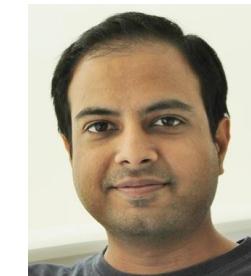


# Recent trends in large language models

Mounika Marreddy<sup>1</sup>, Subba Reddy Oota<sup>2</sup>, Manish Gupta<sup>3</sup>, Lucie Flek<sup>1</sup>

<sup>1</sup>University of Bonn, Germany; <sup>2</sup>Inria Bordeaux, France; <sup>3</sup>Microsoft, India

[mmarredd@uni-bonn.de](mailto:mmarredd@uni-bonn.de), [subba-reddy.oota@inria.fr](mailto:subba-reddy.oota@inria.fr), [gmanish@microsoft.com](mailto:gmanish@microsoft.com), [flek@bit.uni-bonn.de](mailto:flek@bit.uni-bonn.de)



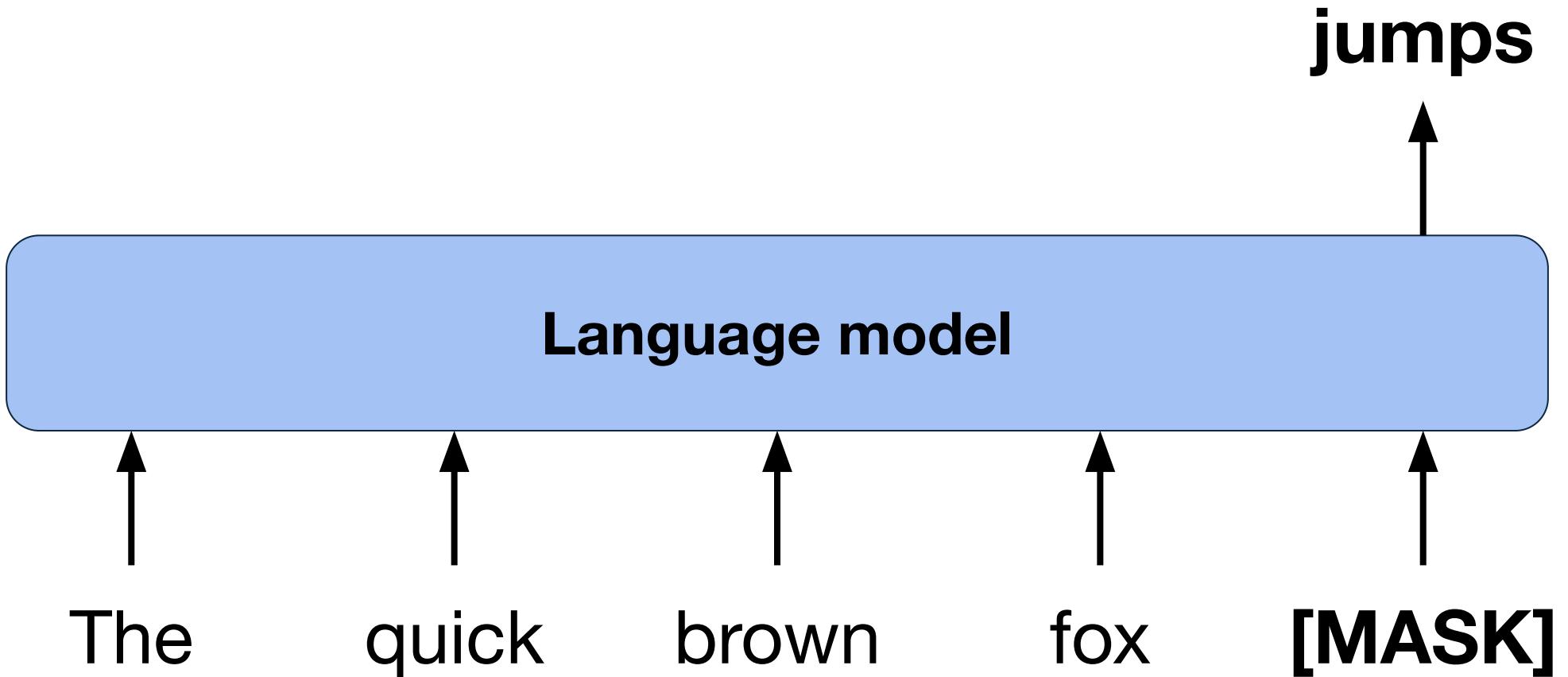
# Agenda

- Recap on small language models [5 min]
- Emerging abilities of language models and why they are effective? [25 min]
- Text-to-Text Transfer Transformer, Prompting, Instruction-tuning [1 hour]
- Lunch [45 min]
- Large language models, Chain-of-Thought reasoning, RAG [1 hour 30 min]
- Coffee break [15 min]
- Hands on session [1 hour]

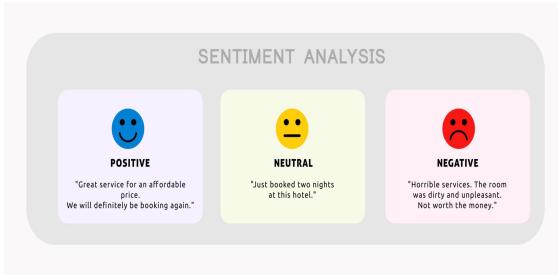
# Agenda

- Recap on small language models [5 min]
- Emerging abilities of language models and why they are effective? [25 min]
- Text-to-Text Transfer Transformer, Prompting, Instruction-tuning [1 hour]
- Lunch [45 min]
- Large language models, Chain-of-Thought reasoning, RAG [1 hour 30 min]
- Coffee break [15 min]
- Hands on session [1 hour]

# LMs are trained to predict missing words



# Language models are everywhere



Small language models



MMLU (Massive Multitask Language Understanding)

BIG-bench A small icon of a wooden chair.

<https://gluebenchmark.com/>

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

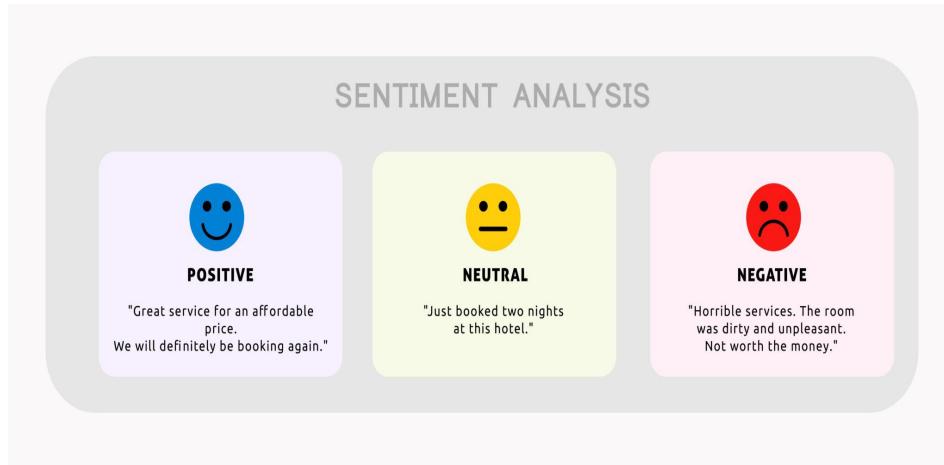
<https://crfm.stanford.edu/helm/lite/latest/>

<https://paperswithcode.com/dataset/mmlu>

<https://paperswithcode.com/dataset/big-bench>

# 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?

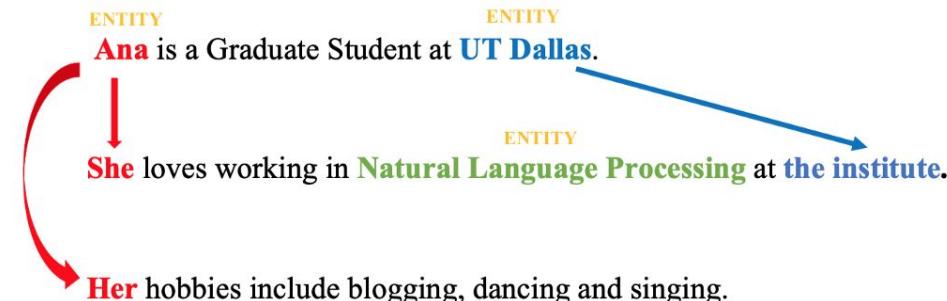
**Answer Candidate**

gravity

Summarization

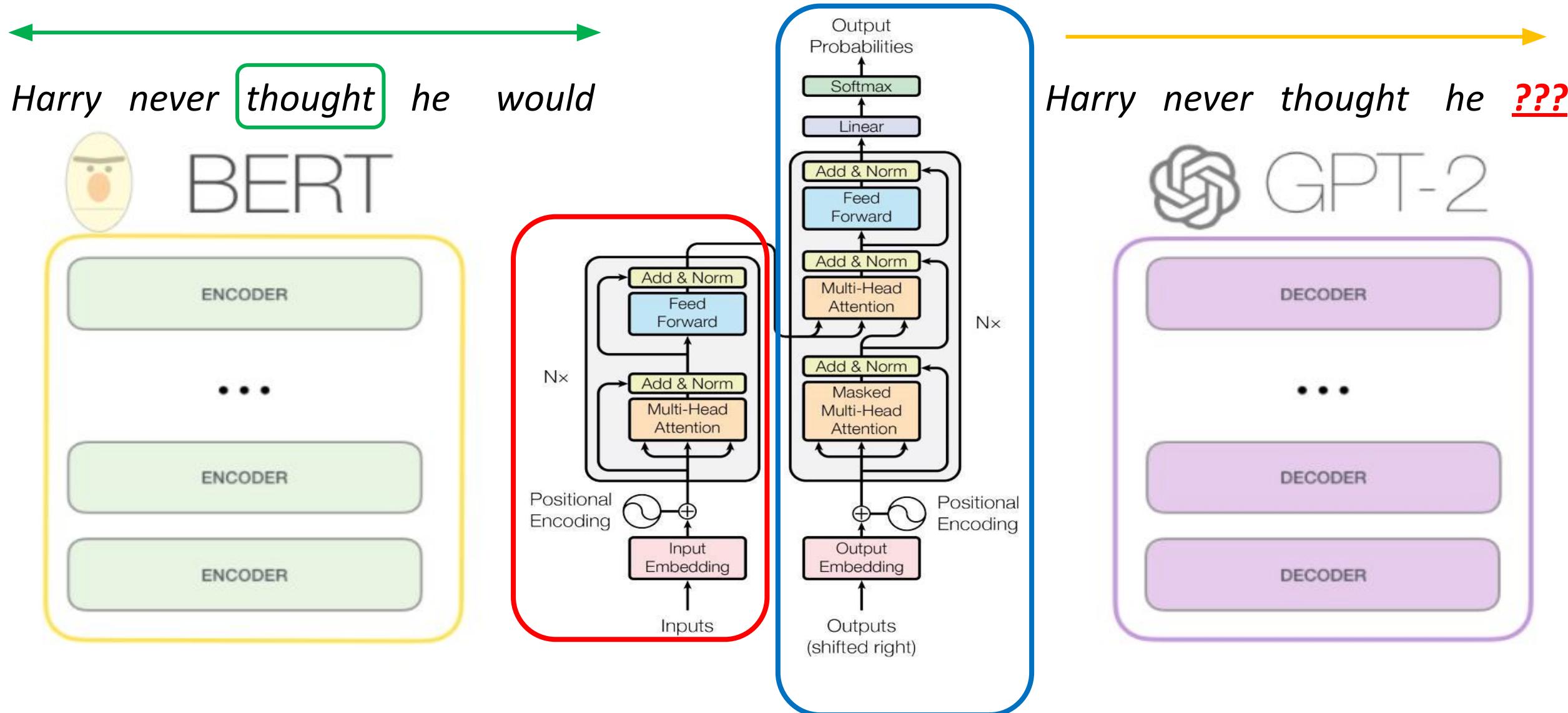
**Input Article**

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. ...



