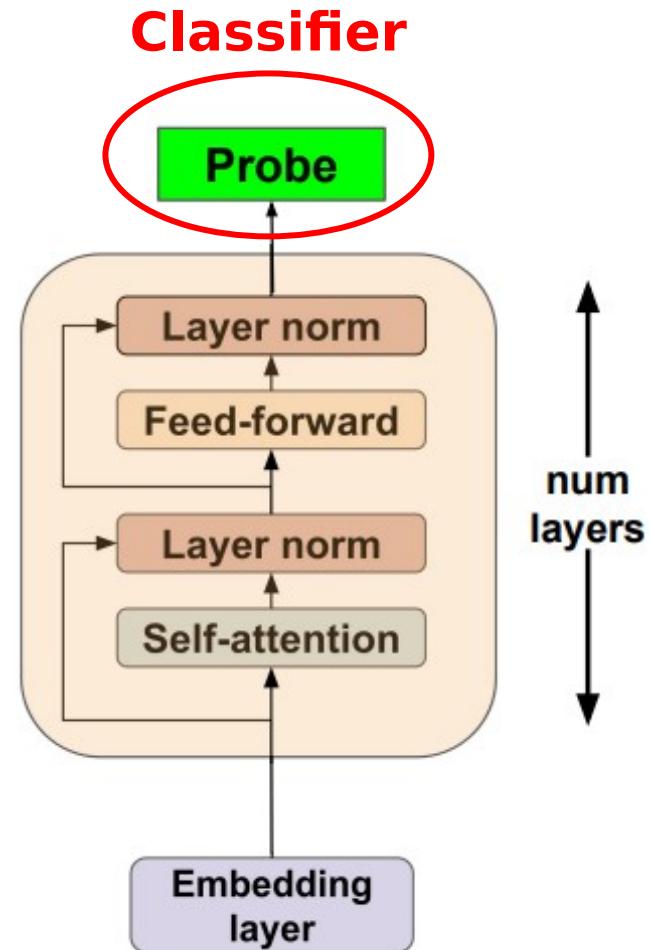
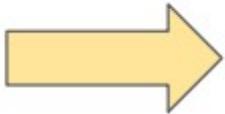
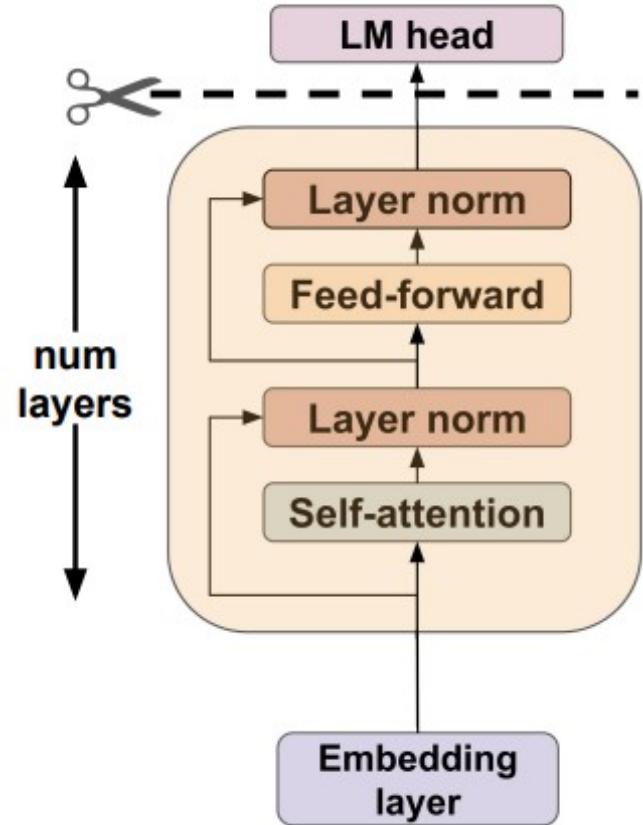


**Analysis of
representations via
probing**

**Bridging language
models and brain**

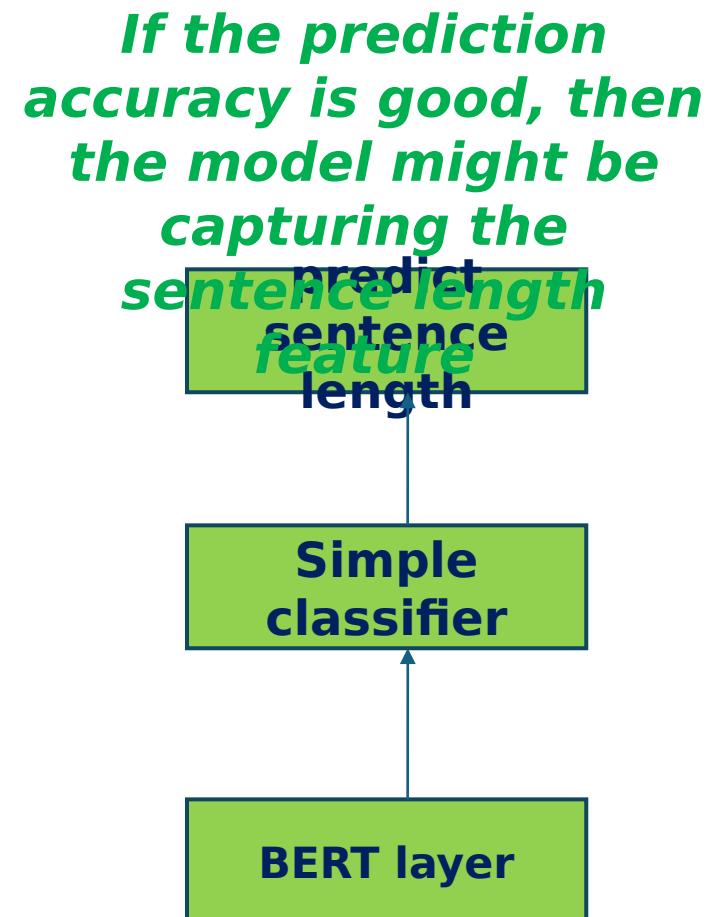
**Robustness
Reverse engineering of
neural computations**

Model Interpretability?



Hierarchy of Linguistic Info - Setting

- Conneau et al., ACL'18 - Build diagnostic classifier to predict if a linguistic property is encoded in the given sentence representation.
- Features:
 - **Surface** – Sentence Length, Word Content
 - **Syntactic** – Bigram shift, Tree depth, Top constituent
 - **Semantic** – Tense, Subject Number, Object Number, Coordination Inversion and Semantic Odd Man Out.

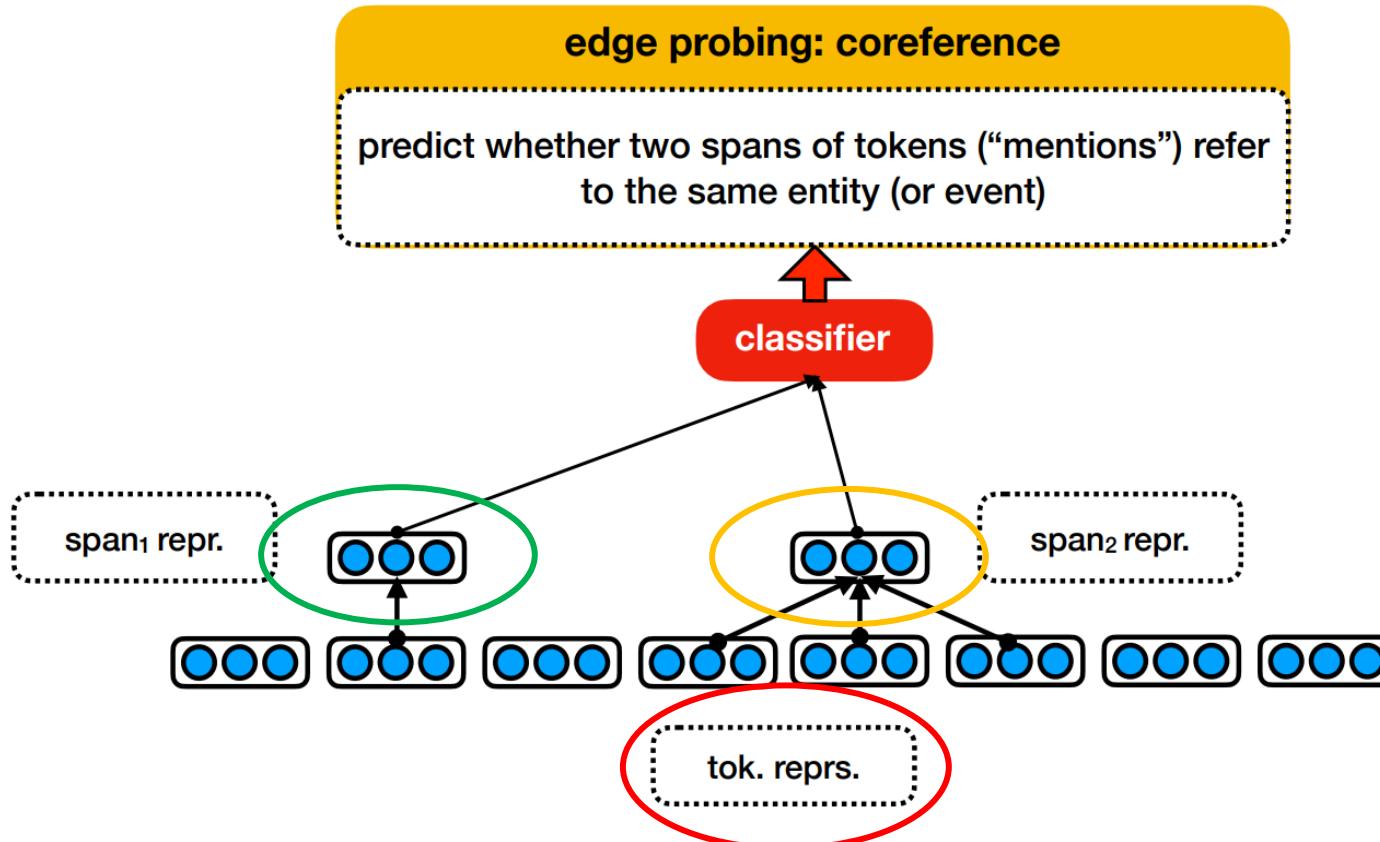


BERT composes a hierarchy of linguistic signals ranging from surface to semantic features