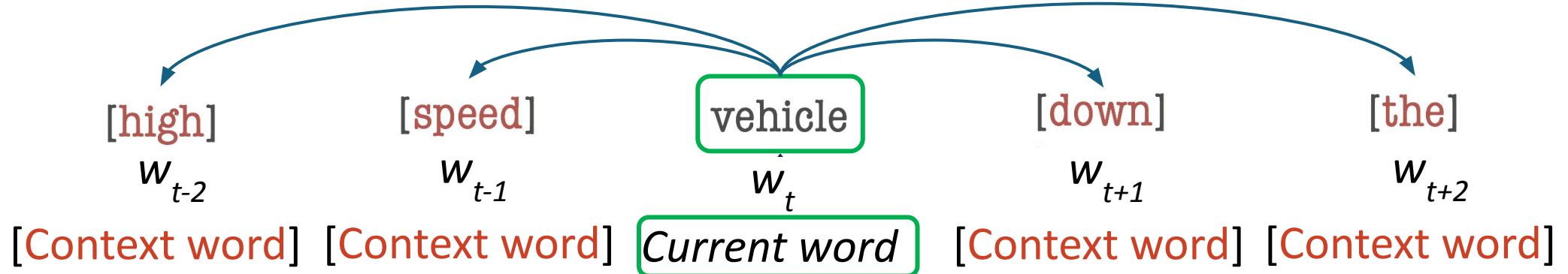
Coreference Resolution

# Transformer

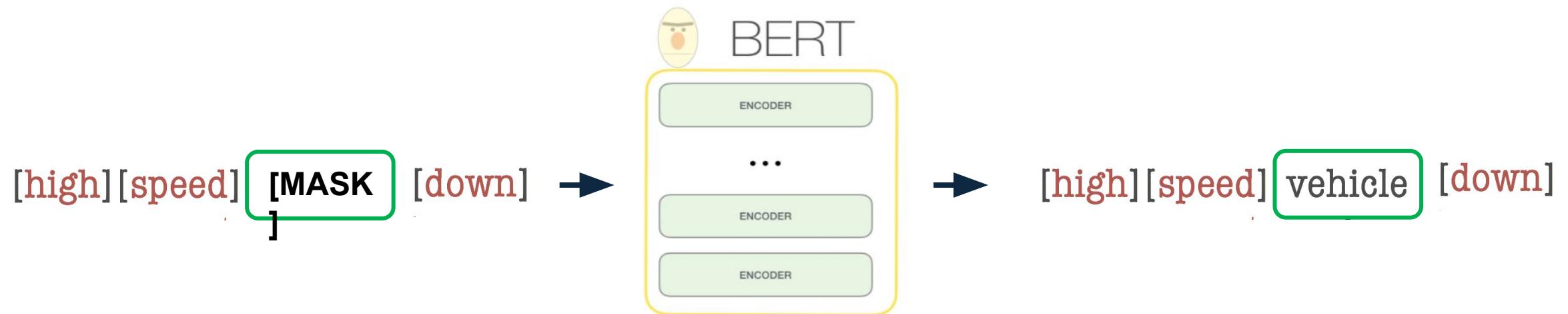


# BERT: Workflow

## 1. Self-attention with Bi-directional context

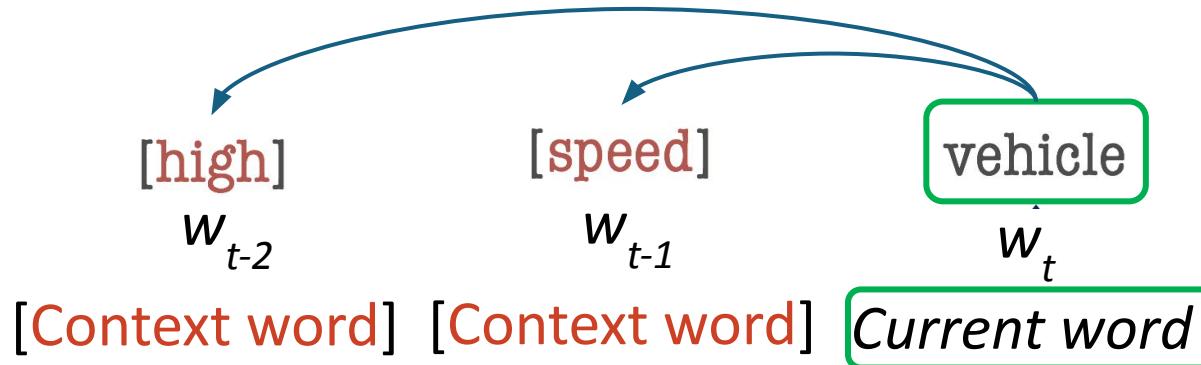


## 2. Masked language modeling (MLM)



# GPT2: Workflow

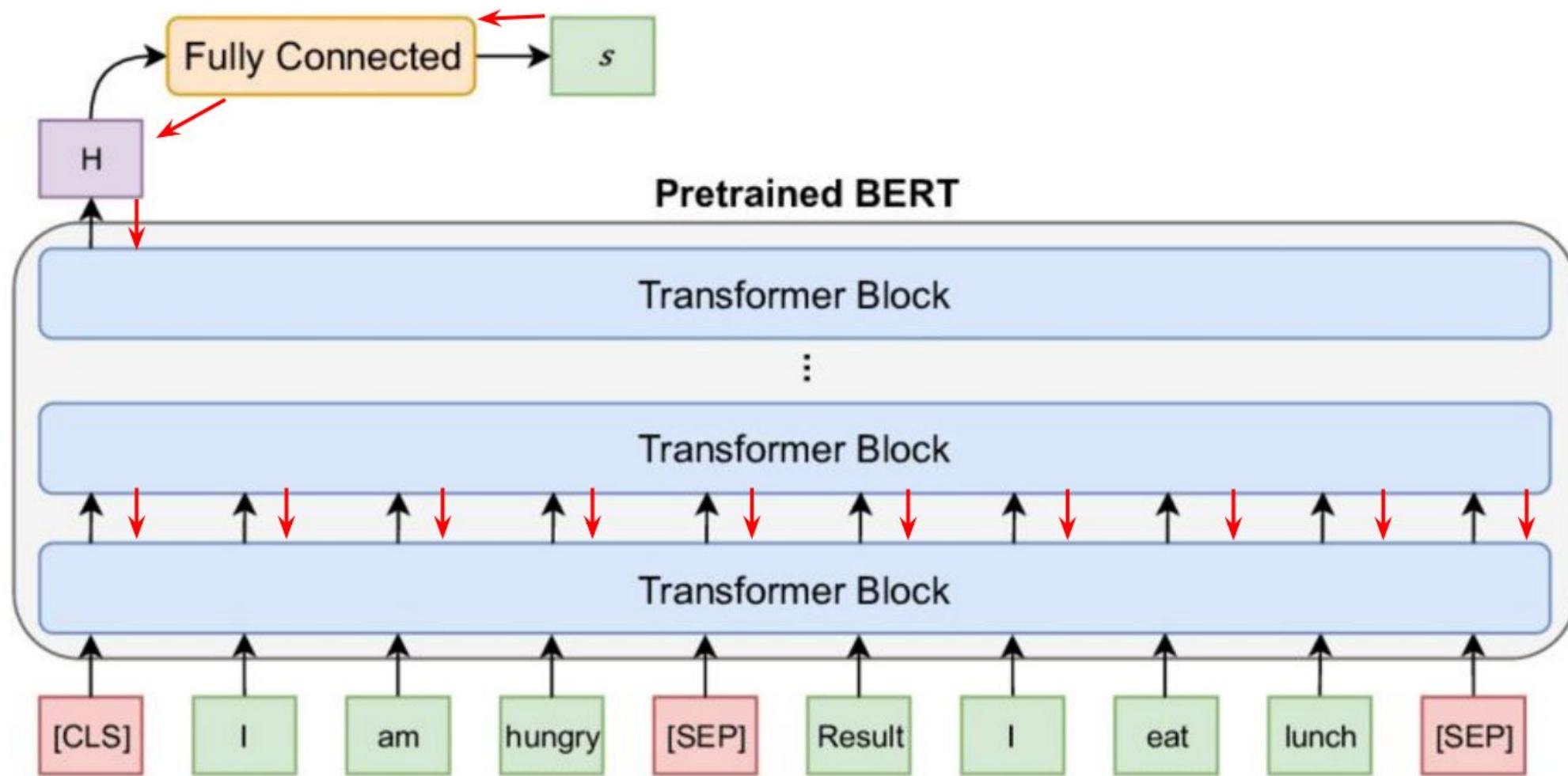
## 1. Self-attention with Uni-directional context



## 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]
  - Analyzing and Interpreting language models
- Text-to-Text Transfer Transformer, Prompting, Instruction-tuning [1 hour]
- Lunch [45 min]
- Large language models, Chain-of-Thought reasoning, RAG [1 hour 30 min]
- Coffee break [15 min]
- Hands on session [1 hour]

# 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, Prompting, Instruction-tuning [1 hour]
- Lunch [45 min]
- Large language models, Chain-of-Thought reasoning, RAG [1 hour 30 min]
- Coffee break [15 min]
- Hands on session [1 hour]

# Analyzing and Interpreting LMs

Model  
Interpretability

**Analysis of representations  
via probing**

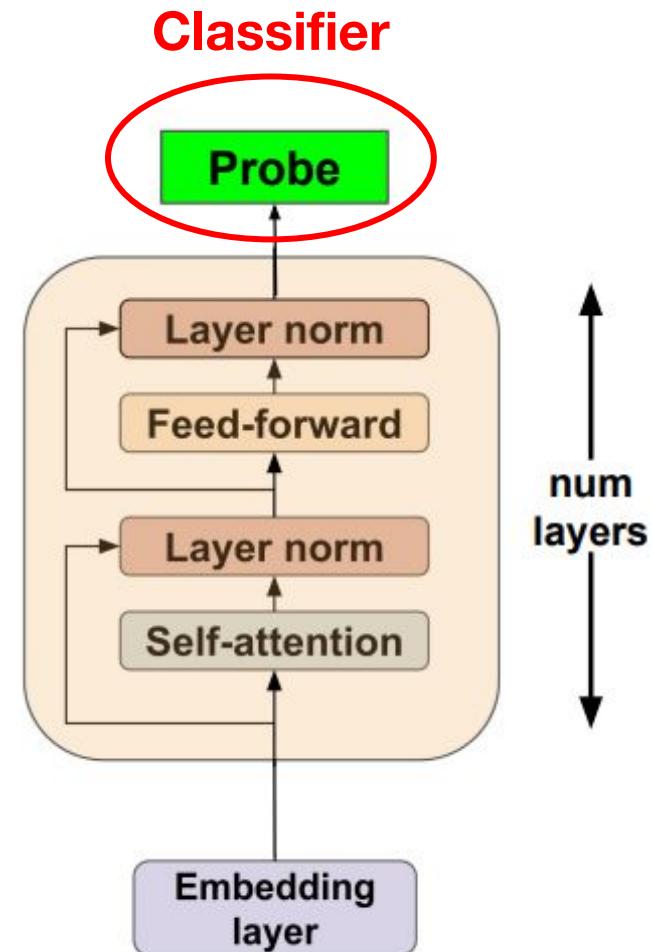
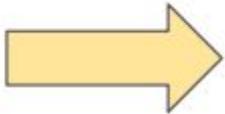
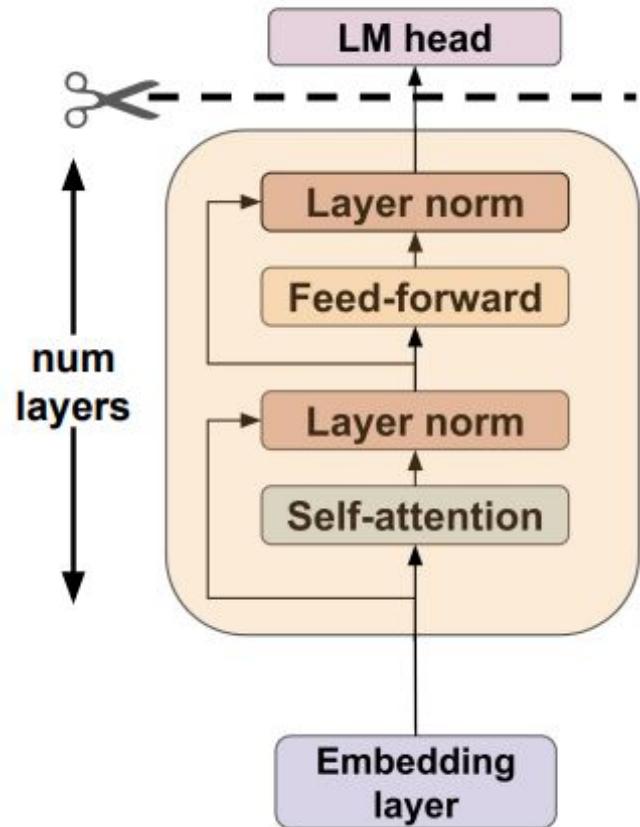
Behavioural  
Interpretability

**Bridging language models  
and brain**

Mechanistic  
Interpretability

**Robustness  
Reverse engineering of  
neural computations**

# Model Interpretability?



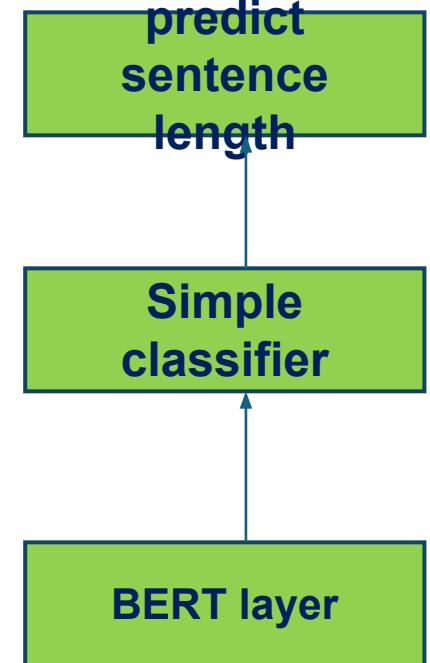
**BERTology**



# 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.

*If the prediction accuracy is good, then the model might be capturing the sentence length feature*

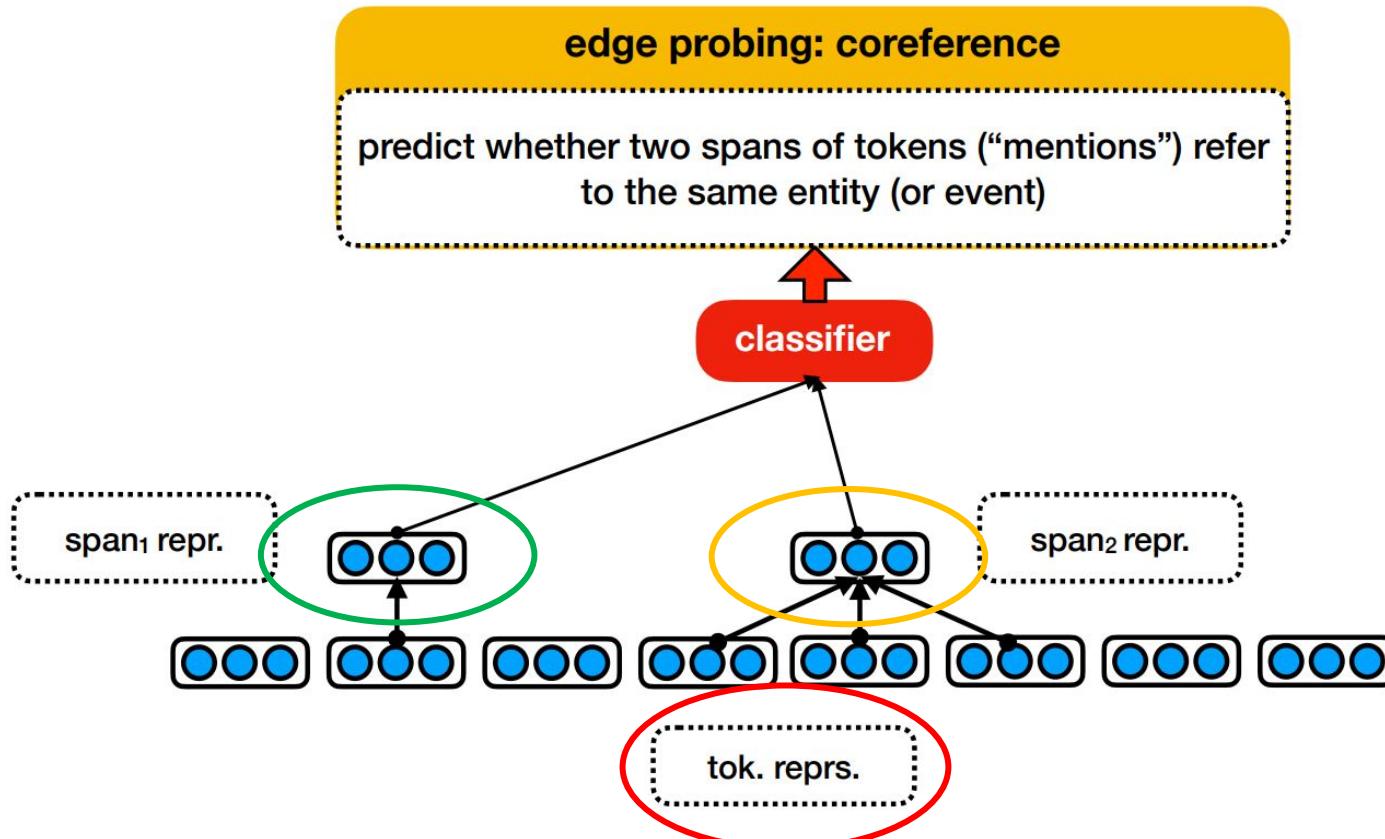


# BERT composes a hierarchy of linguistic signals ranging from surface to semantic features

	<b>Surface</b>	<b>Syntactic</b>			<b>Semantic</b>					
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	<b>96.2 (3.9)</b>	<b>66.5 (66.0)</b>	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)	<b>69.8 (69.6)</b>	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)	<b>41.3 (13.0)</b>	83.3 (36.6)	82.9 (32.9)	89.8 (27.6)	<b>88.1 (21.9)</b>	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)	<b>84.1 (39.5)</b>	83.0 (32.9)	89.9 (27.5)	87.4 (22.2)	<b>82.2 (21.1)</b>	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)	<b>38.5 (10.8)</b>	83.1 (39.8)	<b>87.0 (37.1)</b>	<b>90.0 (28.0)</b>	87.6 (22.9)	81.8 (20.5)	63.4 (13.4)	<b>78.7 (28.9)</b>
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)	<b>65.2 (15.3)</b>	74.9 (25.4)

Jawahar et al. 2019 ACL

# Edge Probing

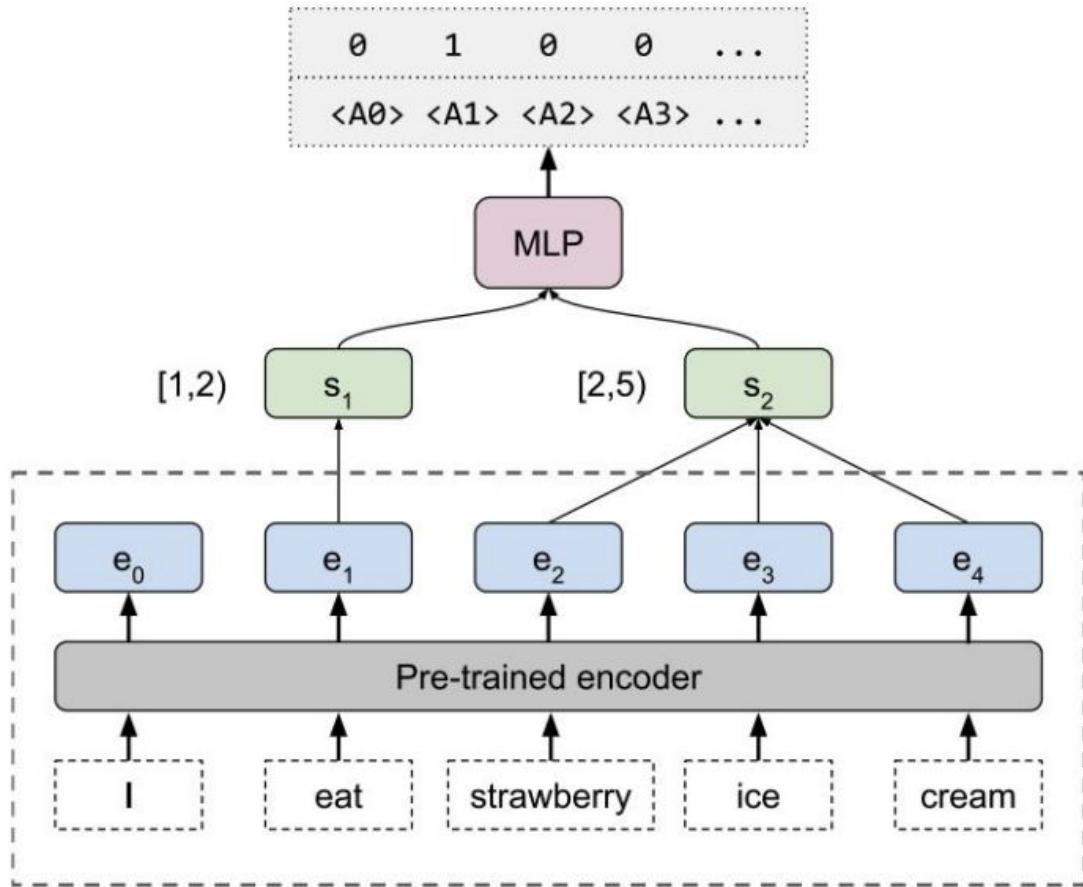


# Edge Probing

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

# Edge Probing

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



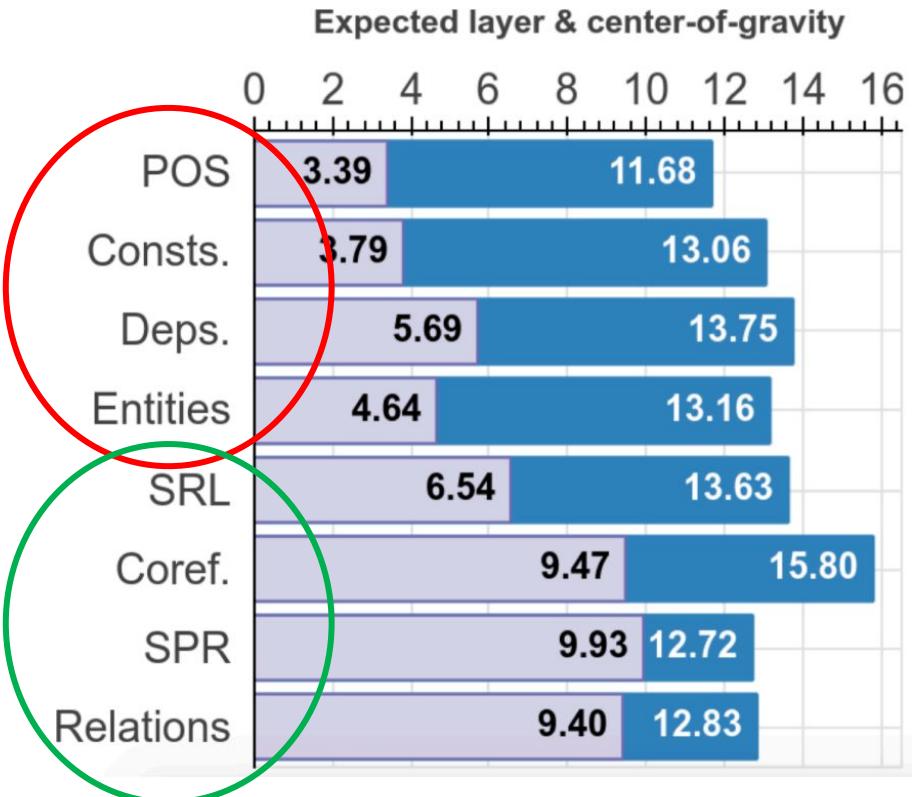
Labels

Binary classifiers

Span representations

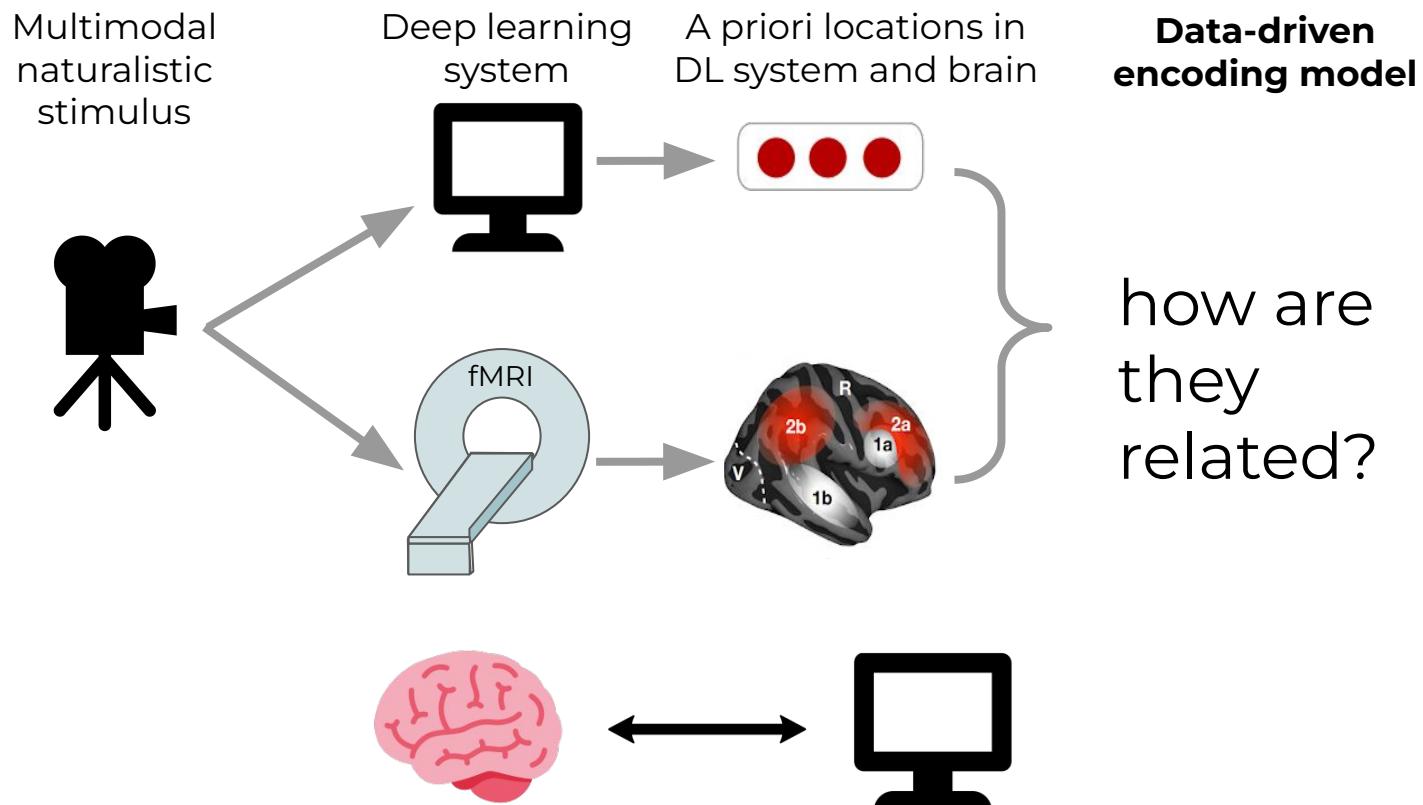
Contextual vectors

Input tokens

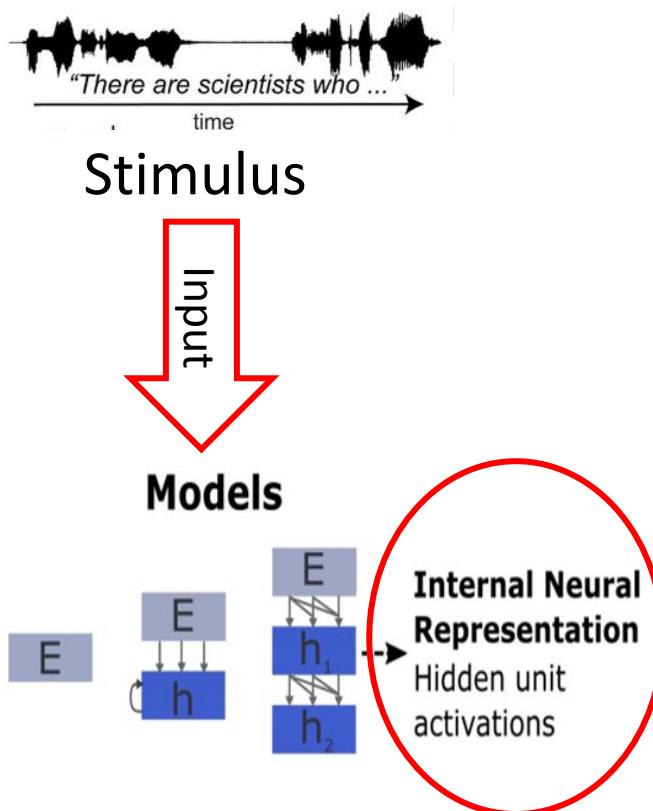


# Behavioural Interpretability

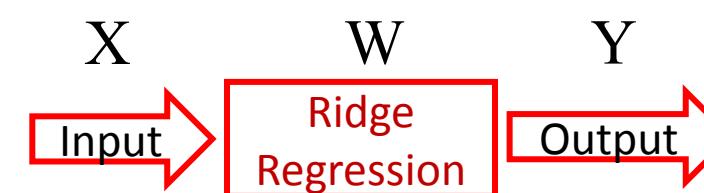
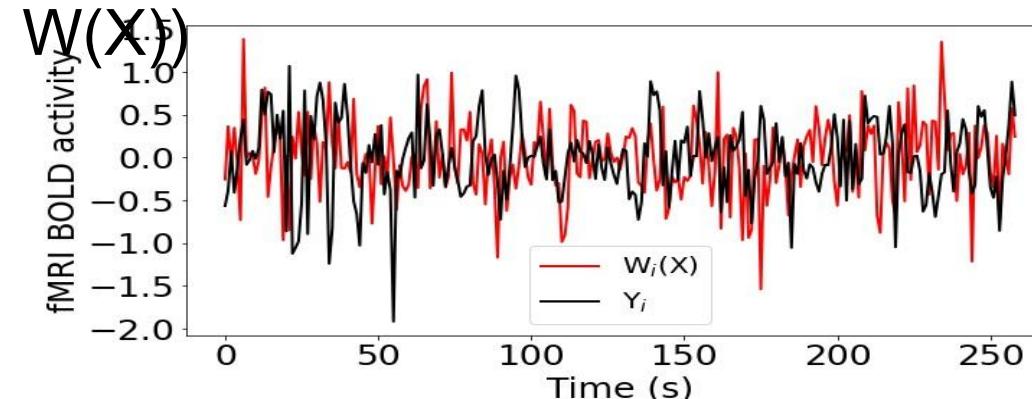
- Data-driven encoding models evaluate the relationships between brains and deep learning models



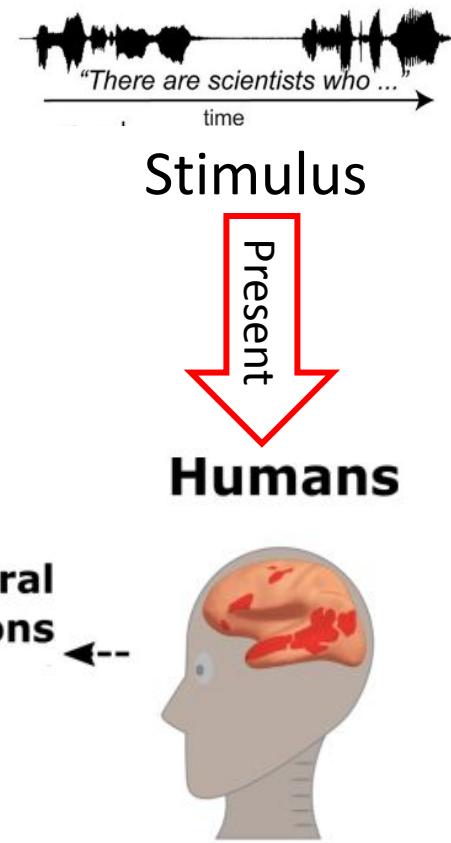
# Brain Encoding?