	Surface	Syntactic				Semantic				
Layer	SentLen (Surface)	WC (Surface)	TreeDepth (Syntactic)	TopConst (Syntactic)	BShift (Syntactic)	Tense (Semantic)	SubjNum (Semantic)	ObjNum (Semantic)	SOMO (Semantic)	CoordInv (Semantic)
1	93.9 (2.0)	24.9 (24.8)	35.9 (6.1)	63.6 (9.0)	50.3 (0.3)	82.2 (18.4)	77.6 (10.2)	76.7 (26.3)	49.9 (-0.1)	53.9 (3.9)
2	95.9 (3.4)	65.0 (64.8)	40.6 (11.3)	71.3 (16.1)	55.8 (5.8)	85.9 (23.5)	82.5 (15.3)	80.6 (17.1)	53.8 (4.4)	58.5 (8.5)
3	96.2 (3.9)	66.5 (66.0)	39.7 (10.4)	71.5 (18.5)	64.9 (14.9)	86.6 (23.8)	82.0 (14.6)	80.3 (16.6)	55.8 (5.9)	59.3 (9.3)
4	94.2 (2.3)	69.8 (69.6)	39.4 (10.8)	71.3 (18.3)	74.4 (24.5)	87.6 (25.2)	81.9 (15.0)	81.4 (19.1)	59.0 (8.5)	58.1 (8.1)
5	92.0 (0.5)	69.2 (69.0)	40.6 (11.8)	81.3 (30.8)	81.4 (31.4)	89.5 (26.7)	85.8 (19.4)	81.2 (18.6)	60.2 (10.3)	64.1 (14.1)
6	88.4 (-3.0)	63.5 (63.4)	41.3 (13.0)	83.3 (36.6)	82.9 (32.9)	89.8 (27.6)	88.1 (21.9)	82.0 (20.1)	60.7 (10.2)	71.1 (21.2)
7	83.7 (-7.7)	56.9 (56.7)	40.1 (12.0)	84.1 (39.5)	83.0 (32.9)	89.9 (27.5)	87.4 (22.2)	82.2 (21.1)	61.6 (11.7)	74.8 (24.9)
8	82.9 (-8.1)	51.1 (51.0)	39.2 (10.3)	84.0 (39.5)	83.9 (33.9)	89.9 (27.6)	87.5 (22.2)	81.2 (19.7)	62.1 (12.2)	76.4 (26.4)
9	80.1 (-11.1)	47.9 (47.8)	38.5 (10.8)	83.1 (39.8)	87.0 (37.1)	90.0 (28.0)	87.6 (22.9)	81.8 (20.5)	63.4 (13.4)	78.7 (28.9)
10	77.0 (-14.0)	43.4 (43.2)	38.1 (9.9)	81.7 (39.8)	86.7 (36.7)	89.7 (27.6)	87.1 (22.6)	80.5 (19.9)	63.3 (12.7)	78.4 (28.1)
11	73.9 (-17.0)	42.8 (42.7)	36.3 (7.9)	80.3 (39.1)	86.8 (36.8)	89.9 (27.8)	85.7 (21.9)	78.9 (18.6)	64.4 (14.5)	77.6 (27.9)
12	69.5 (-21.4)	49.1 (49.0)	34.7 (6.9)	76.5 (37.2)	86.4 (36.4)	89.5 (27.7)	84.0 (20.2)	78.7 (18.4)	65.2 (15.3)	74.9 (25.4)

Jawahar et al. 2019 ACL

Edge Probing

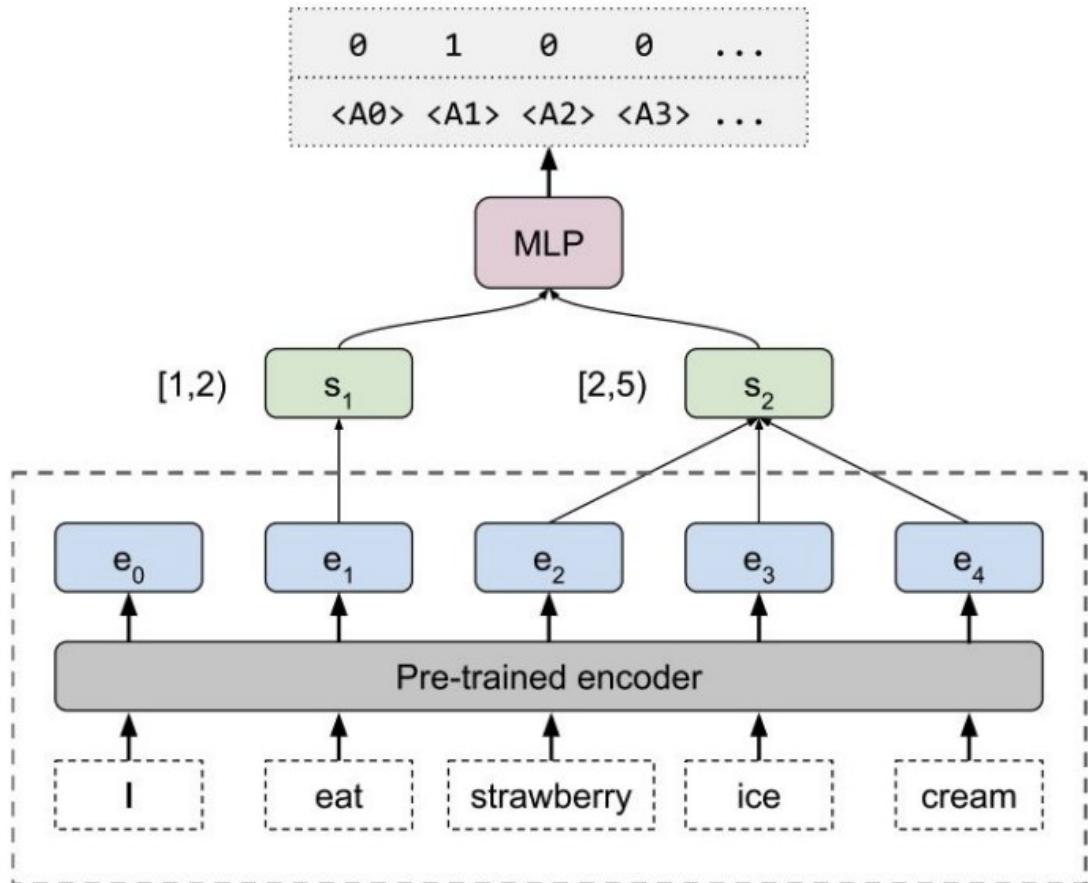


Edge Probing

POS	The important thing about Disney is that it is a global [brand] ₁ . → NN (Noun)
Constit.	The important thing about Disney is that it [is a global brand] ₁ . → VP (Verb Phrase)
Depend.	[Atmosphere] ₁ is always [fun] ₂ → nsubj (nominal subject)
Entities	The important thing about [Disney] ₁ is that it is a global brand. → Organization
SRL	[The important thing about Disney] ₂ [is] ₁ that it is a global brand. → Arg1 (Agent)
SPR	[It] ₁ [endorsed] ₂ the White House strategy... → {awareness, existed_after, ... }
Coref. ^O	The important thing about [Disney] ₁ is that [it] ₂ is a global brand. → True
Coref. ^W	[Characters] ₂ entertain audiences because [they] ₁ want people to be happy. → True Characters entertain [audiences] ₂ because [they] ₁ want people to be happy. → False
Rel.	The [burst] ₁ has been caused by water hammer [pressure] ₂ . → Cause-Effect(e_2, e_1)

Edge Probing

- Local syntax (word-level) captured at initial-middle layers
- High-level semantics captured at later layers