Pearson Correlation ( $R$ ) =  $\text{Corr}(Y,$

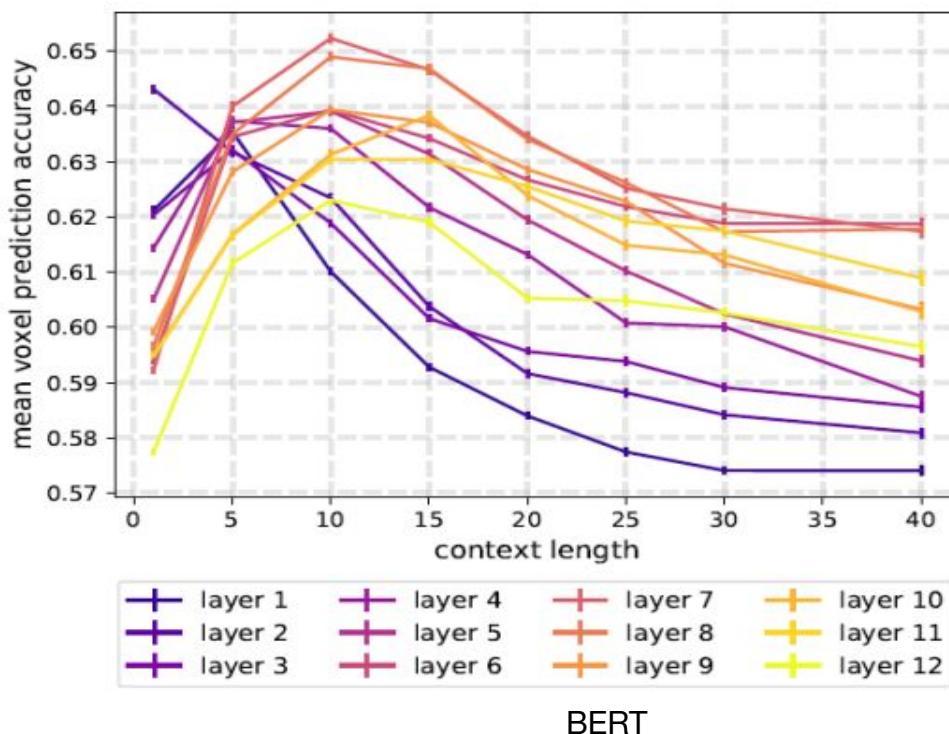


Internal Neural Representations  
fMRI

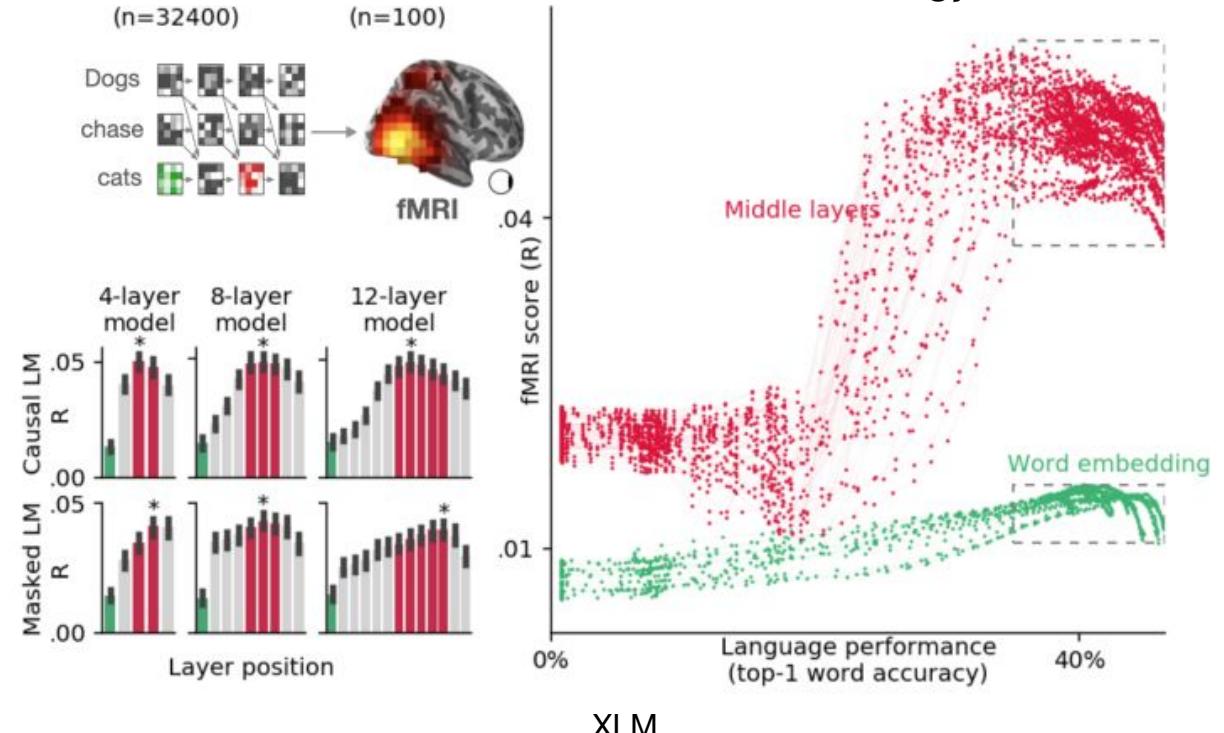


# The strongest alignment with high-level language brain regions has consistently been observed in middle layers

Toneva et al. NeurIPS-2019



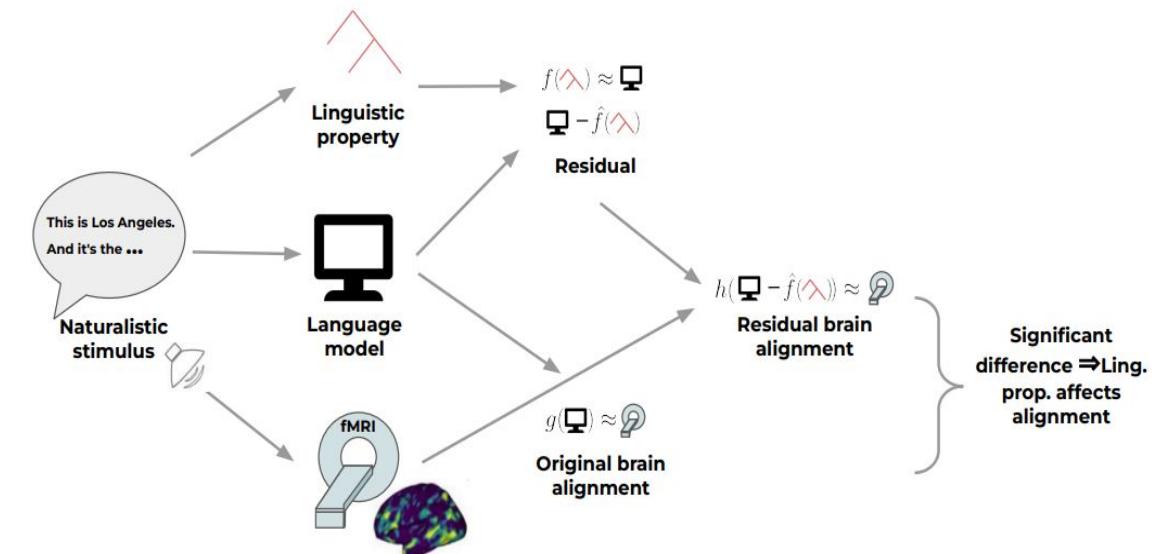
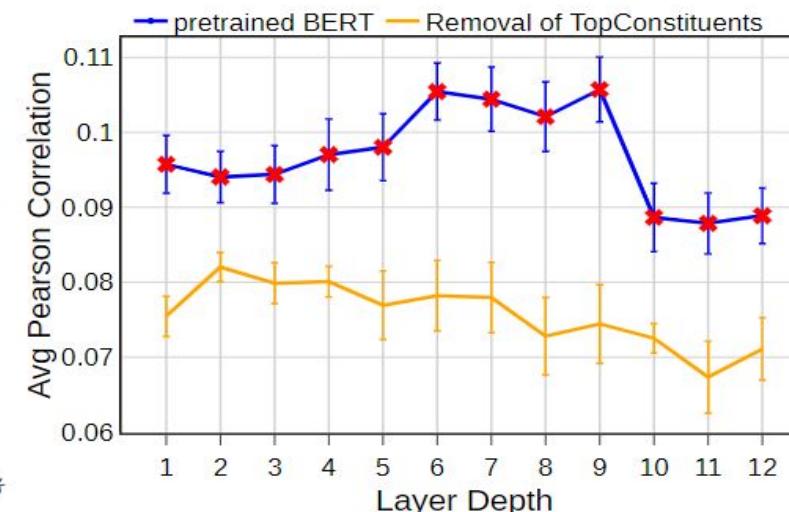
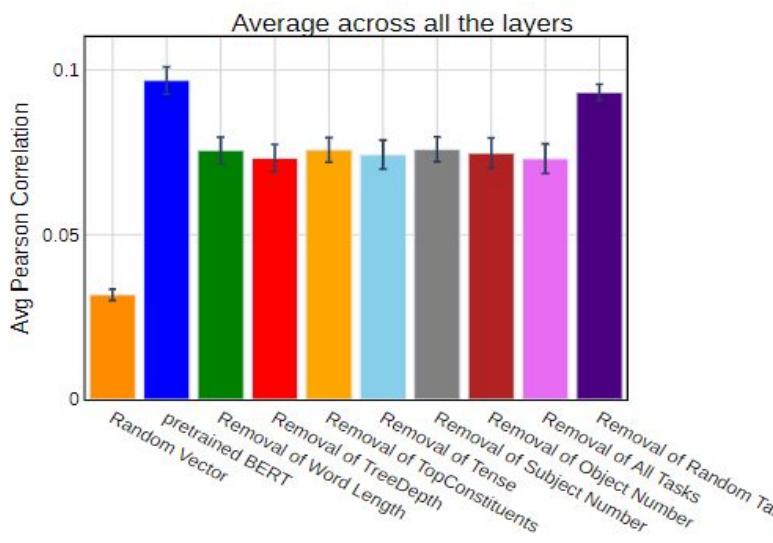
Caucheteux et al. Communication Biology 2022



Across several types of large NLP systems, best alignment with fMRI in middle layers

# Joint processing of linguistic properties in brains and language models

- Stimuli: Narrative Stories
- Stimulus representation: pretrained NLP model and removal of linguistic properties
- Brain recording & modality: fMRI, Listening
- Questions: What linguistic properties underlie brain alignment, across all layers but also specifically in middle layers?

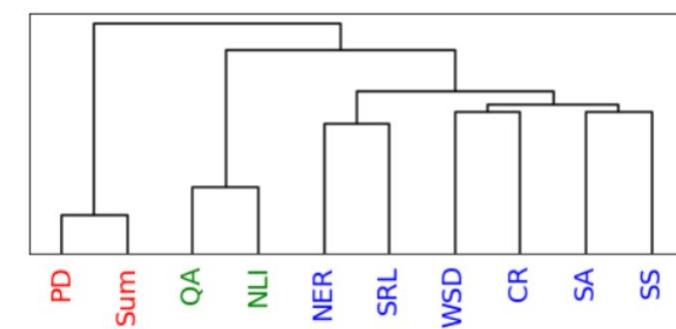
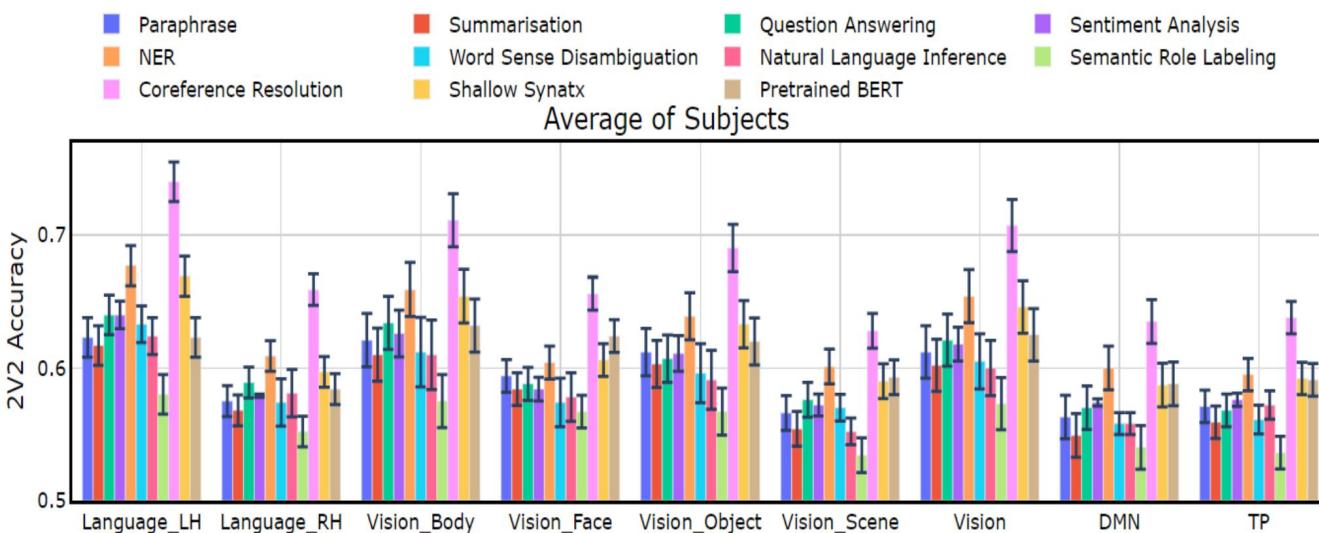


Top constituents and Tree Depth contribute the most to the alignment trend across layers

# Tasks affect processing

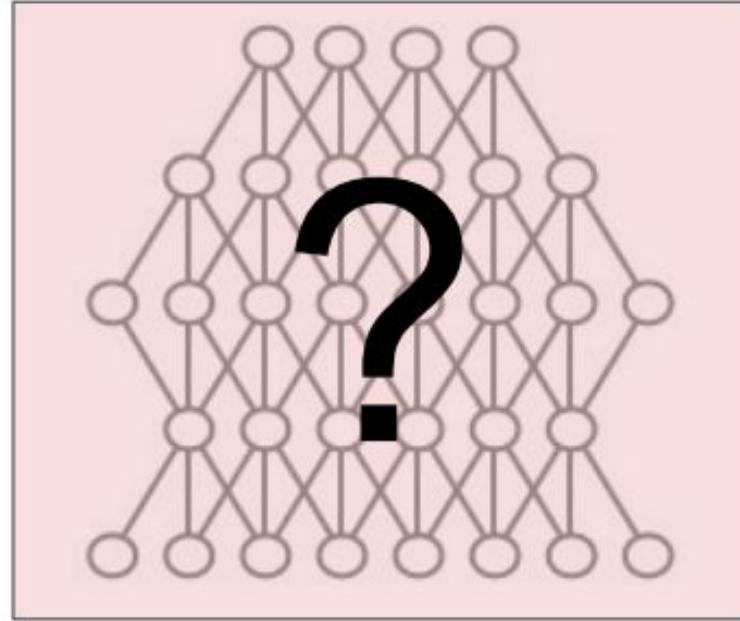
- Stimuli: passages and narratives
- Stimulus representation: task-optimized NLP models for a range of tasks
- Brain recording & modality: fMRI, reading & listening of different stimuli

- Reading fMRI best explained by coref. resolution, NER, shallow syntax parsing
- Listening fMRI best explained by paraphrasing, summarization, NLI



Oota, Subba Reddy, Jashn Arora, Veeral Agarwal, Mounika Marreddy, Manish Gupta, and Bapi Raju Surampudi. "Neural Language Taskonomy: Which NLP Tasks are the most Predictive of fMRI Brain Activity?." *arXiv preprint arXiv:2205.01404* (2022).

# Mechanistic Interpretability



Interpreting Model Predictions

- Why did the model make this prediction?

What if?  
Drop layers..