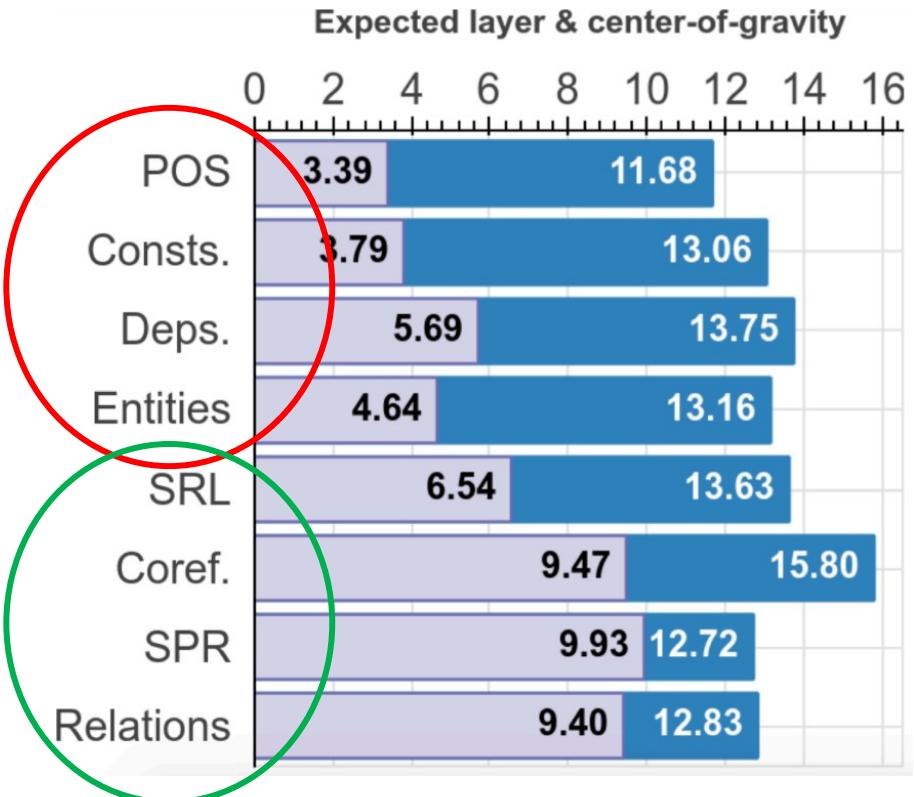
Labels

Binary classifiers

Span representations

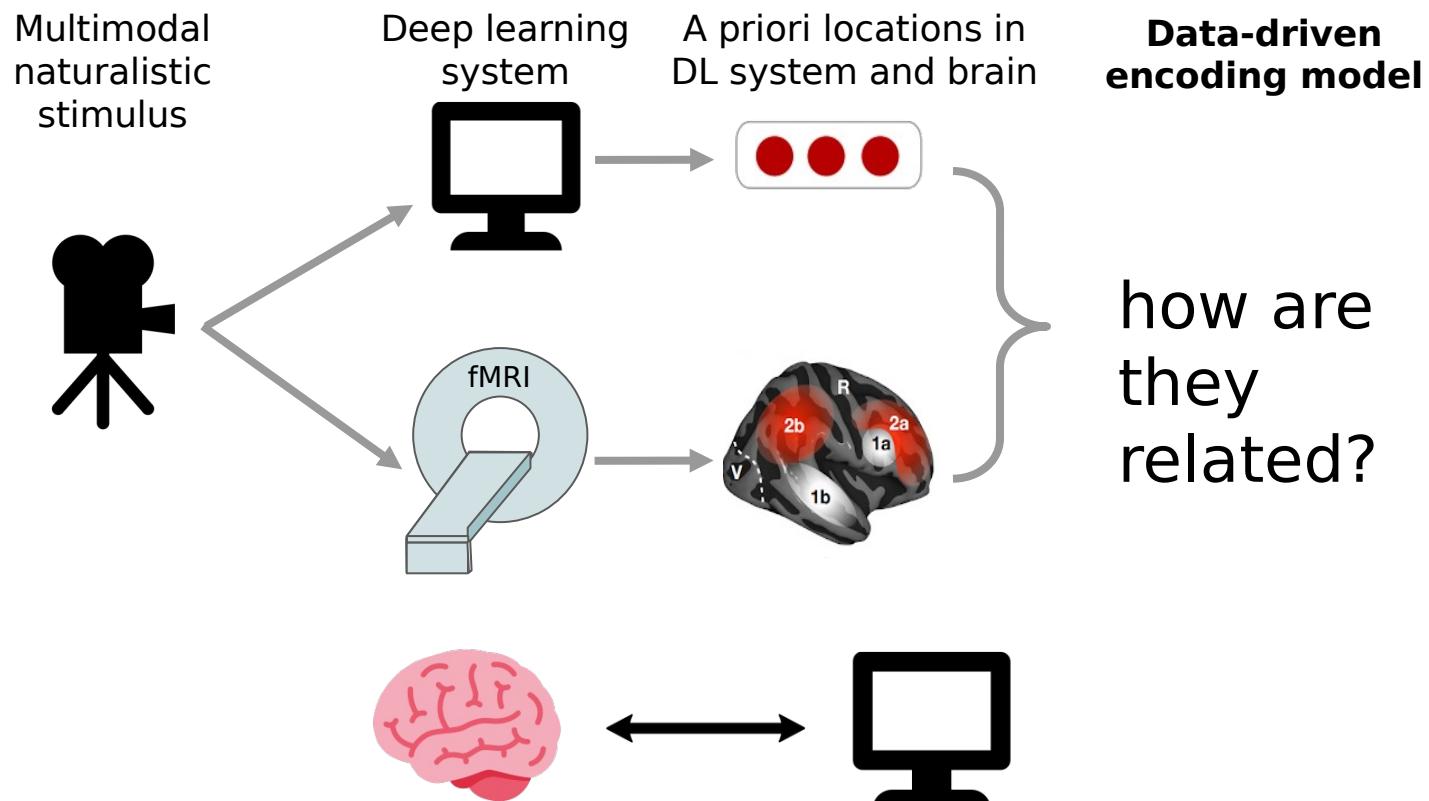
Contextual vectors

Input tokens

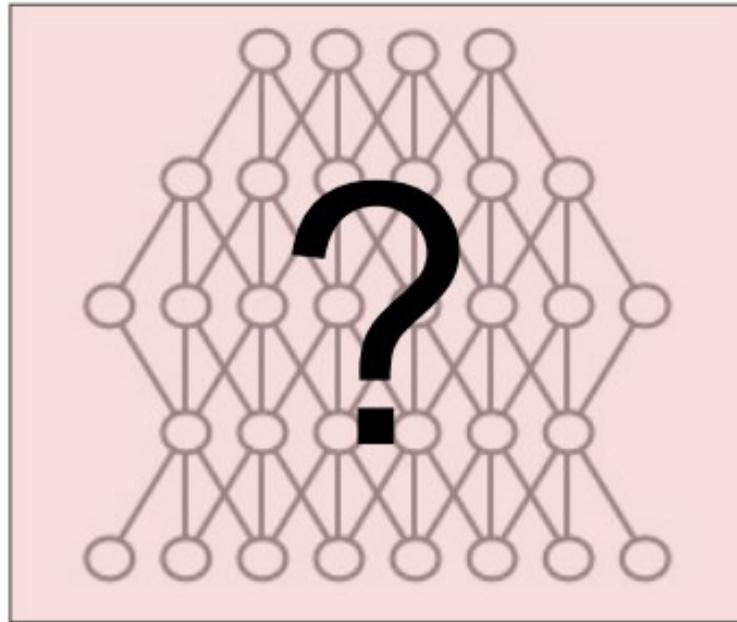


Behavioural Interpretability

- Data-driven encoding models evaluate the relationships between brains and deep learning models



Mechanistic Interpretability



Interpreting Model Predictions

- Why did the model make this prediction?

What if?
Drop layers..

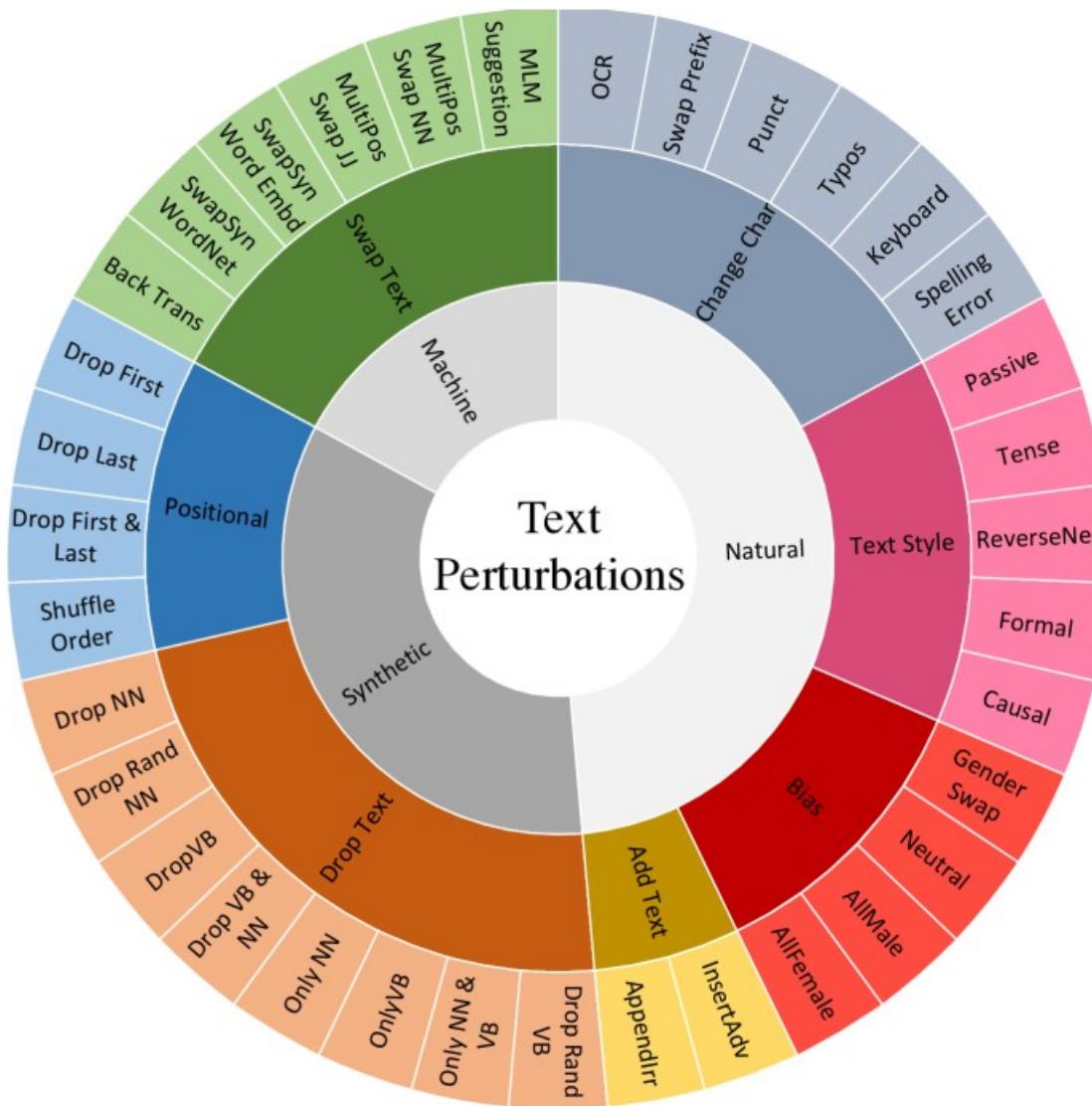
What if?
Change input examples..

What if?
Change of weights..

Interesting questions about BERT, GPT2

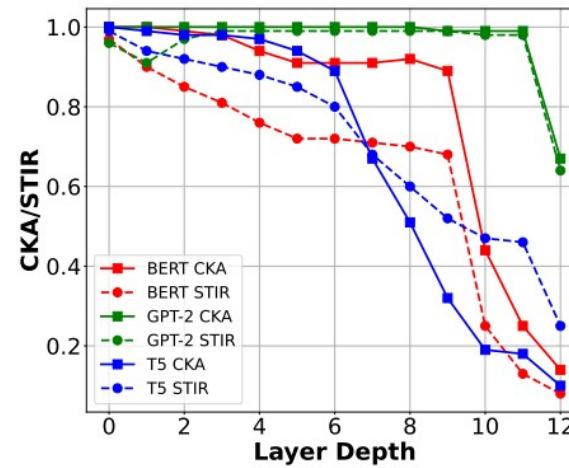
- Finetuning modifies the representations generated by each layer.
 - While fine-tuning these models, what changes across layers with respect to the pre-trained checkpoints?
- Robustness to input perturbations
 - How robust are BERT, GPT2 to input perturbations?
 - Is the effect of finetuning consistent across all models for various NLP tasks?
 - Do these models exhibit varying levels of robustness to input text perturbations when finetuned for different NLP tasks?
- Centred Kernel Alignment (CKA)
 - Compare layer-wise hidden state representations of pre-trained and finetuned models.
 - 1 ↗ perfect sim; 0 ↗ no sim.
- Similarity Through Inverted Representations (STIR) (finetuned|pre-trained)
 - Given pretrained model and data samples, STIR defines how invariant finetuned model is to perturbations of the samples that are imperceptible by , i.e., do not change their representations according to .

Robustness

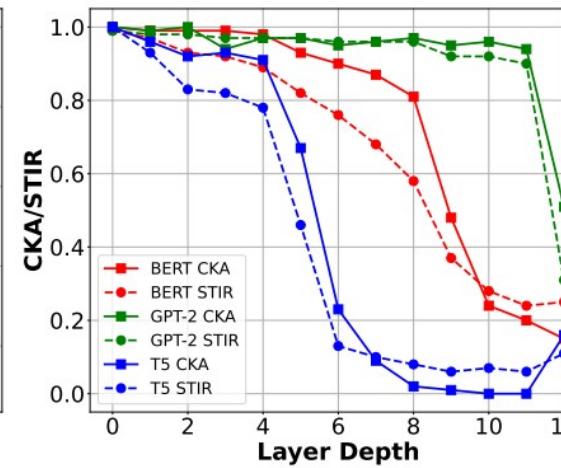


- Text perturbations grouped into seven different categories
 - ChangeChar
 - AddText
 - Bias
 - Positional
 - DropText
 - SwapText
 - TextStyle

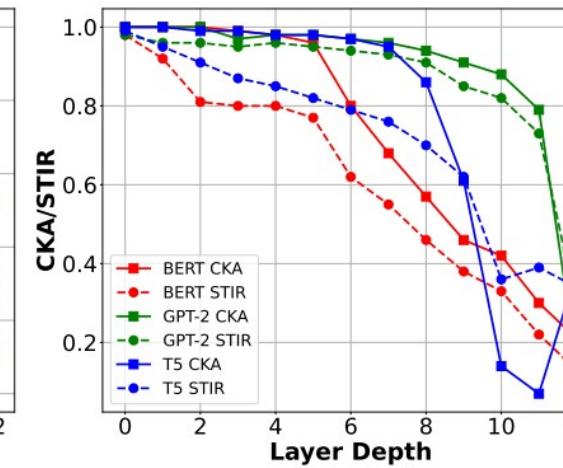
How does finetuning modify the layers representations for different models?



(a) CoLA

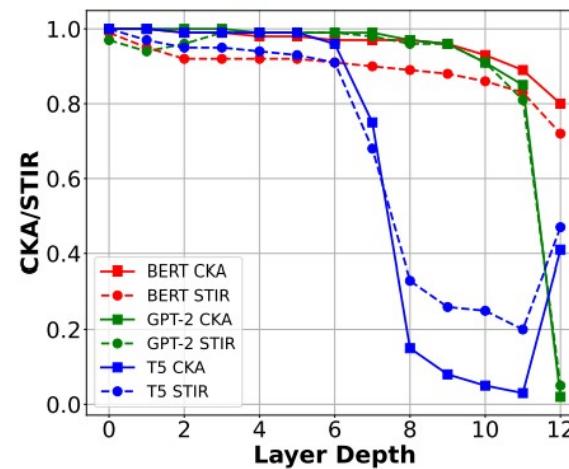


(b) SST-2

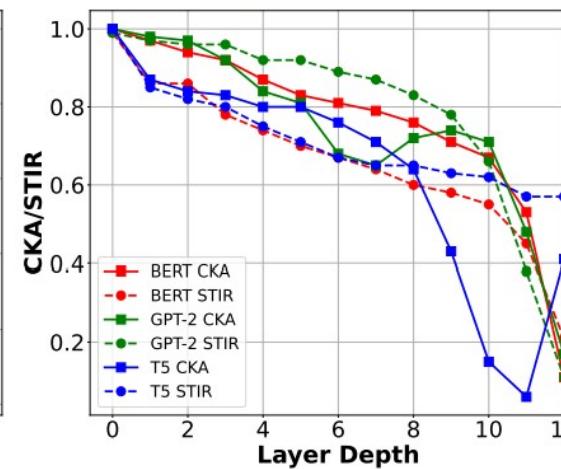


(c) MRPC

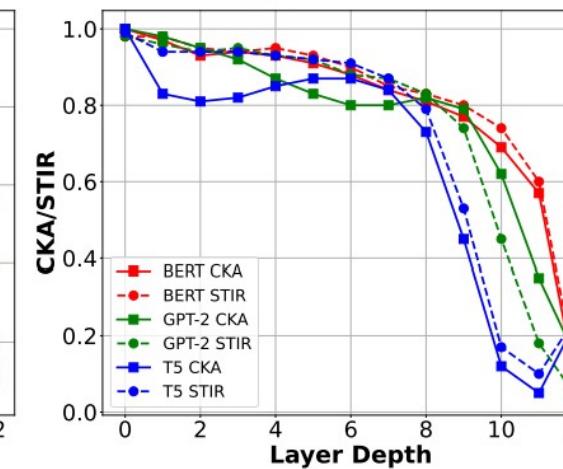
- GPT-2 had fewer affected layers, indicating higher semantic stability.
- CKA and STIR vary consistently across all three models.