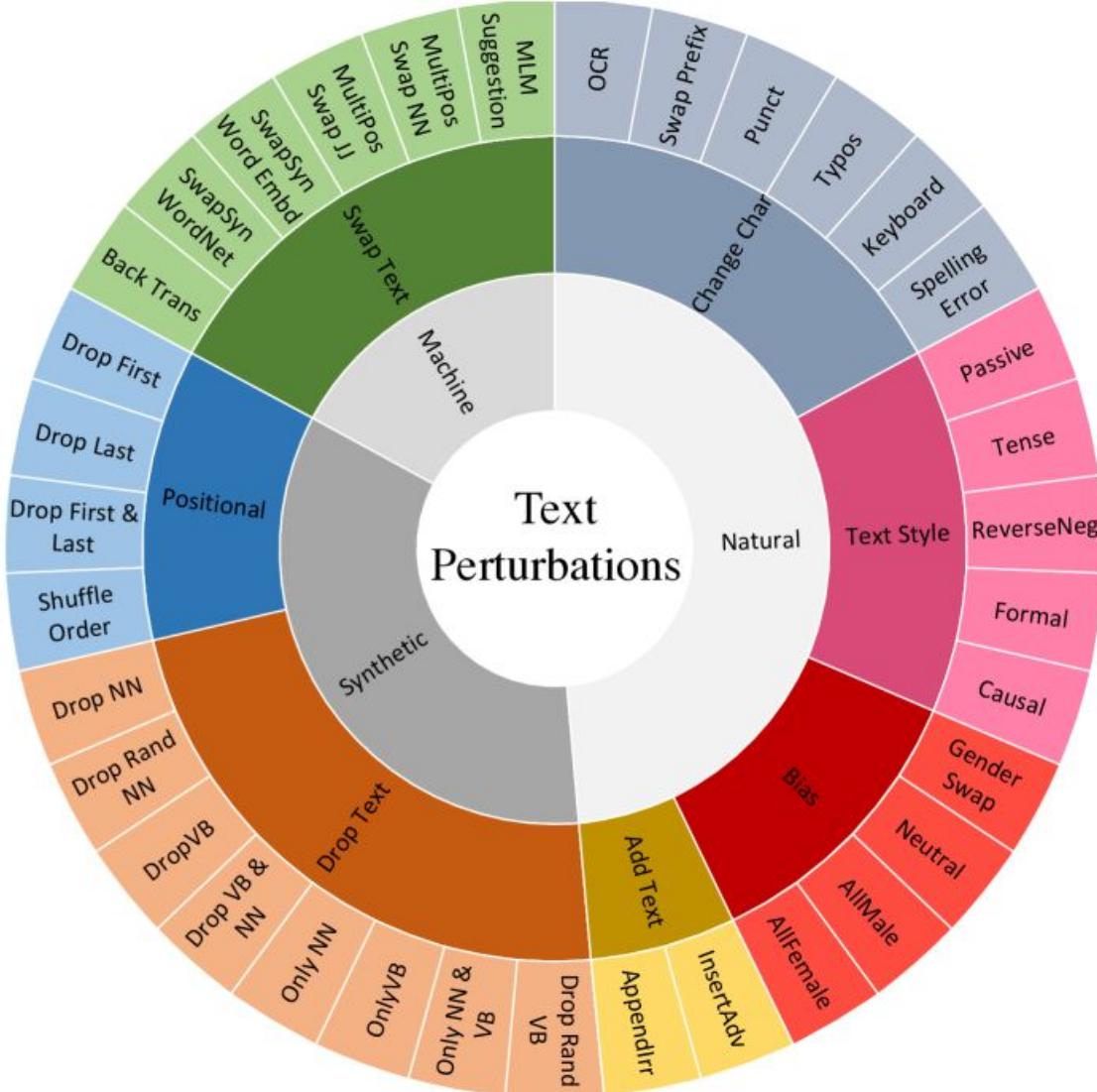
What if?  
Change input  
examples..

What if?  
Change of  
weights..

# Interesting questions about BERT, GPT2

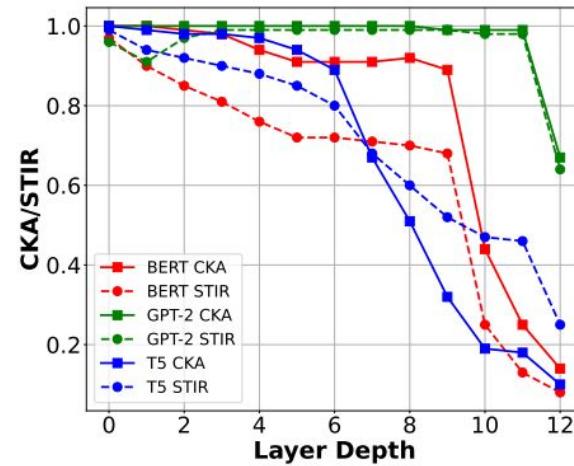
- Finetuning modifies the representations generated by each layer.
  - While fine-tuning these models, what changes across layers with respect to the pre-trained checkpoints?
- Robustness to input perturbations
  - How robust are BERT, GPT2 to input perturbations?
  - Is the effect of finetuning consistent 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  $m_1$  and  $n$  data samples, STIR defines how invariant finetuned model  $m_2$  is to perturbations of the samples that are imperceptible by  $m_1$ , i.e., do not change their representations according to  $m_1$ .

# Robustness

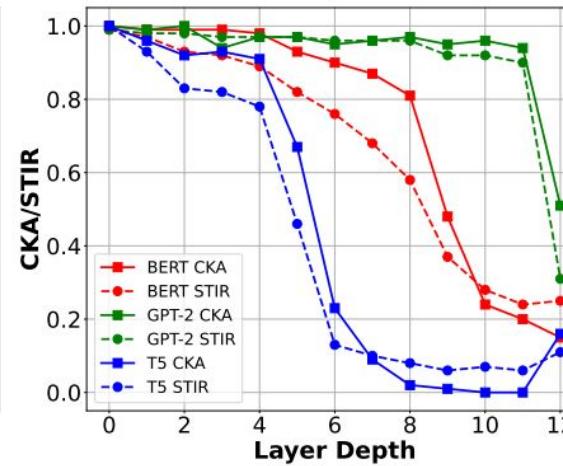


- Text perturbations grouped into seven different categories
  - ChangeChar
  - AddText
  - Bias
  - Positional
  - DropText
  - SwapText
  - TextStyle

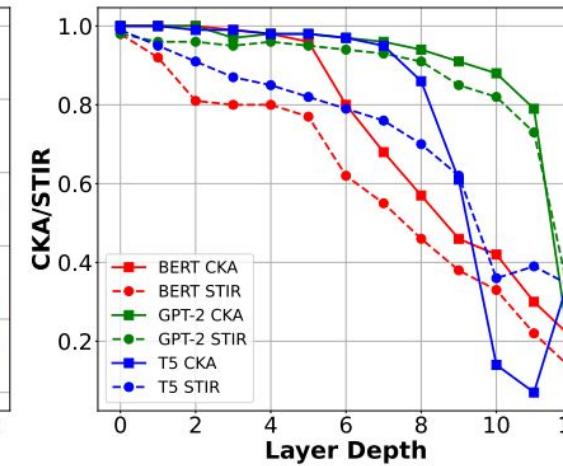
# How does finetuning modify the layers representations for different models?



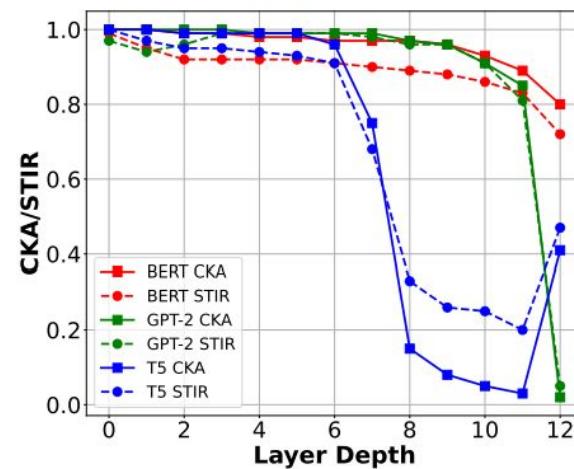
(a) CoLA



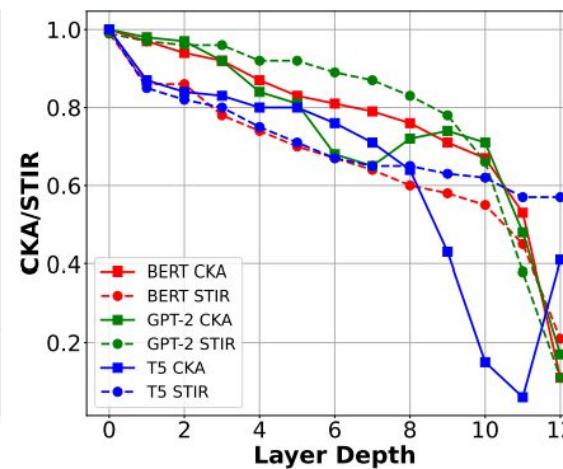
(b) SST-2



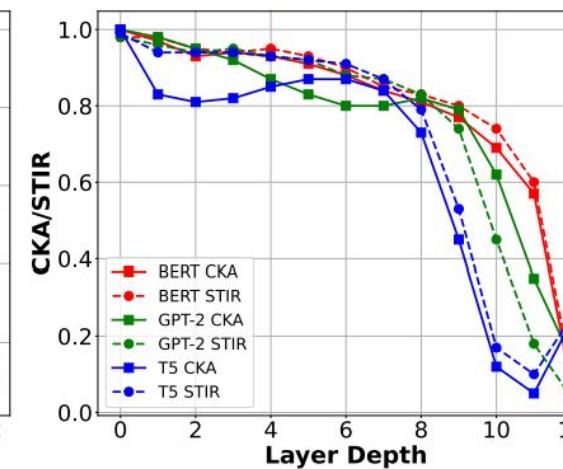
(c) MRPC



(d) STS-B



(e) QQP



(f) MNLI-Matched

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

# 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>	<b>0.24</b>	0.92	<b>0.93</b>	<b>0.93</b>	0.94	<b>0.96</b>	0.94	0.56	0.48	<b>0.57</b>	0.89	<b>0.92</b>	0.89
Drop verbs	<u>0.05</u>	<b>0.24</b>	0.06	<b>0.95</b>	<b>0.95</b>	<b>0.95</b>	0.98	<b>0.99</b>	0.96	<b>0.93</b>	0.92	0.89	<b>0.97</b>	<u>0.96</u>	0.96
Drop first	0.48	<b>0.75</b>	0.54	<b>0.98</b>	0.97	<b>0.98</b>	<b>1.00</b>	0.99	<b>1.00</b>	<b>0.98</b>	0.93	0.94	<b>0.99</b>	0.98	<b>0.99</b>
Drop last	0.34	<b>0.45</b>	0.32	<b>1.00</b>	0.99	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	<b>0.84</b>	0.83	<u>0.83</u>	<u>0.95</u>	<b>0.96</b>	<u>0.95</u>
Swap text	0.13	<b>0.16</b>	0.06	<b>0.98</b>	<b>0.98</b>	0.97	0.99	<b>1.01</b>	0.98	<b>0.98</b>	<u>0.96</u>	0.95	<b>0.97</b>	<b>0.97</b>	0.96
Add text	0.85	<b>0.92</b>	0.86	<b>0.99</b>	<b>0.99</b>	<b>0.99</b>	0.93	<b>1.00</b>	<u>0.96</u>	<b>0.99</b>	<b>0.99</b>	0.98	0.99	<b>1.00</b>	0.99
Change char	0.14	<b>0.29</b>	<b>0.29</b>	0.84	<b>0.86</b>	0.84	0.43	<b>0.97</b>	0.65	<b>0.58</b>	0.52	0.57	<u>0.88</u>	<b>0.95</b>	0.94
Bias	0.95	<b>0.96</b>	0.92	1.00	<b>1.01</b>	1.00	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	<b>0.99</b>	<b>0.99</b>	<b>0.99</b>	1.00	<b>1.01</b>	1.00

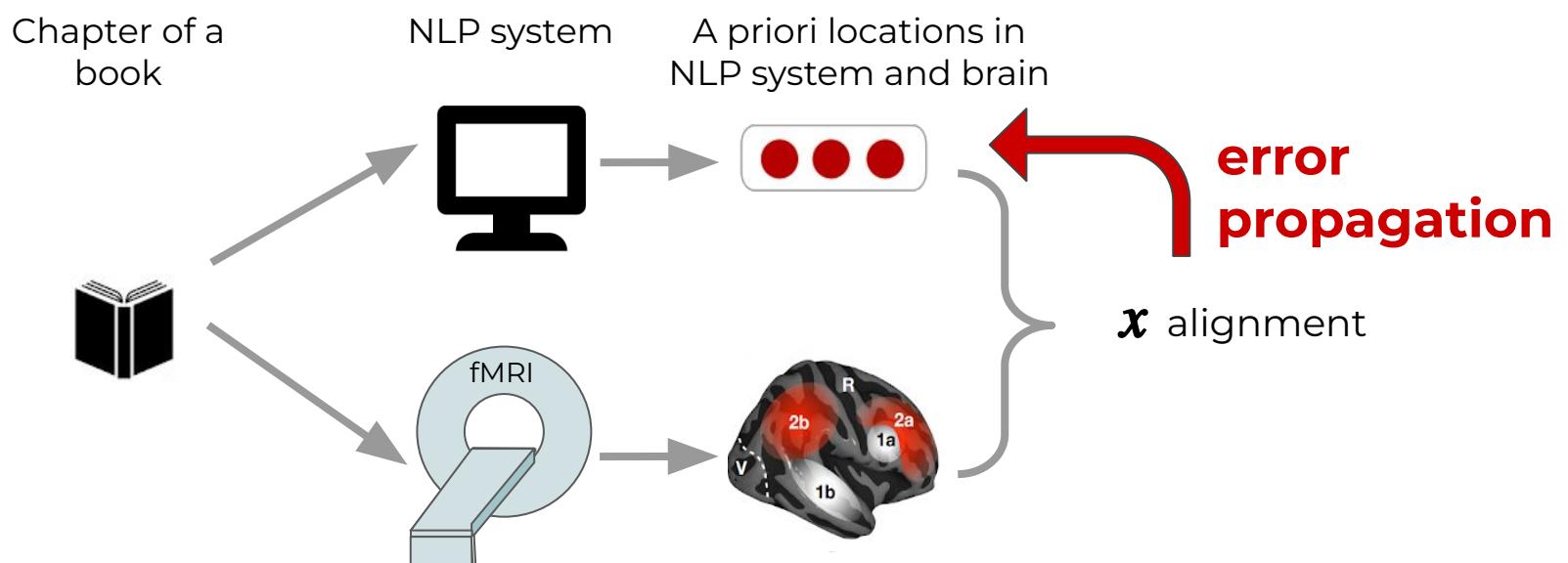
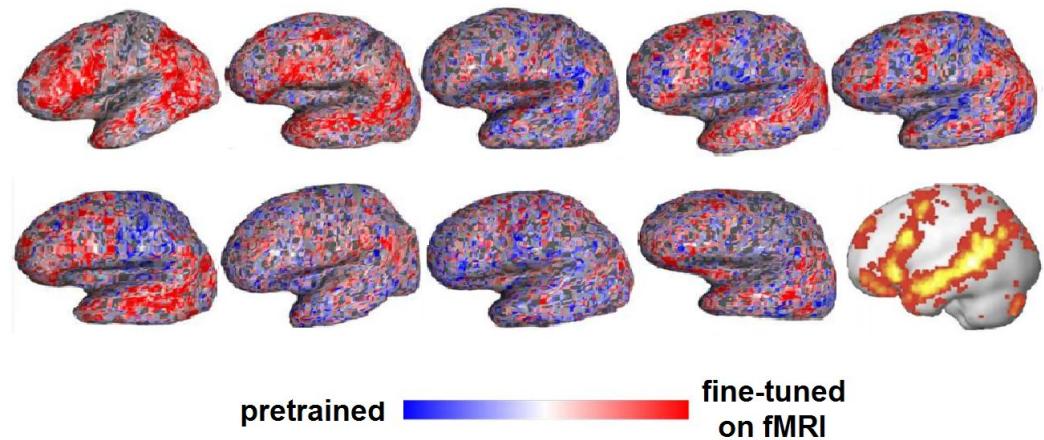
  

Perturbation	MNLI-m (Accuracy)			QNLI (Accuracy)			RTE (Accuracy)			WNLI (Accuracy)		
	BERT	GPT-2	T5	BERT	GPT-2	T5	BERT	GPT-2	T5	BERT	GPT-2	T5
Drop nouns	0.83	<b>0.85</b>	0.83	0.82	<b>0.87</b>	0.82	0.84	<b>1.01</b>	0.89	1.00	1.00	<b>1.05</b>
Drop verbs	0.89	<b>0.90</b>	<b>0.90</b>	<b>0.96</b>	0.94	0.94	0.98	<b>1.01</b>	<u>0.96</u>	1.00	1.01	<b>1.03</b>
Drop first	0.94	0.94	<b>0.95</b>	0.97	<b>0.98</b>	0.97	0.95	<b>1.00</b>	<b>1.00</b>	1.00	<u>0.99</u>	<b>1.01</b>
Drop last	0.89	<b>0.90</b>	0.89	0.97	<b>0.98</b>	0.97	0.97	<b>1.01</b>	0.98	<b>1.00</b>	<u>0.99</u>	<b>1.00</b>
Swap text	0.94	<b>0.95</b>	0.94	<b>0.97</b>	<b>0.97</b>	<b>0.97</b>	<b>0.98</b>	<u>0.97</u>	0.97	1.00	<b>1.01</b>	1.00
Add text	<b>0.95</b>	<b>0.95</b>	<b>0.95</b>	0.99	<b>1.00</b>	0.99	<b>1.00</b>	<u>0.99</u>	0.97	1.00	<b>1.03</b>	<b>1.03</b>
Change char	0.67	0.66	<b>0.68</b>	<b>0.77</b>	0.75	<b>0.77</b>	0.82	<b>0.99</b>	<u>0.83</u>	1.00	<b>1.03</b>	1.01
Bias	0.99	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	1.00	<b>1.02</b>	1.00	1.00	<b>1.02</b>	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.