(d) STS-B



(e) QQP



(f) MNLI-Matched

Is the impact of input text perturbations on finetuned models task-dependent?

Perturbation	CoLA (Matthews CC)			SST-2 (Accuracy)			MRPC (Accuracy)			STS-B (PearsonCC)			QQP (Accuracy)		
	BERT	GPT-2	T5	BERT	GPT-2	T5	BERT	GPT-2	T5	BERT	GPT-2	T5	BERT	GPT-2	T5
Drop nouns	0.18	<u>0.10</u>	0.24	0.92	0.93	0.93	0.94	0.96	0.94	0.56	0.48	0.57	0.89	0.92	0.89
Drop verbs	<u>0.05</u>	0.24	0.06	<u>0.95</u>	0.95	0.95	0.98	0.99	0.96	0.93	0.92	0.89	0.97	<u>0.96</u>	0.96
Drop first	0.48	0.75	0.54	0.98	0.97	0.98	1.00	0.99	1.00	0.98	0.93	0.94	0.99	0.98	0.99
Drop last	0.34	0.45	0.32	1.00	0.99	1.00	1.00	1.00	1.00	<u>0.84</u>	0.83	<u>0.83</u>	<u>0.95</u>	0.96	<u>0.95</u>
Swap text	0.13	0.16	0.06	0.98	0.98	0.97	0.99	1.01	0.98	0.98	<u>0.96</u>	0.95	0.97	0.97	0.96
Add text	0.85	0.92	0.86	0.99	0.99	0.99	0.93	1.00	<u>0.96</u>	0.99	0.99	0.98	0.99	1.00	0.99
Change char	0.14	0.29	0.29	<u>0.84</u>	0.86	0.84	0.43	0.97	0.65	0.58	0.52	0.57	<u>0.88</u>	0.95	0.94
Bias	0.95	0.96	0.92	1.00	1.01	1.00	1.00	1.00	1.00	0.99	0.99	0.99	1.00	1.01	1.00
Perturbation	MNLI-m (Accuracy)			QNLI (Accuracy)			RTE (Accuracy)			WNLI (Accuracy)					
	BERT	GPT-2	T5	BERT	GPT-2	T5	BERT	GPT-2	T5	BERT	GPT-2	T5	BERT	GPT-2	T5
Drop nouns	0.83	0.85	0.83	0.82	0.87	0.82	0.84	1.01	0.89	1.00	1.00	1.00	1.00	1.00	1.05
Drop verbs	0.89	0.90	0.90	0.96	0.94	0.94	0.98	1.01	<u>0.96</u>	1.00	1.01	1.01	1.00	1.01	1.03
Drop first	0.94	0.94	0.95	0.97	0.98	0.97	0.95	1.00	1.00	1.00	0.99	1.01	1.00	<u>0.99</u>	1.01
Drop last	0.89	0.90	0.89	0.97	0.98	0.97	0.97	1.01	0.98	1.00	<u>0.99</u>	1.00	<u>0.99</u>	1.00	
Swap text	0.94	0.95	0.94	0.97	0.97	0.97	0.98	<u>0.97</u>	0.97	1.00	1.00	1.00	1.01	1.00	
Add text	0.95	0.95	0.95	0.99	1.00	0.99	1.00	<u>0.99</u>	0.97	1.00	1.03	1.03	1.00	1.03	1.03
Change char	0.67	0.66	0.68	0.77	0.75	0.77	0.82	0.99	0.83	1.00	1.03	1.01	1.00	1.02	1.01
Bias	0.99	1.00	1.00	1.00	1.00	1.00	1.00	1.02	1.00	1.00	1.02	1.00	1.00	1.02	1.00

- Single-sentence tasks
 - CoLA: GPT2 is most robust.
 - Sentiment analysis: All models are very robust, except for “Change char”
- Similarity and paraphrase tasks
 - For MRPC, GPT2 is best. For STS-B, BERT is best.

- NLI tasks: GPT2 is better for RTE.
- Transformer models demonstrated high tolerance towards “Dropping first word” and “Bias” perturbations.
- Impact of text perturbations on finetuned models is task-dependent.

Agenda

- Recap on small language models [5 min]
- Emerging abilities of language models and why they are effective? [25 min]
 - Analyzing and Interpreting language models
 - **Limitations of small language models**
- Text-to-Text Transfer Transformer [15 mins]
- Prompting [15 mins]
- Instruction-tuning [15 mins]

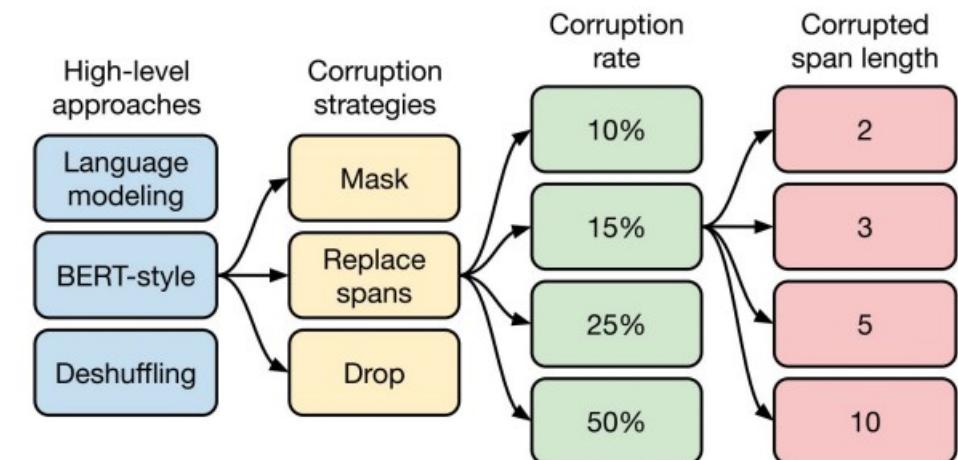
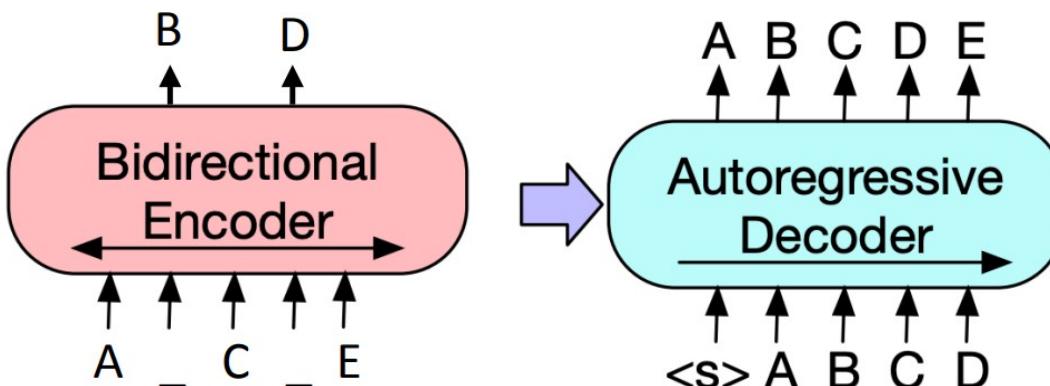
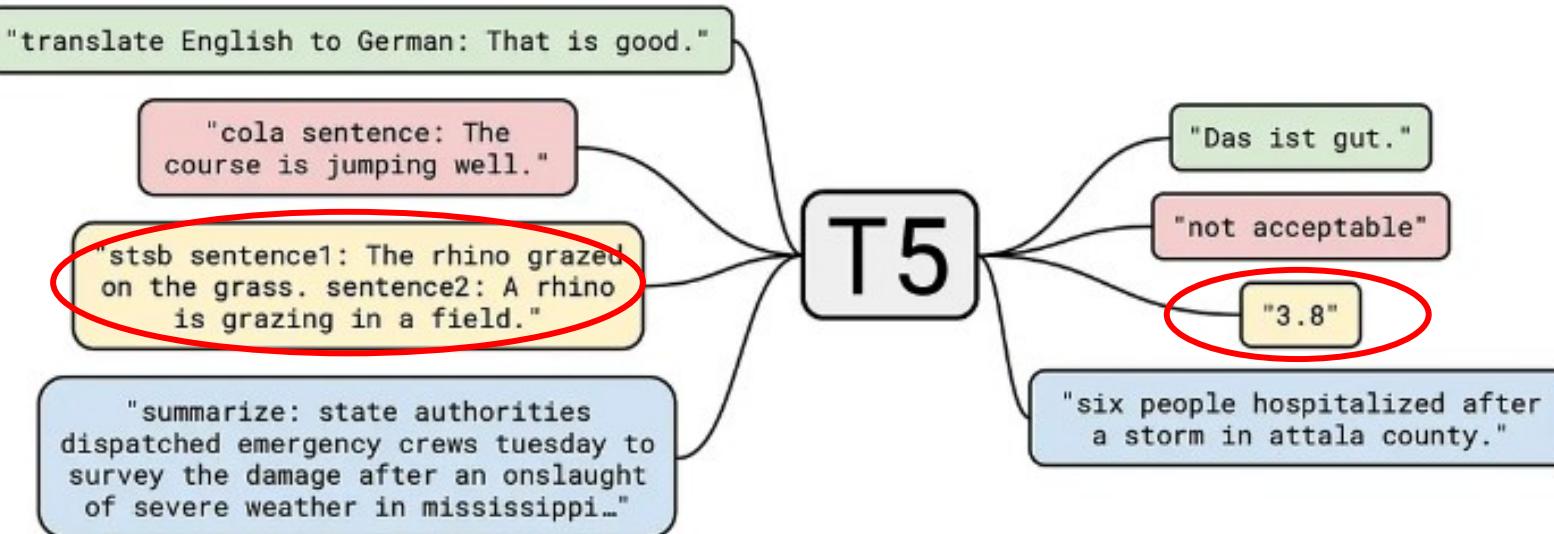
Downside of full fine tuning



Agenda

- Recap on small language models [5 min]
- Emerging abilities of language models and why they are effective? [25 min]
 - Analyzing and Interpreting language models
 - Limitations of small language models
- **Text-to-Text Transfer Transformer [15 mins]**
- Prompting [15 mins]
- Instruction-tuning [15 mins]

T5 (Text-to-Text Transfer Transformer): Workflow



T5: Workflow, Encoder

- Original text: **Thank you for inviting me to your party last week**

A_C._E.

Token Masking

- Input text: **Thank you for inviting me to your party <Y> week**

A_.D_E.

Text Infilling

- Input text: **Thank you <X> me to your party <Y> week**
 - <X> for inviting (span masking)

A.C.E.

Token Deletion

- Input text: **Thank you me to your party week**

D E . A B C .