# Training DL models using brain recordings

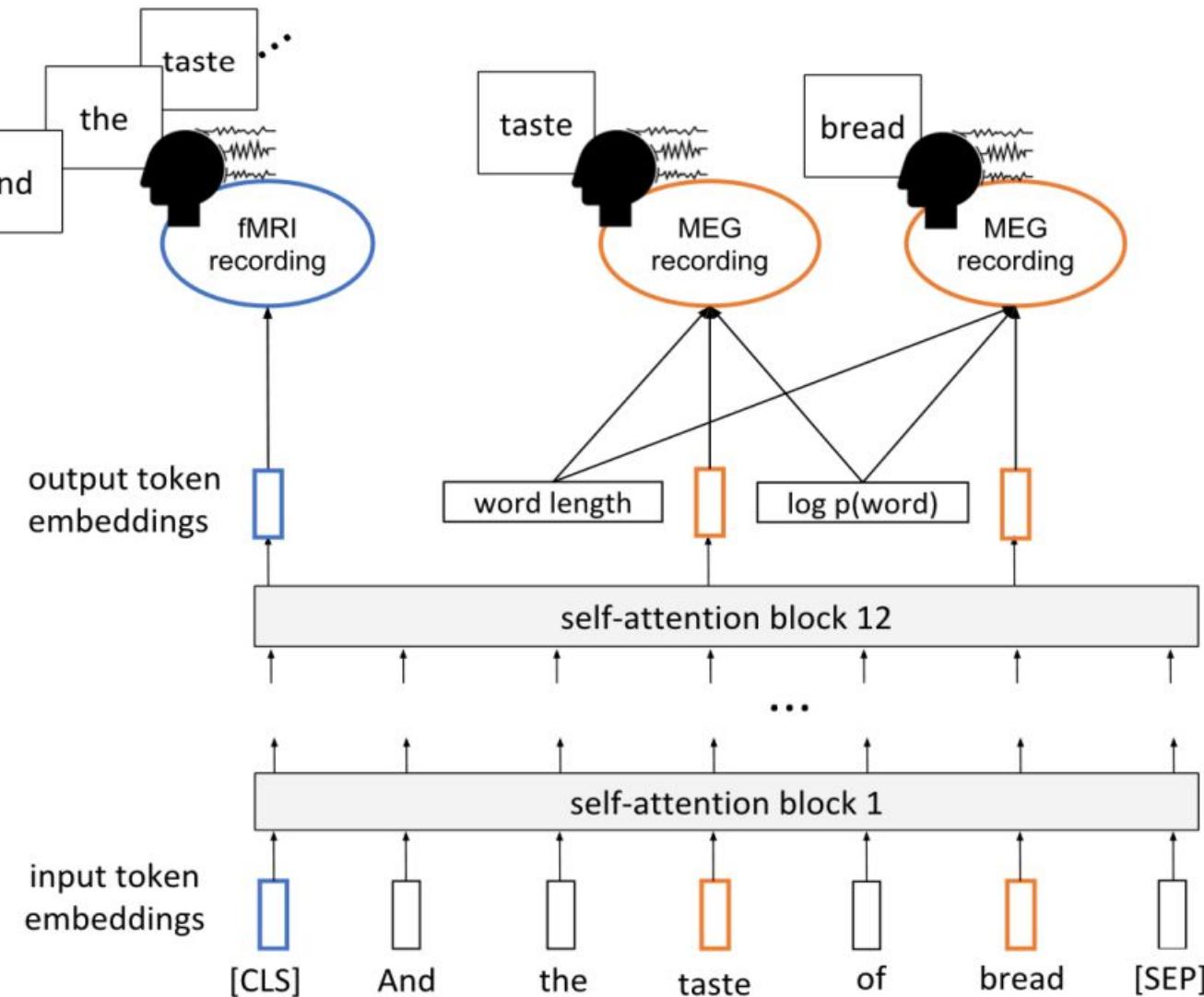
- Stimuli: one chapter of Harry Potter
- Stimulus representation: brain-optimized NLP model
- Brain recording & modality: fMRI & MEG, reading



Brain-optimized NLP  
model predicts unseen  
fMRI recordings better,  
especially in canonical  
language regions

Schwartz, Dan, Mariya Teneva, and Leila Wehbe. "Inducing brain-relevant bias in natural language processing models." *Advances in neural information processing systems* 32 (2019).

# Inducing Brain Relevant Bias

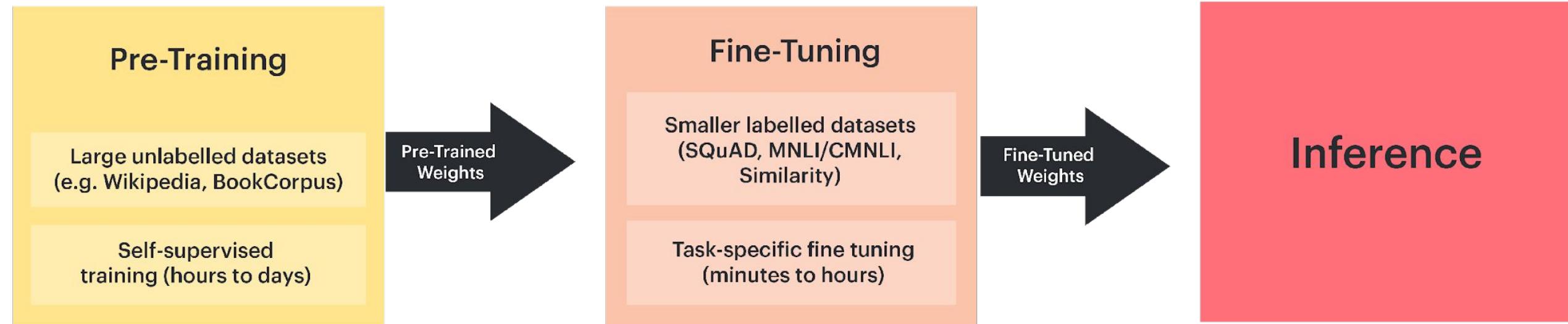


Metric	Vanilla	MEG	Joint
CoLA	57.29	57.63	<b>57.97</b>
SST-2	93.00	<b>93.23</b>	91.62
MRPC (Acc.)	83.82	83.97	<b>84.04</b>
MRPC (F1)	88.85	<b>88.93</b>	88.91
STS-B (Pears.)	<b>89.70</b>	89.32	88.60
STS-B (Spear.)	<b>89.37</b>	88.87	88.23
QQP (Acc.)	90.72	<b>91.06</b>	90.87
QQP (F1)	87.41	<b>87.91</b>	87.69
MNLI-m	83.95	<b>84.26</b>	84.08
MNLI-mm	84.39	84.65	<b>85.15</b>
QNLI	89.04	<b>91.73</b>	91.49
RTE	61.01	<b>65.42</b>	62.02
WNLI	53.52	<b>53.80</b>	51.97

# 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, Prompting, Instruction-tuning [1 hour]
- Lunch [45 min]
- Large language models, Chain-of-Thought reasoning, RAG [1 hour 30 min]
- Coffee break [15 min]
- Hands on session [1 hour]

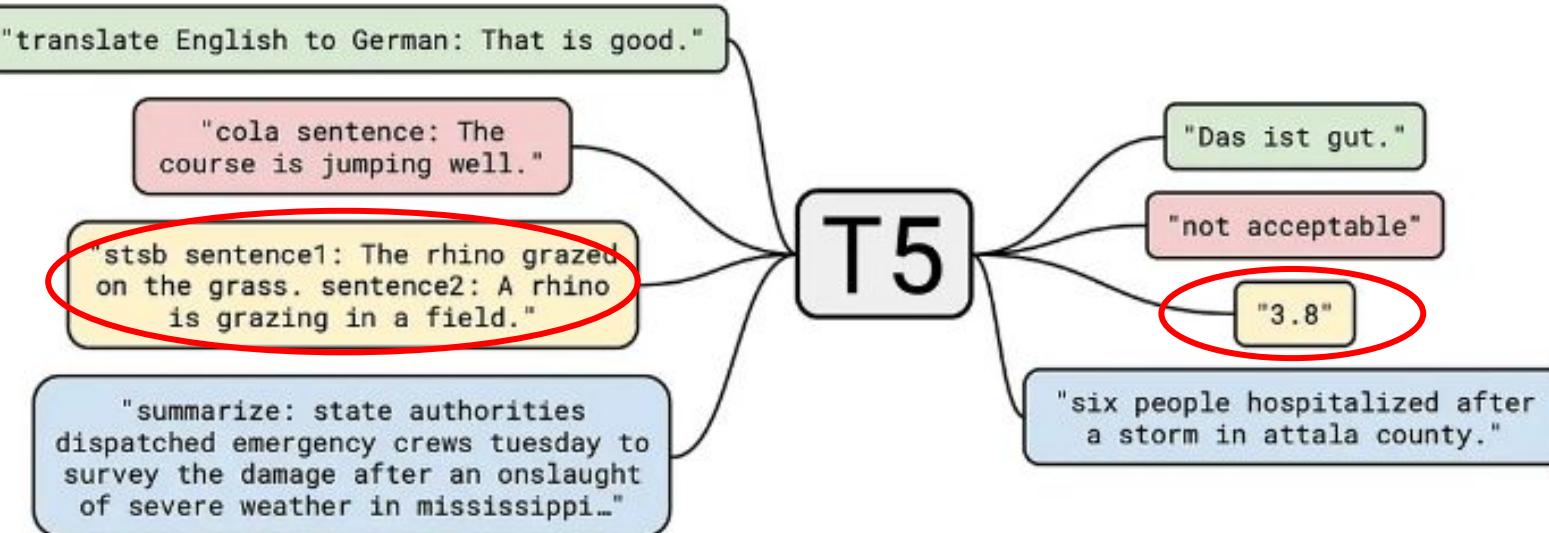
# Downside of full fine tuning



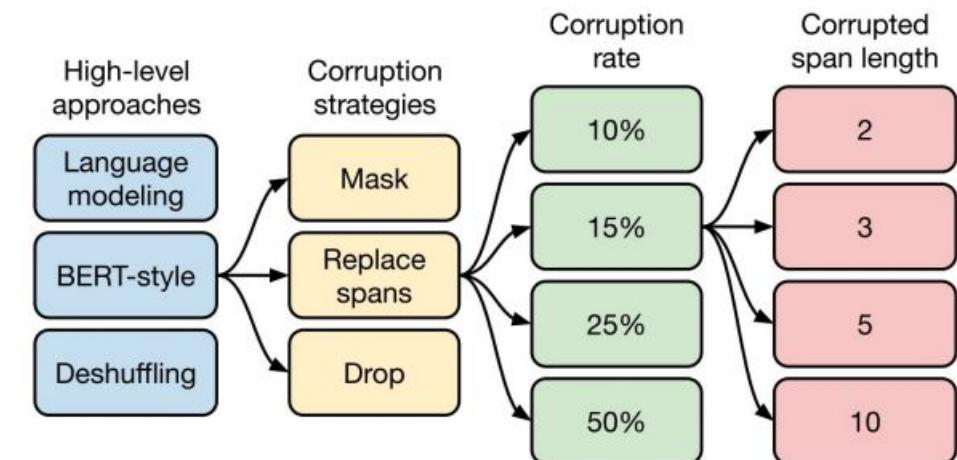
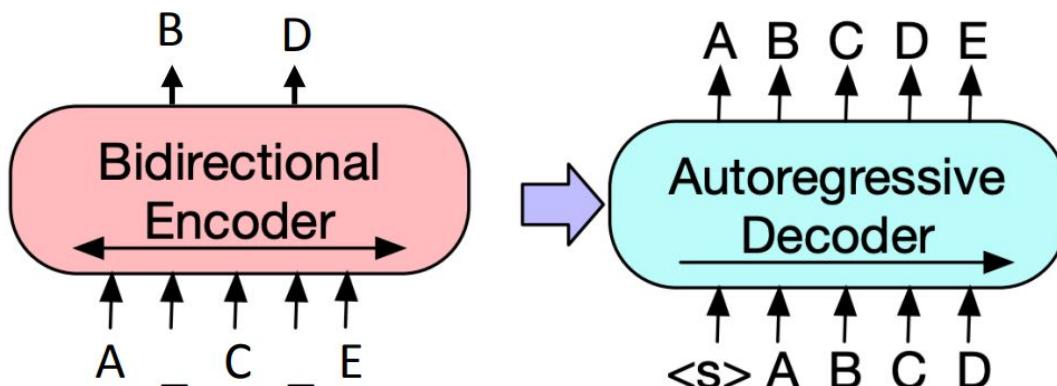
# Agenda

- Recap on small language models [5 min]
- Emerging abilities of language models and why they are effective? [25 min]
  - Analyzing and Interpreting language models
- **Text-to-Text Transfer Transformer, Prompting, Instruction-tuning [1 hour]**
- Lunch [45 min]
- Large language models, Chain-of-Thought reasoning, RAG [1 hour 30 min]
- Coffee break [15 min]
- Hands on session [1 hour]

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



- Text-to-text transformer
- Encoder-decoder model
- Reformulates all tasks (during pretraining and finetuning) with a text-to-text format



# 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 .

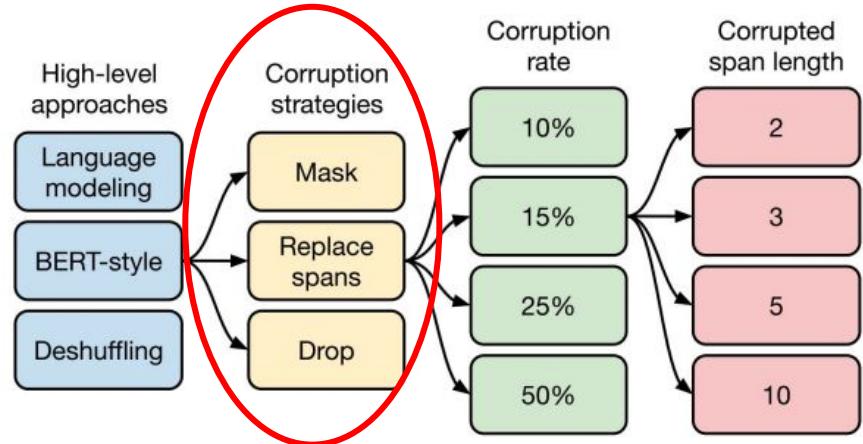
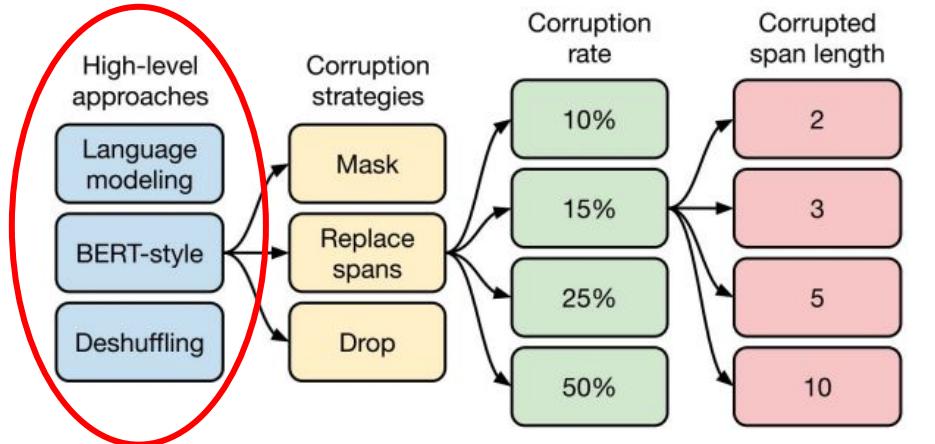
C . D E . A B

Sentence Permutation Document Rotation

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

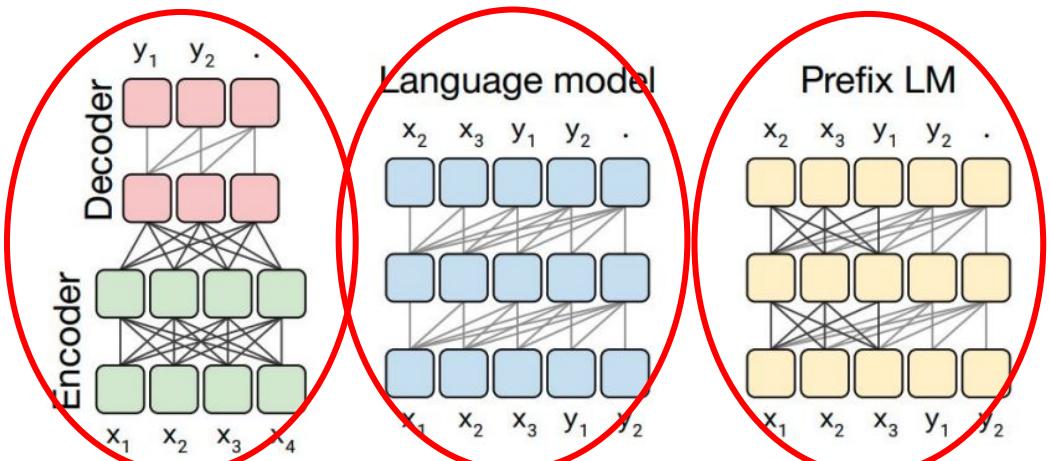
# T5: Different unsupervised objectives

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



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

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 .	(original text) <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
- **Text-to-Text Transfer Transformer, Prompting, Instruction-tuning [1 hour]**
  - Behavioural interpretability
- Lunch [45 min]
- Large language models, Chain-of-Thought reasoning, RAG [1 hour 30 min]
- Coffee break [15 min]
- Hands on session [1 hour]

# Aspects of NLP models and Brain datasets

- Stimuli: Narrative Stories
- Stimulus representation: pretrained NLP model and speech models
- Brain recording & modality: fMRI, Reading, Listening
  
- **Questions:** Is the choice of stimulus modality (reading vs. listening) important for the study of brain alignment?
- Are all naturalistic fMRI datasets equally good for brain encoding?
- How does the type of model (text vs. speech and encoder vs. decoder) affect the resulting alignment?

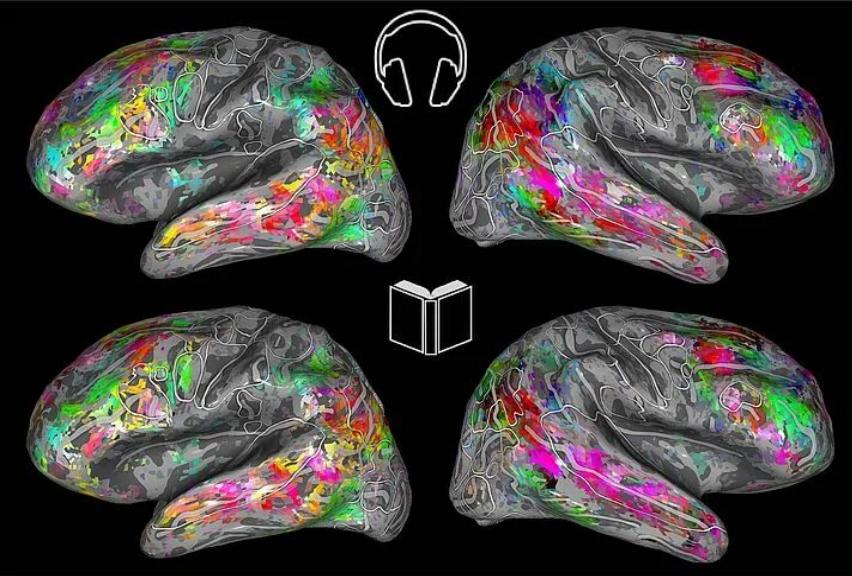


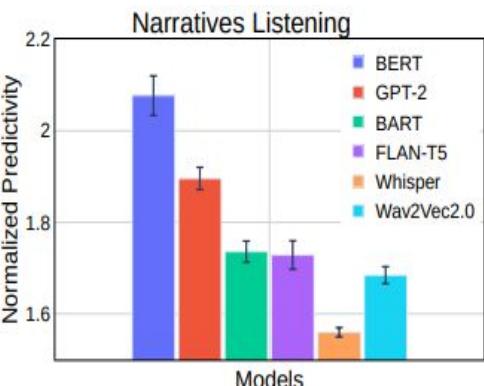
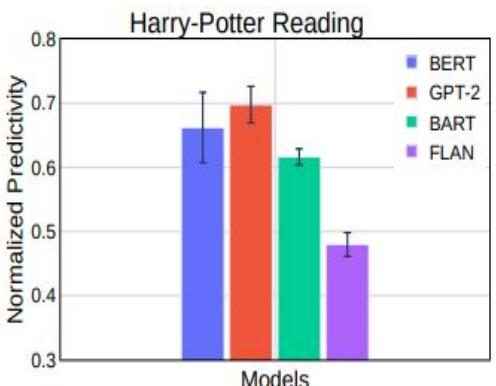
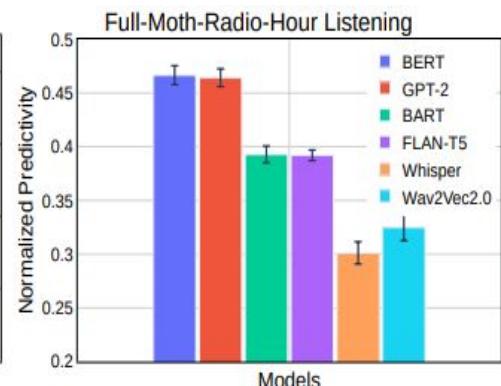
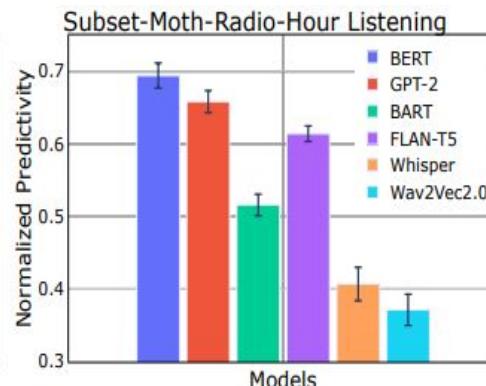
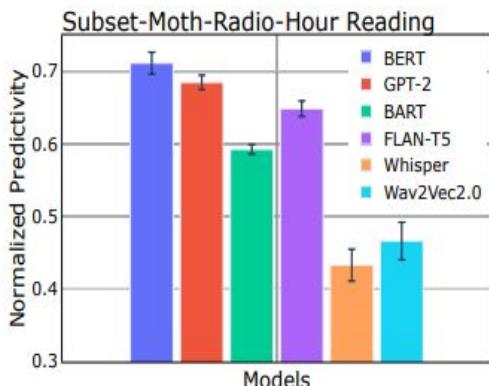
Table 1: Naturalistic Stories Datasets

Dataset	Modality	Subj	1-TR	# TRs
Full-Moth-Radio-Hour	Listening	8	2.0045s	9932
Subset-Moth-Radio-Hour	Reading	6	2.0045s	4028
Subset-Moth-Radio-Hour	Listening	6	2.0045s	4028
Narratives (21 <sup>st</sup> -Year)	Listening	18	1.5s	2250
Harry-Potter	Reading	8	2s	1211

Table 2: Neural Pretrained Transformer Models

Model Name	Pretraining	Type	Layers
BERT-base-uncased	Text	Encoder (Bidirectional)	12
GPT-2-Small	Text	Decoder (Unidirectional)	12
BART-base	Text	Encoder-Decoder	12
FLAN-T5-base	Text	Encoder-Decoder	24
Wav2Vec2.0-base	Speech	Encoder	12
Whisper-small	Speech	Encoder-Decoder	24

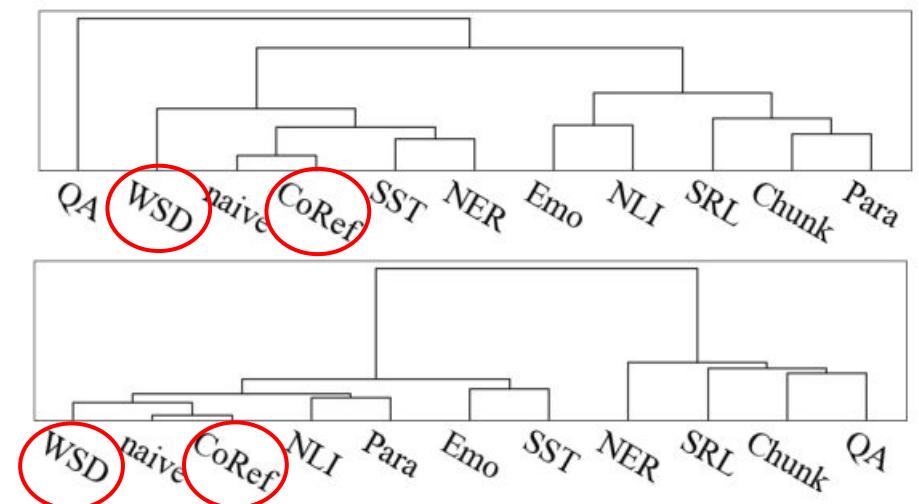
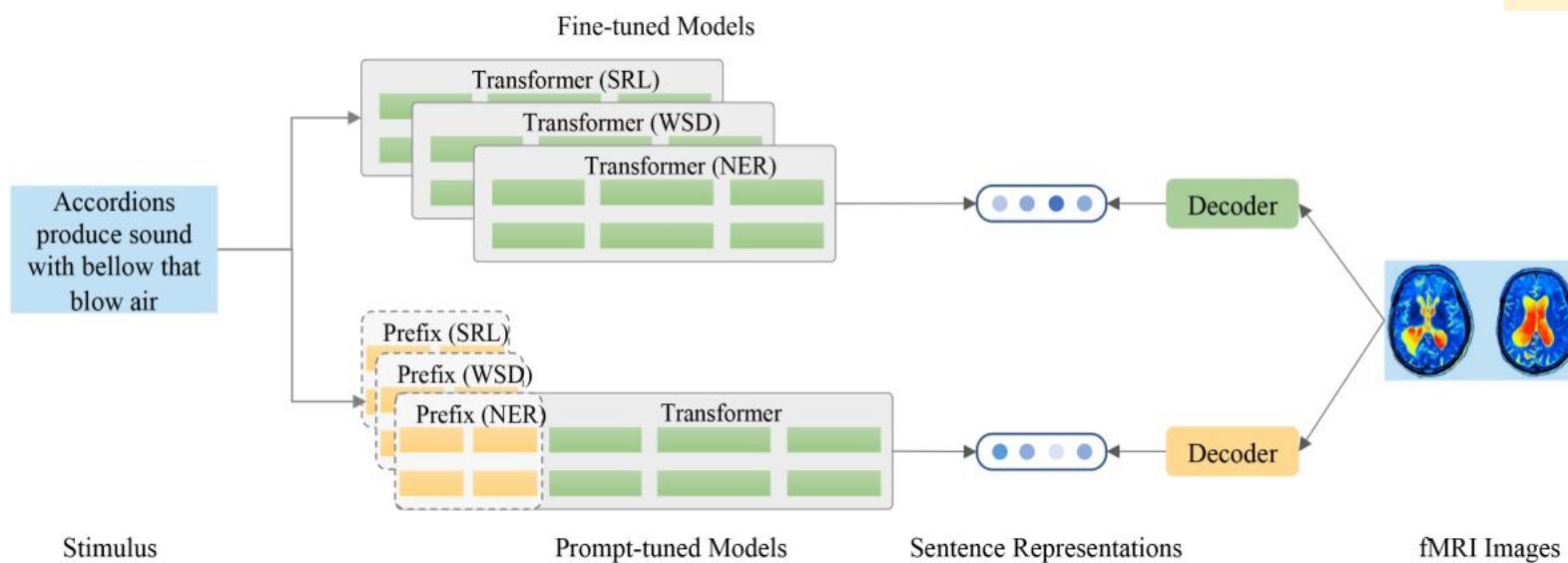
- Text models predict fMRI recordings significantly better than speech models



# Tasks affect processing: Prompting vs. finetuning

- Stimuli: passages
- Stimulus representation: Prompt-tuned NLP models for a range of tasks
- Brain recording & modality: fMRI, reading of different stimuli

- Full fine-tuning is cognitively inconsistent with the human brain's mechanism of language representation.
- Reading fMRI best explained by CoRef and WSD tasks irrespective of fine-tuning or prompt-tuning.
- Deep level semantic information possibly taking a larger proportion of brains language representation than the shallow syntactic information.