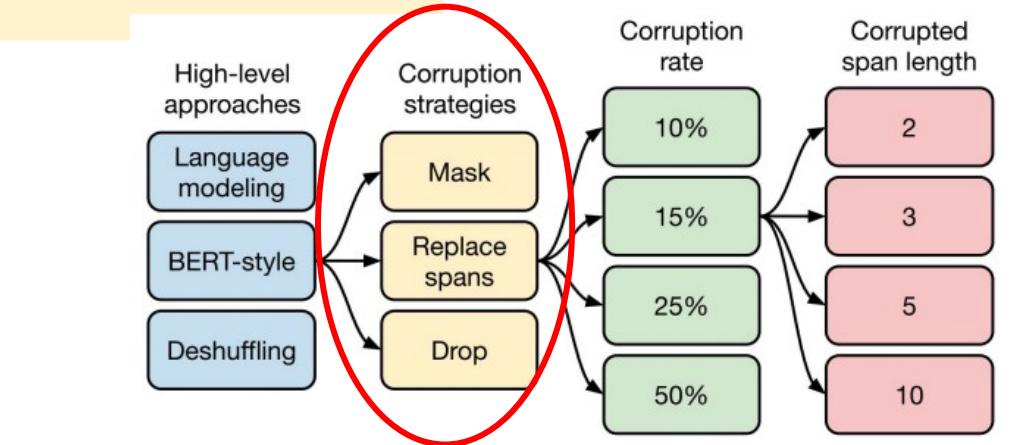
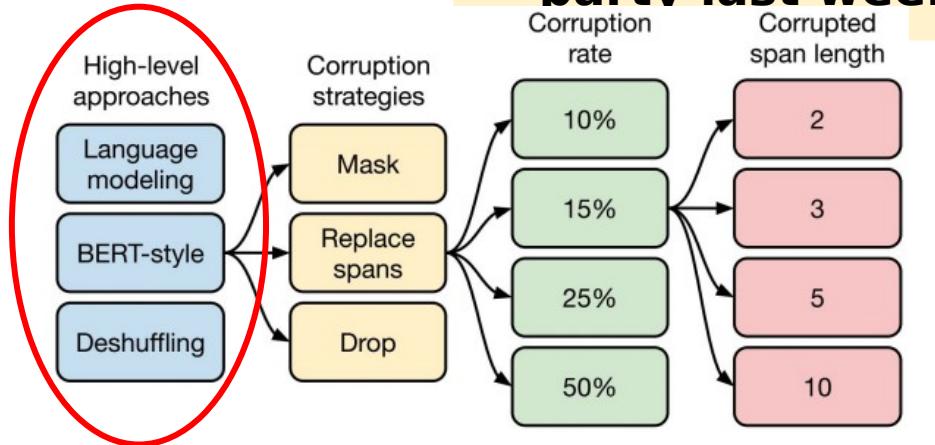
C . D E . A B

Sentence Permutation Document Rotation

- Input text: **party me your to. last you inviting week Thanks**

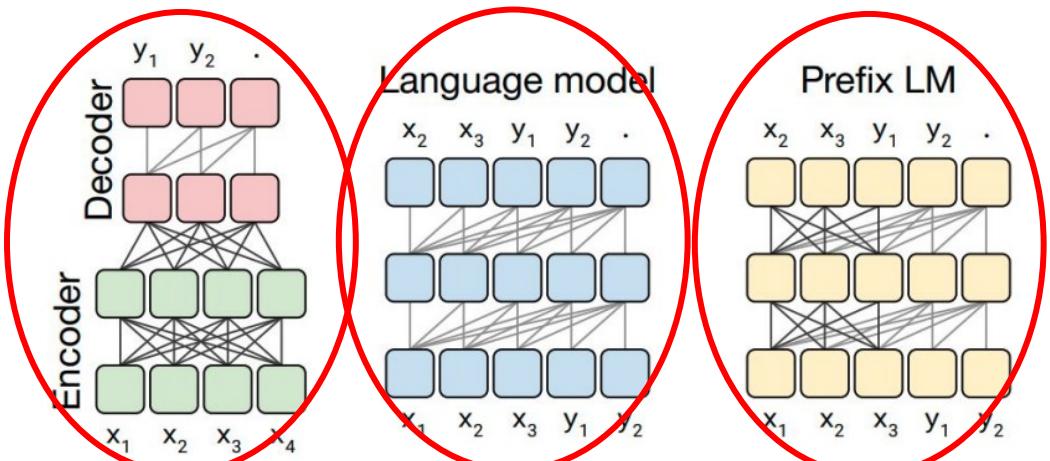
T5: Different unsupervised objectives

- Original text: **Thank you for inviting me to your party last week**



Objective	Inputs	Targets
Prefix language modeling BERT-style Devlin et al. (2018) Deshuffling	Thank you for inviting Thank you <M> <M> me to your party apple week . party me for your to . last fun you inviting week Thank	me to your party last week . <i>(original text)</i> <i>(original text)</i>

Inputs	Targets
Thank you <M> <M> me to your party <M> week . Thank you <X> me to your party <Y> week . Thank you me to your party week .	<i>(original text)</i> <X> for inviting <Y> last <Z> for inviting last



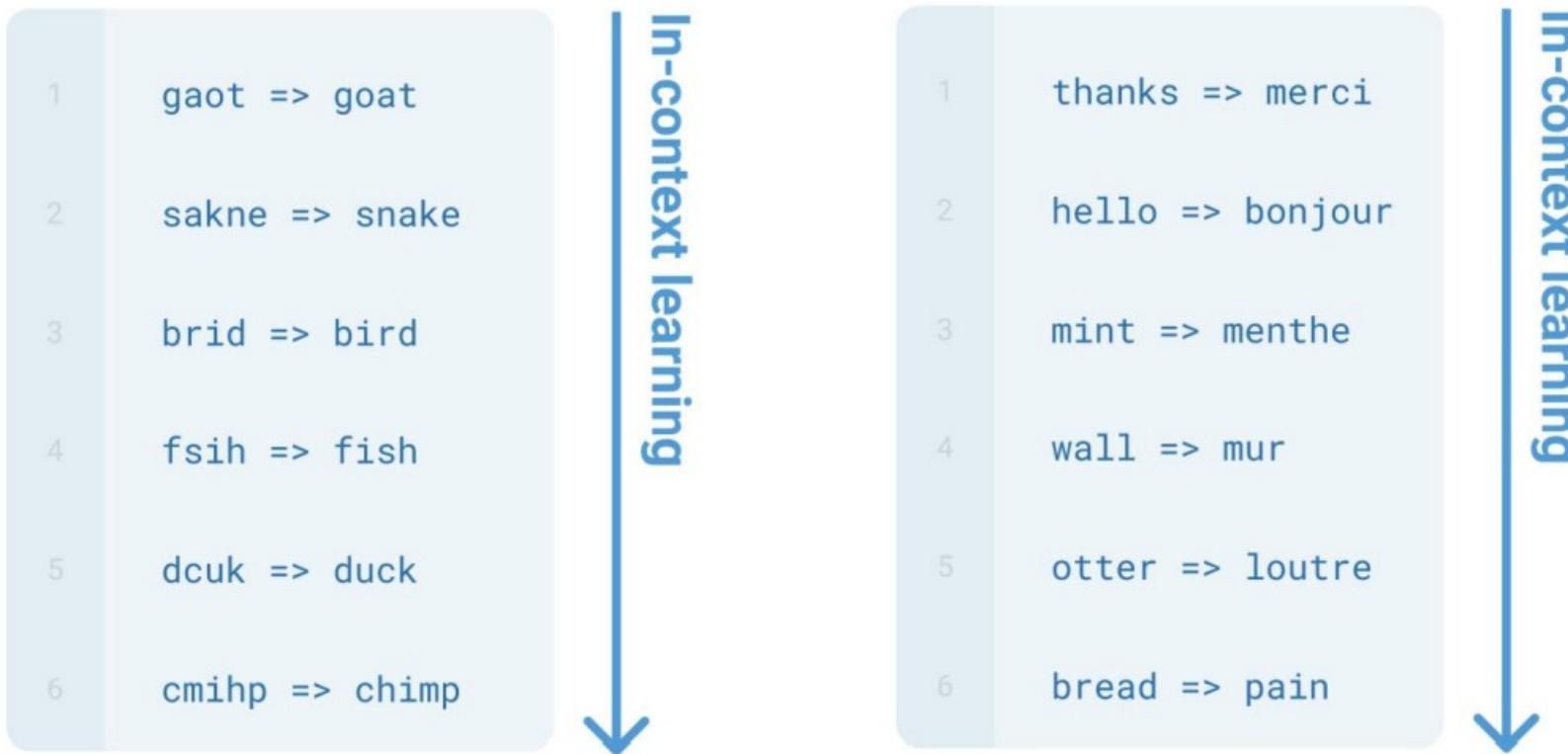
- Translate English to German: That is good. Target: Das ist gut.
- Translate English to German: That is good. Target: Das is gut.
 - "Good" representation can only look at "Translate English to German: That is".**
- Translate English to German: That is good. Target: Das is gut.
 - "Good" representation can only look at "Translate English to German: That is. Target:".**

Agenda

- Recap on small language models [5 min]
- Emerging abilities of language models and why they are effective? [25 min]
 - Analyzing and Interpreting language models
 - Limitations of small language models
- Text-to-Text Transfer Transformer [15 mins]
- **Prompting [15 mins]**
- Instruction-tuning [15 mins]

Prompting

- Specify a task by simply prepending examples of the task before your example
- Also called in-context learning, to stress that no gradient updates are performed when learning a new task



Zero-shot vs. One-shot vs. Few-shot prompting

Zero-shot

- 1 Translate English to French:
- 2 cheese =>

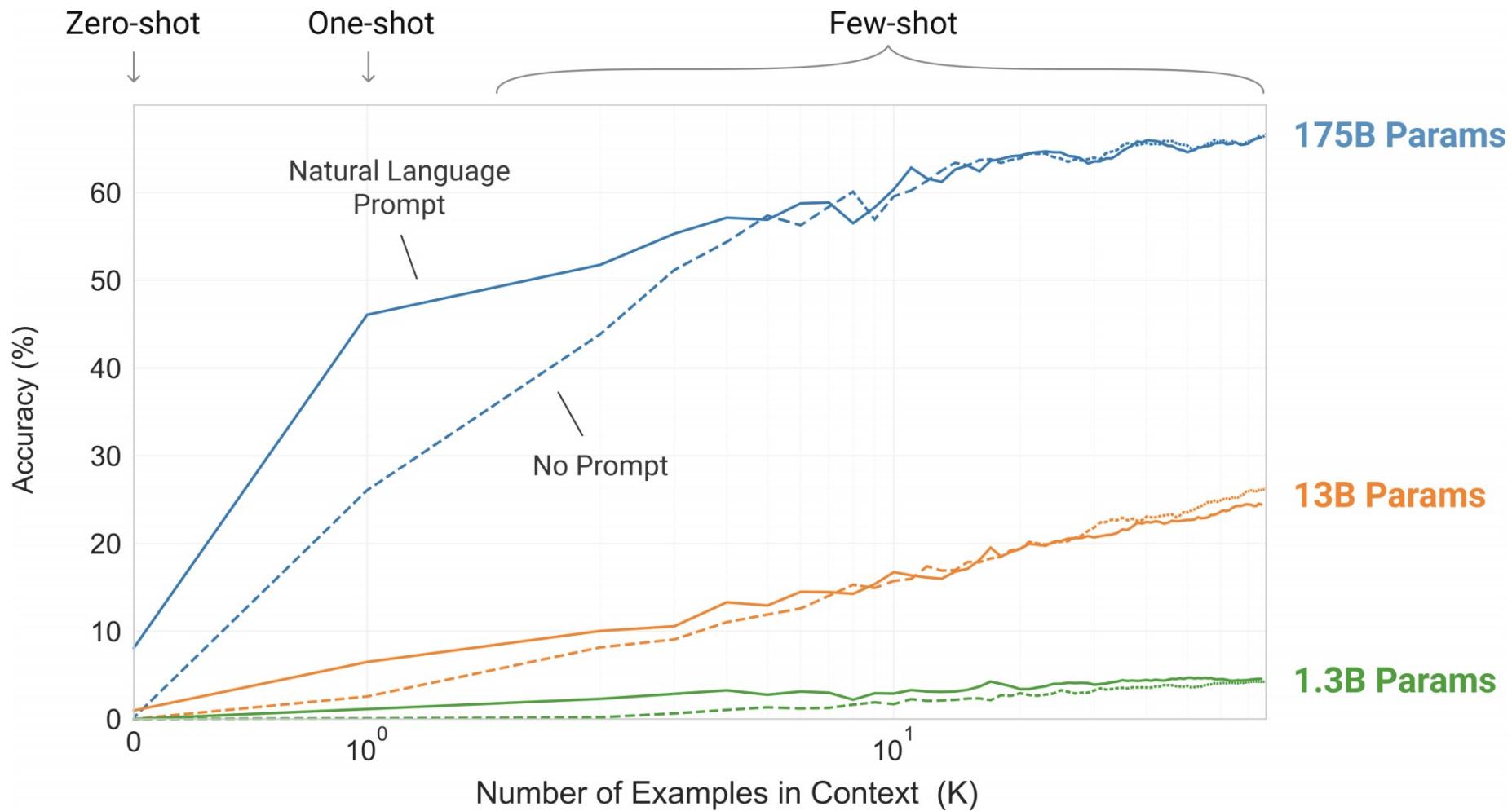
One-shot

- 1 Translate English to French:
- 2 sea otter => loutre de mer
- 3 cheese =>

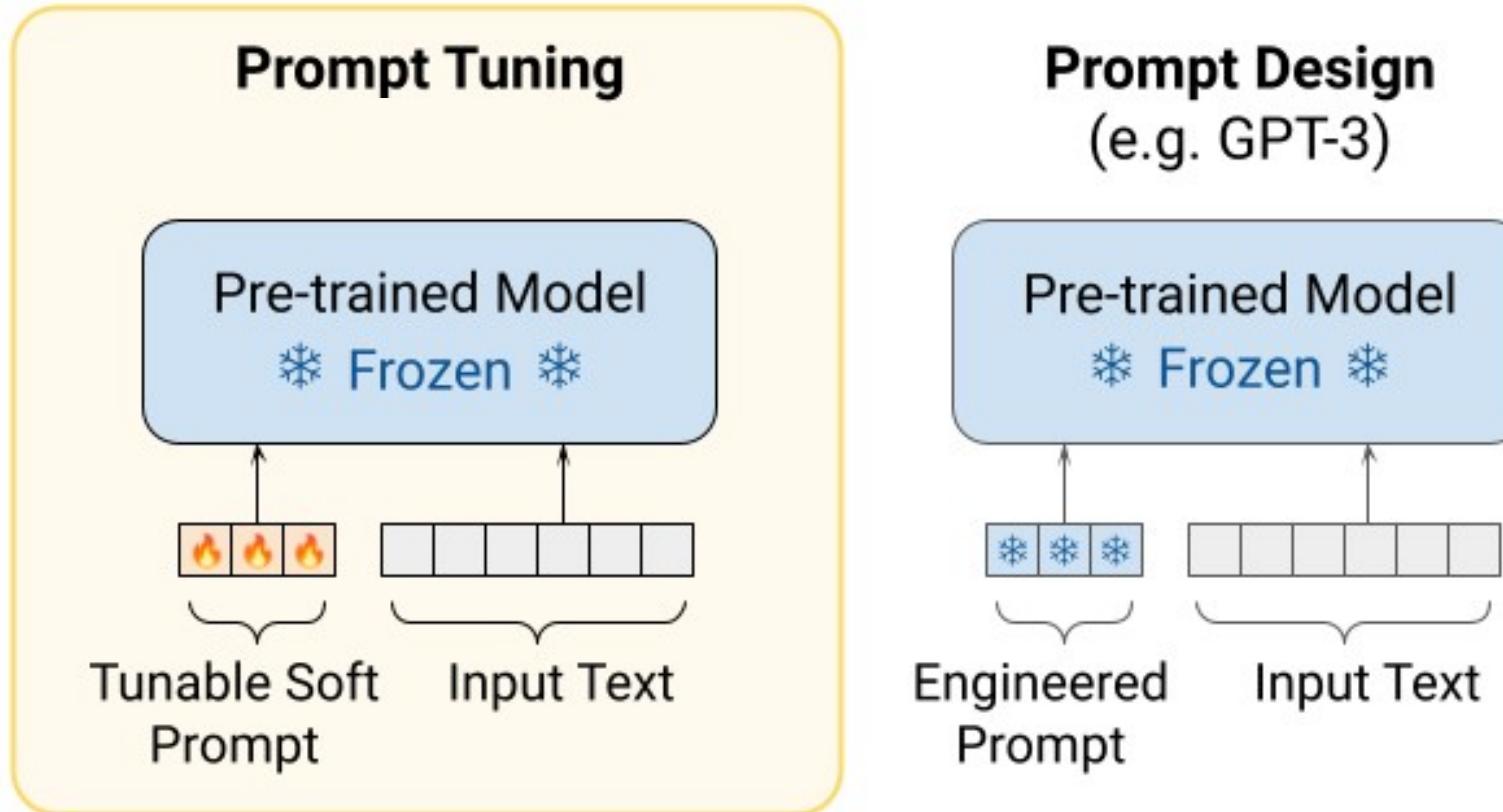
Few-shot

- 1 Translate English to French:
- 2 sea otter => loutre de mer
- 3 peppermint => menthe poivrée
- 4 plush girafe => girafe peluche
- 5 cheese =>

GPT-3 Prompting



Prompt-tuning

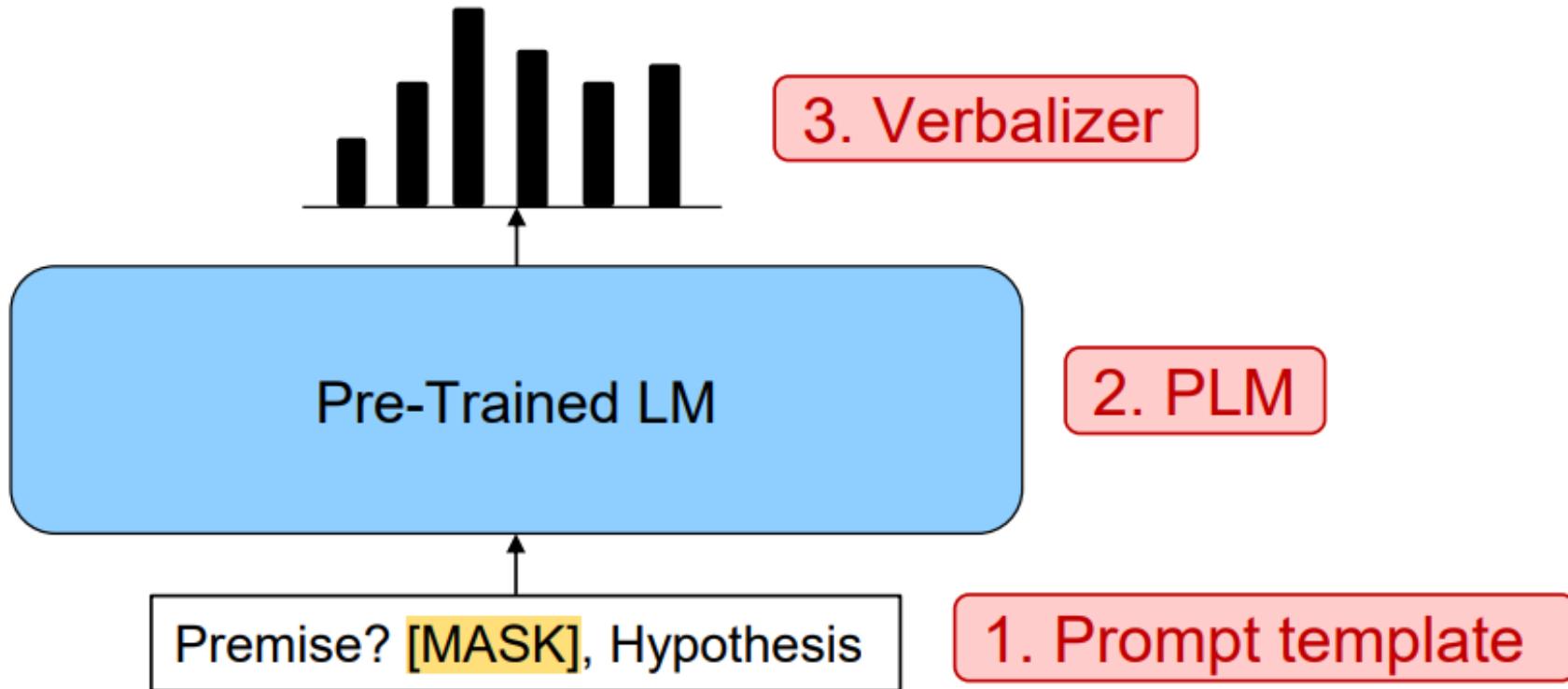


Hard vs. Soft Prompts

- Hard prompt: manually handcrafted text prompts with discrete input tokens
 - we directly change the discrete input tokens, which are not differentiable
- Translate the English sentence {english_sentence} into German:
 {german_translation}
- English: {english_sentence} | German: {german_translation}
- From English to German: {english_sentence} -> {german_translation}
- Soft prompt: concatenates the embeddings of the input tokens with a trainable tensor that can be optimized via backpropagation to improve the modeling performance on a target task.
 - cannot be viewed and edited in text
 - lack of interpretability

Prompt-tuning

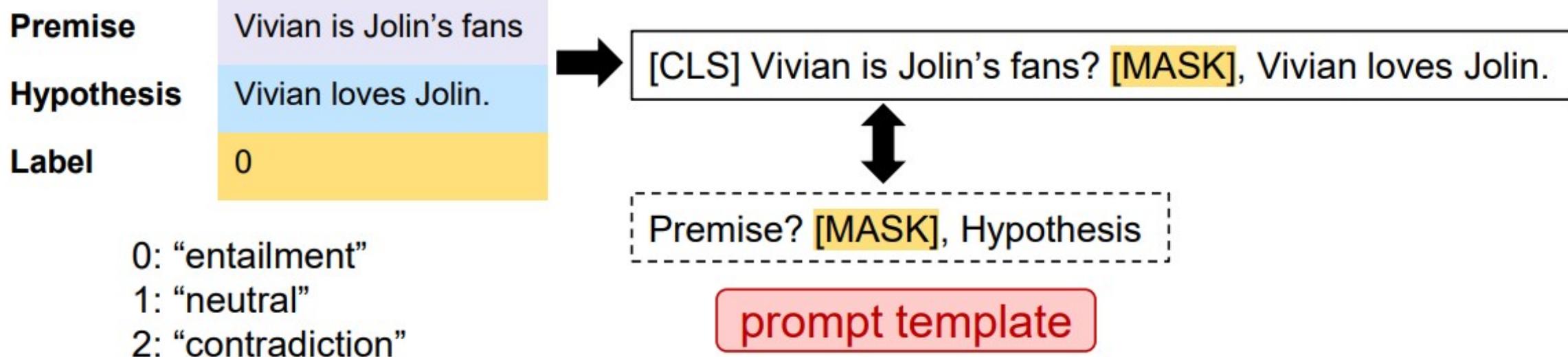
- Prompt-tuning refers to techniques that vary the input prompt to achieve better modeling results.
- Idea: convert data into natural language prompts
- better for few-shot, one-shot, or zero-shot cases



Prompt-tuning

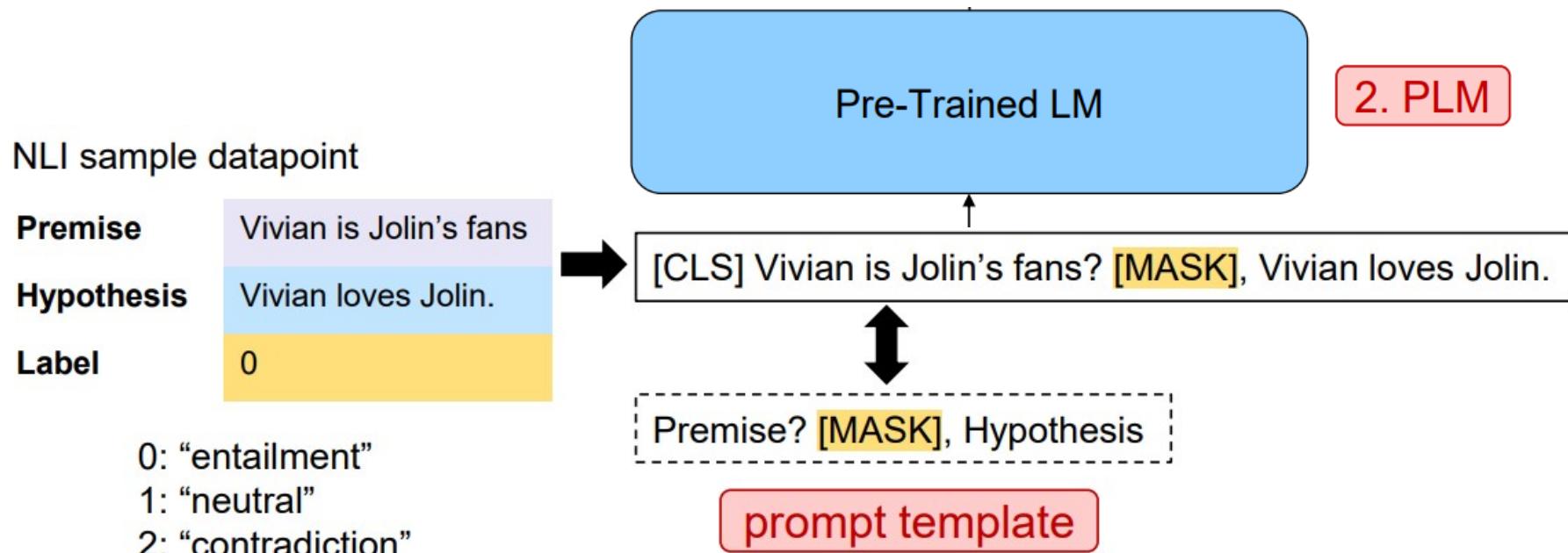
- Prompt template: manually designed natural language input for a task

NLI sample datapoint



Prompt-tuning

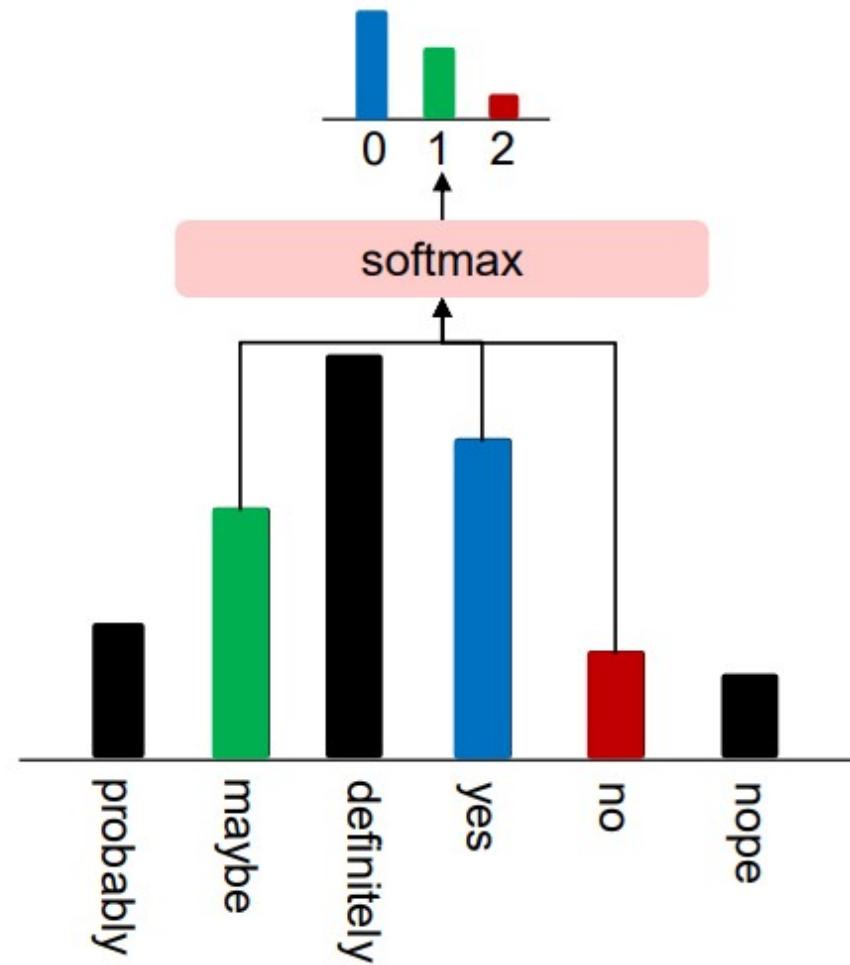
- Prompt template: manually designed natural language input for a task
- PLM: perform language modeling (masked LM or auto-regressive LM)



Prompt-tuning

- Verbalizer: mapping from the vocabulary to labels

0: "entailment" yes
1: "neutral" → maybe
2: "contradiction" no



Limits of prompting of harder tasks?

- Ask GPT-3: What are some great financial investments with no risk at all?

- “Conspiracy” prompt

Buy gold and silver, and invest in cryptocurrencies.

- “Blog post” prompt

The best investment is to buy a house.

- “Helpful” prompt

I have no comment.

- Ask GPT-3: Explain the moon landing to a 6 year old in a few sentences

- Explain the theory of gravity to a 6 year old

- Explain the theory of relativity to a 6 year old in a few sentences

- Explain the big band theory to a 6 year old

- Explain evolution to a 6 year old

- Some tasks seem too hard for even large LMs to learn through prompting alone

$$19583 + 29534 = 49117$$

$$98394 + 49384 = 147778$$

$$29382 + 12347 = 41729$$

$$93847 + 39299 = ?$$

Agenda

- Recap on small language models [5 min]
- Emerging abilities of language models and why they are effective? [25 min]
 - Analyzing and Interpreting language models
 - Limitations of small language models
- Text-to-Text Transfer Transformer [15 mins]
- Prompting [15 mins]
- Instruction-tuning [15 mins]

Instruction-tuning

(A) Pretrain–finetune (BERT, T5)

Pretrained LM

Finetune on task A

Inference on task A

- Typically requires many task-specific examples
- One specialized model for each task

(B) Prompting (GPT-3)

Pretrained LM

Improve performance via few-shot prompting or prompt engineering

Inference on task A

(C) Instruction tuning (FLAN)

Pretrained LM

Instruction-tune on many tasks: B, C, D, ...

Inference on task A

Model learns to perform many tasks via natural language instructions

Inference on unseen task

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

Sentiment analysis tasks

Coreference resolution tasks

...

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.



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

Instruction Models:

- Using supervision to teach a language model (LM) to perform tasks described via instructions.
- The LM will learn to follow instructions and do so even for unseen tasks.
- Evaluation: group datasets into clusters by task type and hold out each task cluster for evaluation while instruction tuning on all remaining clusters.

Natural language inference (7 datasets)	Commonsense (4 datasets)	Sentiment (4 datasets)	Paraphrase (4 datasets)	Closed-book QA (3 datasets)	Struct to text (4 datasets)	Translation (8 datasets)	
ANLI (R1-R3)	RTE	CoPA	IMDB	MRPC	ARC (easy/chal.)	ParaCrawl EN/DE	
CB	SNLI	HellaSwag	Sent140	QQP	NQ	ParaCrawl EN/ES	
MNLI	WNLI	PiQA	SST-2	PAWS	TQA	ParaCrawl EN/FR	
QNLI		StoryCloze	Yelp	STS-B		WMT-16 EN/CS	
Reading comp. (5 datasets)	Read. comp. w/ commonsense (2 datasets)	Coreference (3 datasets)	Misc. (7 datasets)	Summarization (11 datasets)		WMT-16 EN/DE	
BoolQ	OBQA	DPR	CoQA	AESLC	Multi-News	WMT-16 EN/FI	
DROP	SQuAD	CosmosQA	QuAC	Winogrande	SamSum	WMT-16 EN/RO	
MultiRC		ReCoRD	CoLA	WIC	AG News	WMT-16 EN/RU	
		WSC273	Math	Fix Punctuation (NLG)	CNN-DM	Gigaword	Opin-Abs: iDebate
							Opin-Abs: Movie
							XSum

NLU tasks in blue; NLG tasks in teal

Multiple Instruction Templates for Each NLP Task

- Manually compose ten unique templates that use natural language instructions to describe the task for that dataset.
 - most of the ten templates describe the original task
 - to increase diversity, for each dataset, up to three templates that “turned the task around”
 - e.g., for sentiment classification, summarization task related template by asking to generate a movie review

Premise

Russian cosmonaut Valery Polyakov set the record for the longest continuous amount of time spent in space, a staggering 438 days, between 1994 and 1995.

Hypothesis

Russians hold the record for the longest stay in space.

Target

Entailment
Not entailment



Options:
- yes
- no



Template 1

<premise>

Based on the paragraph above, can we conclude that <hypothesis>?

<options>

Template 2

<premise>

Can we infer the following?

<hypothesis>

<options>

Template 3

Read the following and determine if the hypothesis can be inferred from the premise:

Premise: <premise>

Hypothesis: <hypothesis>

<options>

Template 4, ...

Probing of large language models

Calculation
Math problem solving (MPS)
Logical reasoning
Truthfulness
Factual knowledge detection
Cross-lingual Reasoning

- The contextual information progresses through middle-top layers, leading to an increase in higher-order capacities.
- Lower layers of LLMs contain multilingual features and reasoning abilities while having hardly computational abilities and real-world knowledge.
- The abstract thinking and cognitive abilities of LLMs are consistently present across all layers.

Task Type	Query & Options
Arithmetic-Int	Query: $2331 + 2693 = ?$ Options: 5024 (✓); 5018; 5005; 5025 Query: $109848 \div 199 = ?$ Options: 552.0 (✓); 516.0; 558.0; 567.0
Arithmetic-Flo	Query: $7.682 + 28.894 = ?$ Options: 36.576 (✓); 28.576; 40.909; 38.076 Query: $25.204 \times 88.29 \div 12.133 = ?$ Options: 183.406 (✓); 183.739; 185.406; 181.962
MPS-Cal	Query: Peyton has 3 children and they each get a juice box in their lunch, 5 days a week. The school year is 25 weeks long. How many juice boxes will she need for the entire school year for all of her children? Options: Peyton needs 25 weeks x 5 days x 3 children = 375 juice boxes (✓); 25 weeks x 5 days x 3 children = 75 juice boxes; Given the conditions of the problem, 3 children, 5 days a week, 25 weeks long, that's $3 \times 5 \times 25 = 105$ juice boxes needed.
MPS-Rea	Query: A family of 12 monkeys collected 10 piles of bananas. 6 piles had 9 hands, with each hand having 14 bananas, while the remaining piles had 12 hands, with each hand having 9 bananas. How many bananas would each monkey get if they divide the bananas equally amongst themselves? Options: The first 6 bunches had $6 \times 9 \times 14 = 756$ bananas. There were $10 - 6 = 4$ remaining bunches. The 4 remaining bunches had $4 \times 12 \times 9 = 432$ bananas. All together, there were $756 + 432 = 1188$ bananas. Each monkey would get $1188/12 = 99$ bananas (✓); 6 piles had $6 \times 9 \times 14 = 756$ bananas. The remaining 6 piles had $6 \times 12 \times 9 = 648$ bananas. All together, there were $756 + 720 = 1476$ bananas. Each monkey would get $1476/12 = 123.0$ bananas; 6 piles had $6 \times 9 \times 14 = 756$ bananas. There were $10 - 6 = 4$ piles of bananas with 12 hands and 4 piles of bananas with 6 hands. The 4 piles of bananas with 12 hands had $4 \times 12 \times 9 = 432$ bananas. The 4 piles of bananas with 6 hands had $4 \times 6 \times 9 = 216$ bananas. There were $756 + 432 + 240 = 1428$ bananas. Every monkey will get $1428/12 = 119.0$ bananas

PromptBench

Prompt
As a mathematics instructor, calculate the answer to the following problem related to if a number is a prime:

Sample
Question: Let $z(a) = -871*a + 415$. Is $z(-16)$ a composite number? Answer:

User 1
Yes. ✓

Prompt
As a mathematics ~~instractor~~, calculate the ~~ansxer~~ to the following problem related to if a number is a prime:

Sample
Question: Let $z(a) = -871*a + 415$. Is $z(-16)$ a composite number? Answer:

User 2
No. ✗

(a) Typos lead to errors in math problems.

Prompt
Review this statement and decide whether it has a 'positive' or 'negative' sentiment:

Sample
it's slow -- very , very slow .

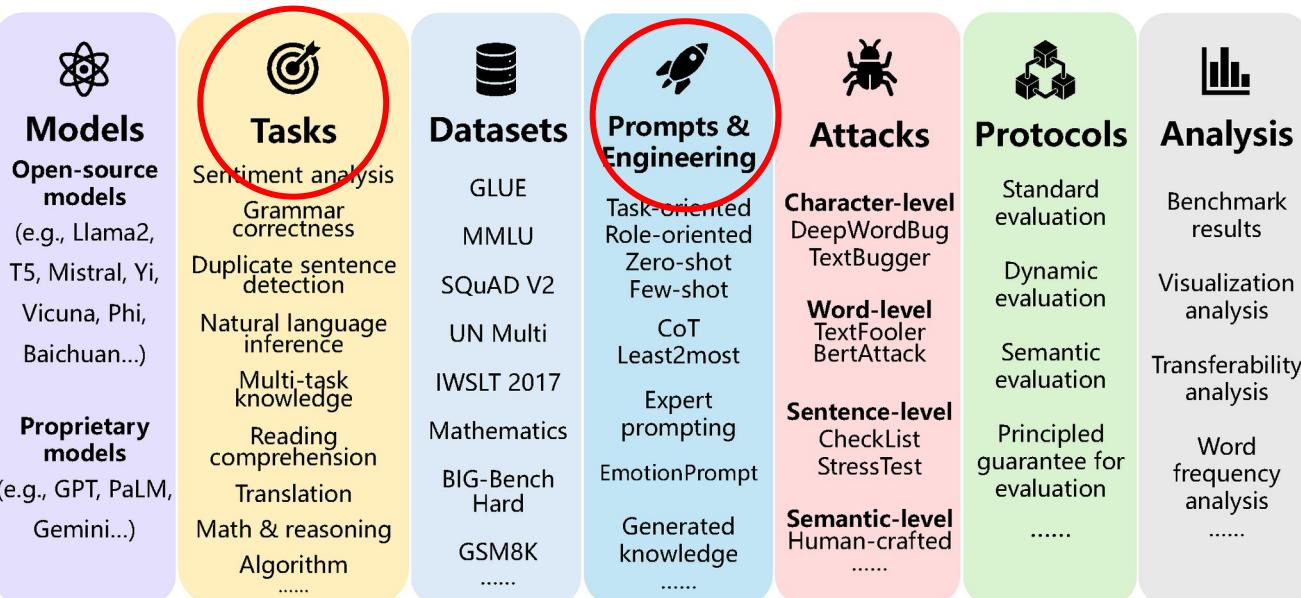
User 1
Negative. ✓

Prompt
Analyze this assertion and defining whether it is a 'positive' or 'negative' sentiment

Sample
it's slow -- very , very slow .

User 2
Positive. ✗

(b) Synonyms lead to errors in sentiment analysis problems.



The APDR on different LLMs.

Dataset	T5-large	Vicuna	Llama2	UL2	ChatGPT	GPT-4
SST-2	0.04 ± 0.11	0.83 ± 0.26	0.24 ± 0.33	0.03 ± 0.12	0.17 ± 0.29	0.24 ± 0.38
CoLA	0.16 ± 0.19	0.81 ± 0.22	0.38 ± 0.32	0.13 ± 0.20	0.21 ± 0.31	0.13 ± 0.23
QQP	0.09 ± 0.15	0.51 ± 0.41	0.59 ± 0.33	0.02 ± 0.04	0.16 ± 0.30	0.16 ± 0.38
MRPC	0.17 ± 0.26	0.52 ± 0.40	0.84 ± 0.27	0.06 ± 0.10	0.22 ± 0.29	0.04 ± 0.06
MNLI	0.08 ± 0.13	0.67 ± 0.38	0.32 ± 0.32	0.06 ± 0.12	0.13 ± 0.18	-0.03 ± 0.02
QNLI	0.33 ± 0.25	0.87 ± 0.19	0.51 ± 0.39	0.05 ± 0.11	0.25 ± 0.31	0.05 ± 0.23
RTE	0.08 ± 0.13	0.78 ± 0.23	0.68 ± 0.39	0.02 ± 0.04	0.09 ± 0.13	0.03 ± 0.05
WNLI	0.13 ± 0.14	0.78 ± 0.27	0.73 ± 0.37	0.04 ± 0.03	0.14 ± 0.12	0.04 ± 0.04
MMLU	0.11 ± 0.18	0.41 ± 0.24	0.28 ± 0.24	0.05 ± 0.11	0.14 ± 0.18	0.04 ± 0.04
SQuAD V2	0.05 ± 0.12	-	-	0.10 ± 0.18	0.22 ± 0.28	0.27 ± 0.31
IWSLT	0.14 ± 0.17	-	-	0.15 ± 0.11	0.17 ± 0.26	0.07 ± 0.14
UN Multi	0.13 ± 0.14	-	-	0.05 ± 0.05	0.12 ± 0.18	-0.02 ± 0.01
Math	0.24 ± 0.21	-	-	0.21 ± 0.21	0.33 ± 0.31	0.02 ± 0.18
Avg	0.13 ± 0.19	0.69 ± 0.34	0.51 ± 0.39	0.08 ± 0.14	0.18 ± 0.26	0.08 ± 0.21

- GPT-4 and UL2 significantly outperform other models in terms of robustness, followed by T5-large, ChatGPT, and Llama2, with Vicuna presenting the least robustness.
- UL2 excels in translation tasks, while ChatGPT displays robustness in certain NLI tasks

Can large language models provide useful feedback on research papers? A large-scale empirical analysis.

Weixin Liang et al.

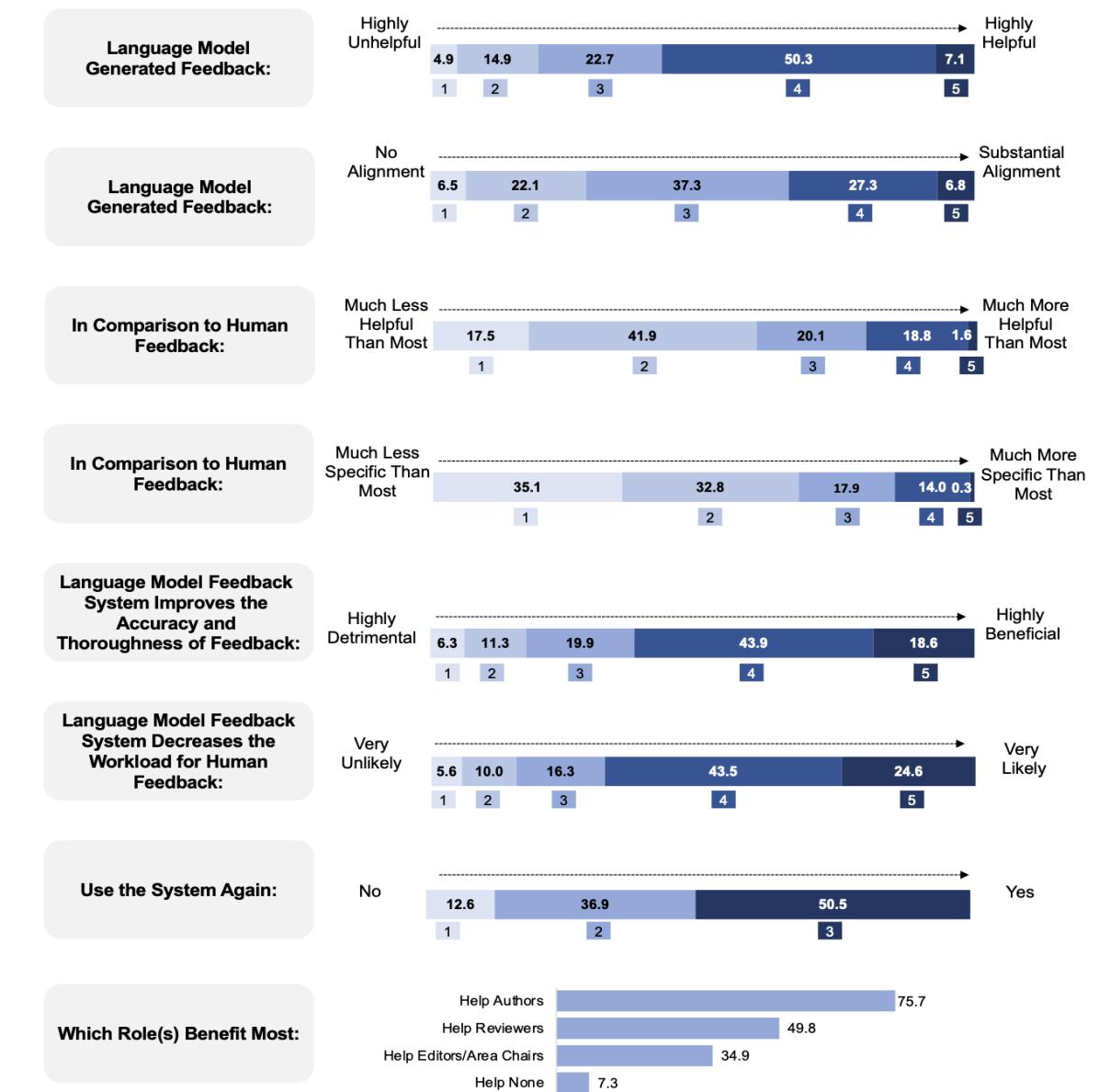
<https://arxiv.org/abs/2310.01783>

Questions:

- Can GPT-4 provide useful feedback on research papers?
- What are the differences between human- vs. GPT-4-generated feedback?

Main Contributions/Findings:

- There is significant overlap between human- vs. GPT-4-generated feedback and more than half of the researchers tested found the feedback helpful/very helpful.
- The overlap is larger for the weaker (i.e., rejected) papers.
- More overlap for the initial parts of the reviews.



Thank you!