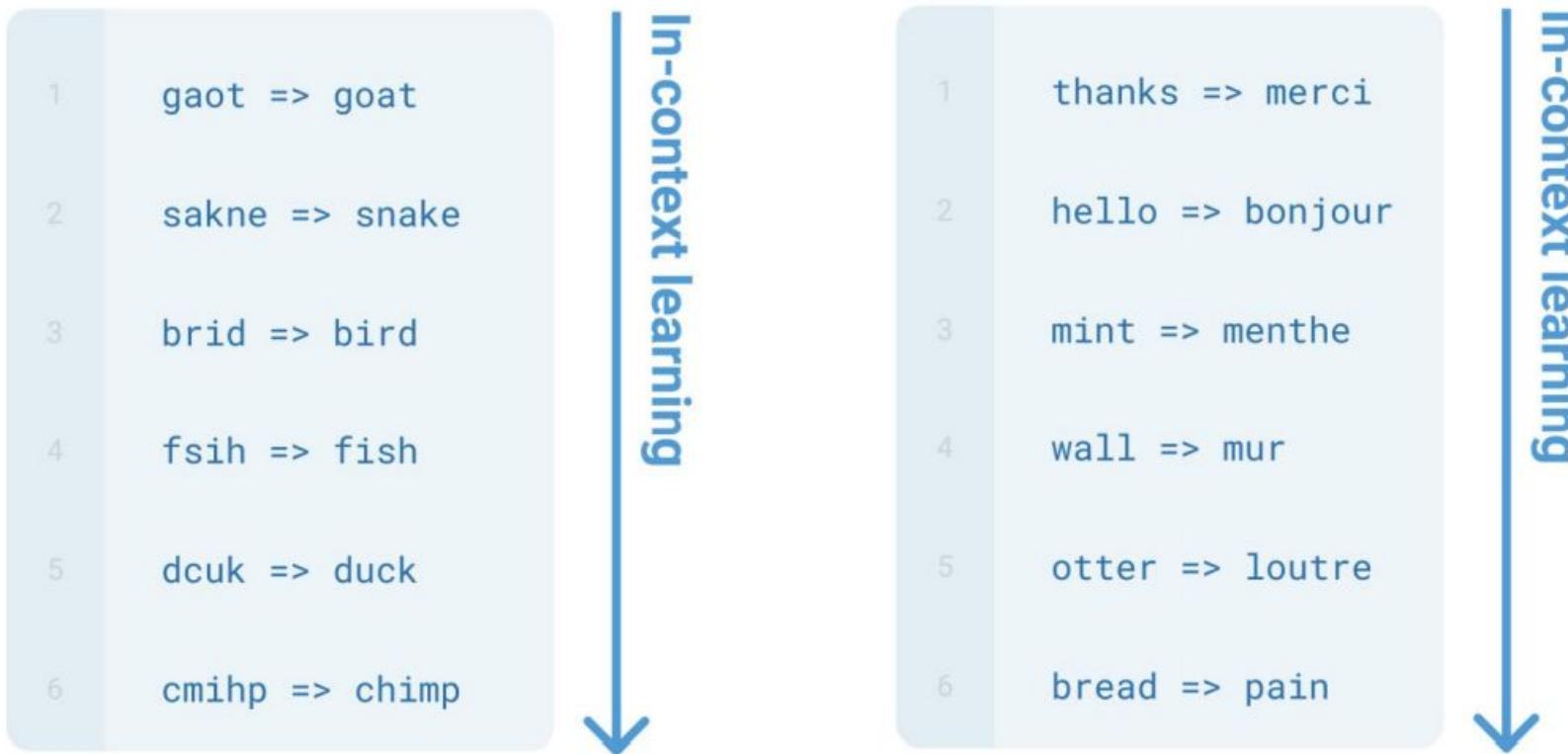


# Agenda

- 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, **Prompting**, Instruction-tuning [1 hour]
- Lunch [45 min]
- Large language models, Chain-of-Thought reasoning, RAG [1 hour 30 min]
- Coffee break [15 min]
- Hands on session [1 hour]

# 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 => .....

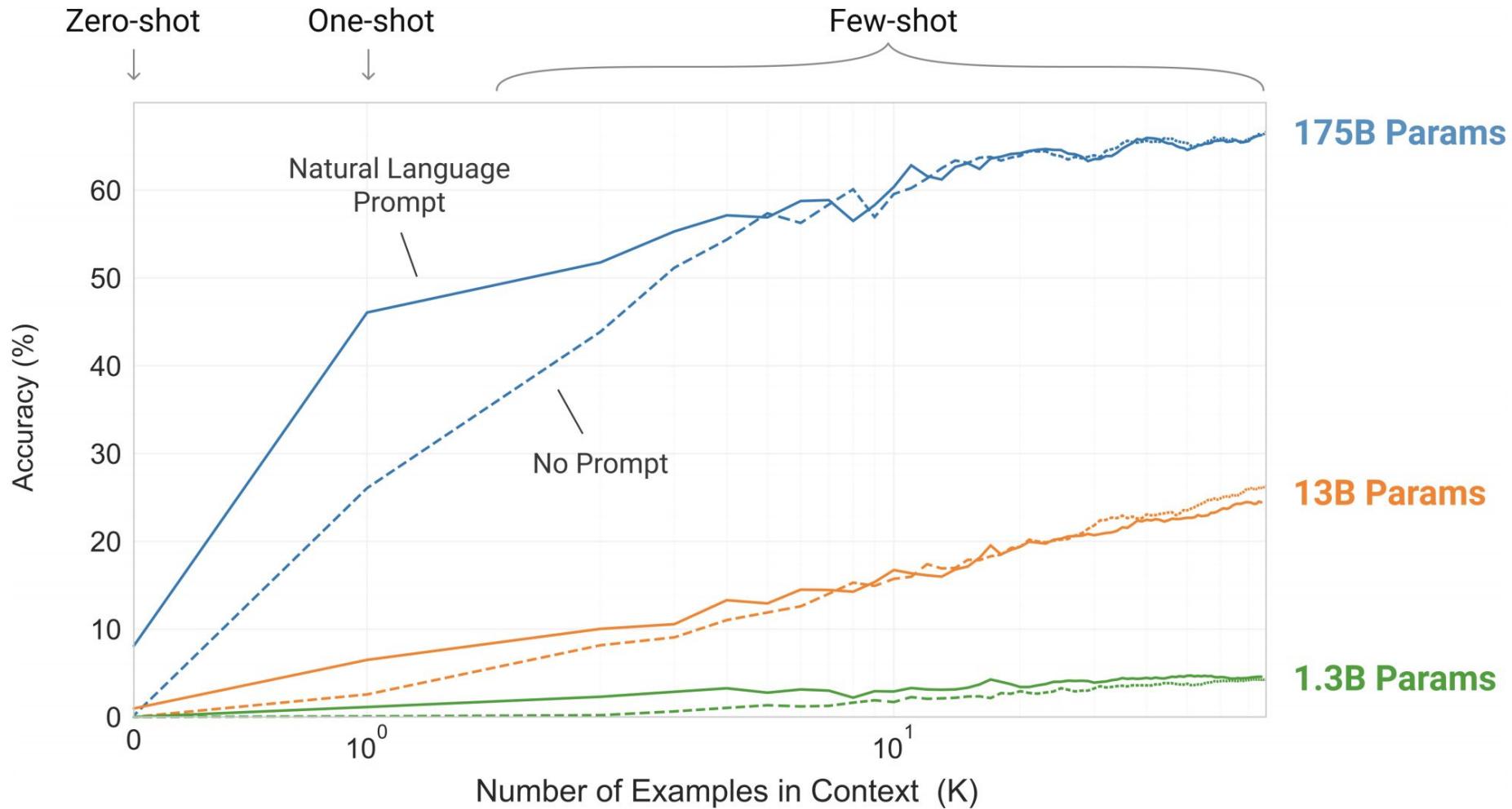
## One-shot

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

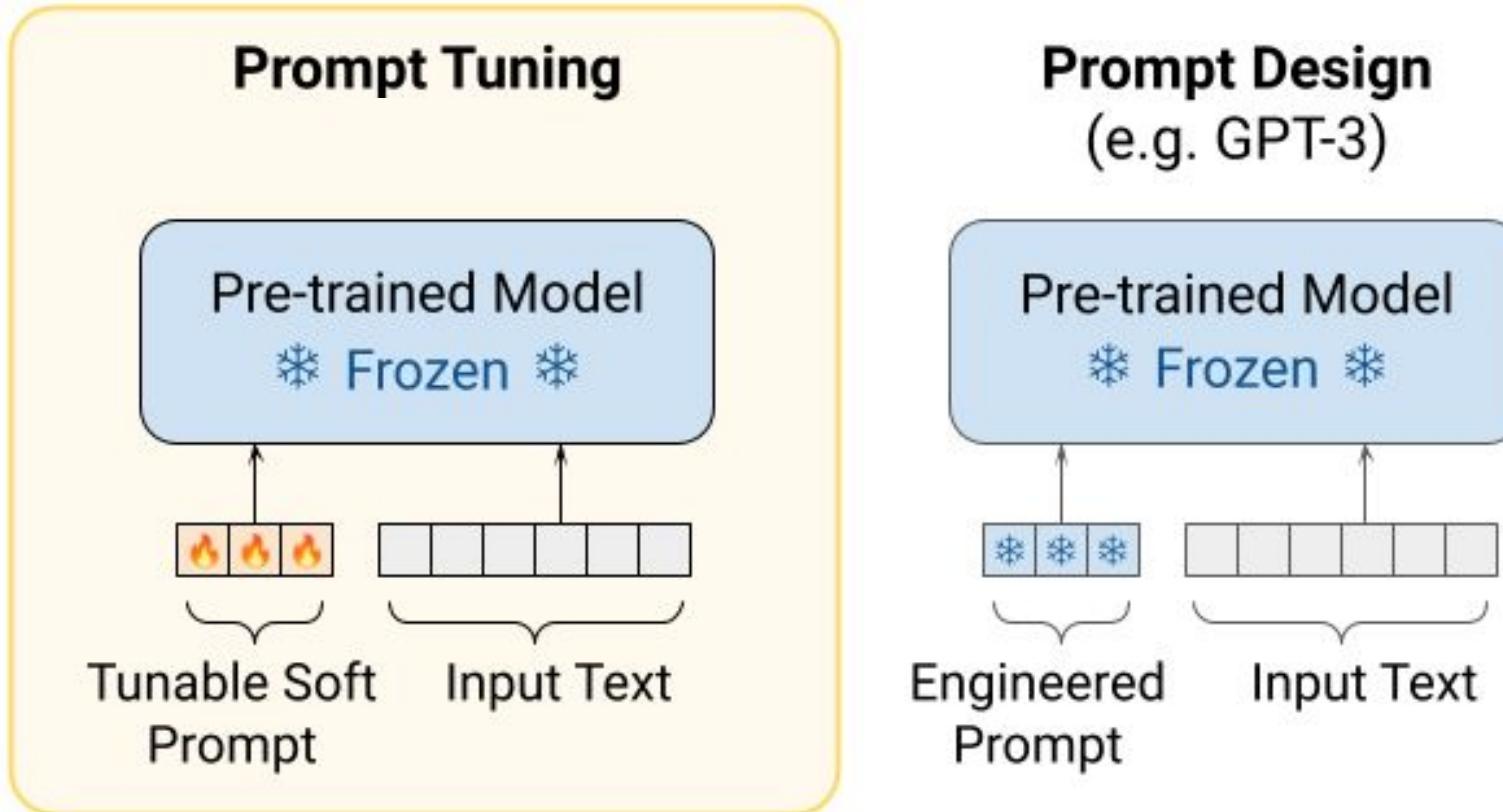
## Few-shot

- 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



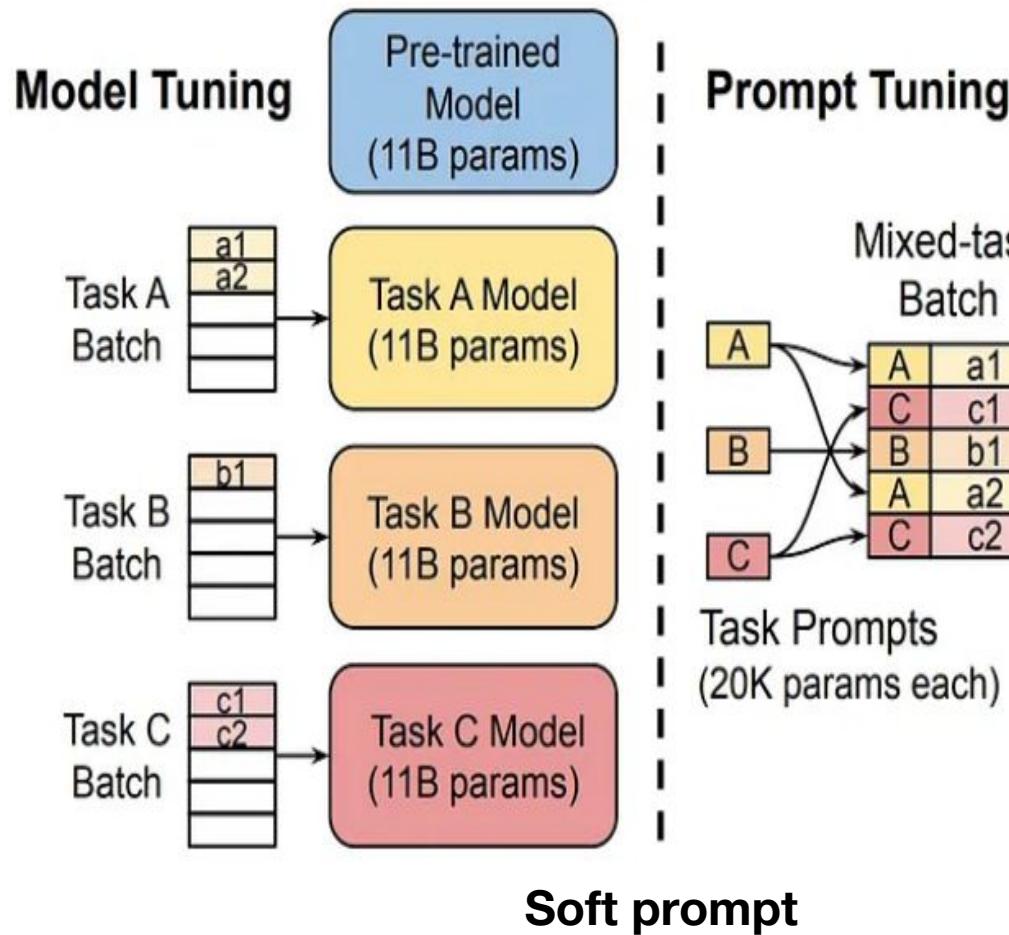
# Prompt-tuning



# Hard vs. Soft Prompts

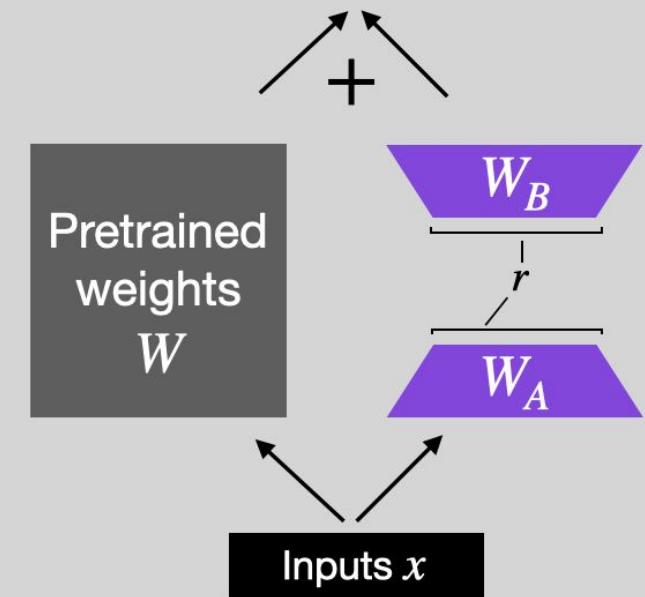
- Hard prompt: manually handcrafted text prompts with discrete input tokens
  - we directly change the discrete input tokens, which are not differentiable
  - Translate the English sentence {english\_sentence} into German: {german\_translation}
  - English: {english\_sentence} | German: {german\_translation}
  - 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

# Soft Prompt-tuning vs. Adapters



Forward pass with **updated** model

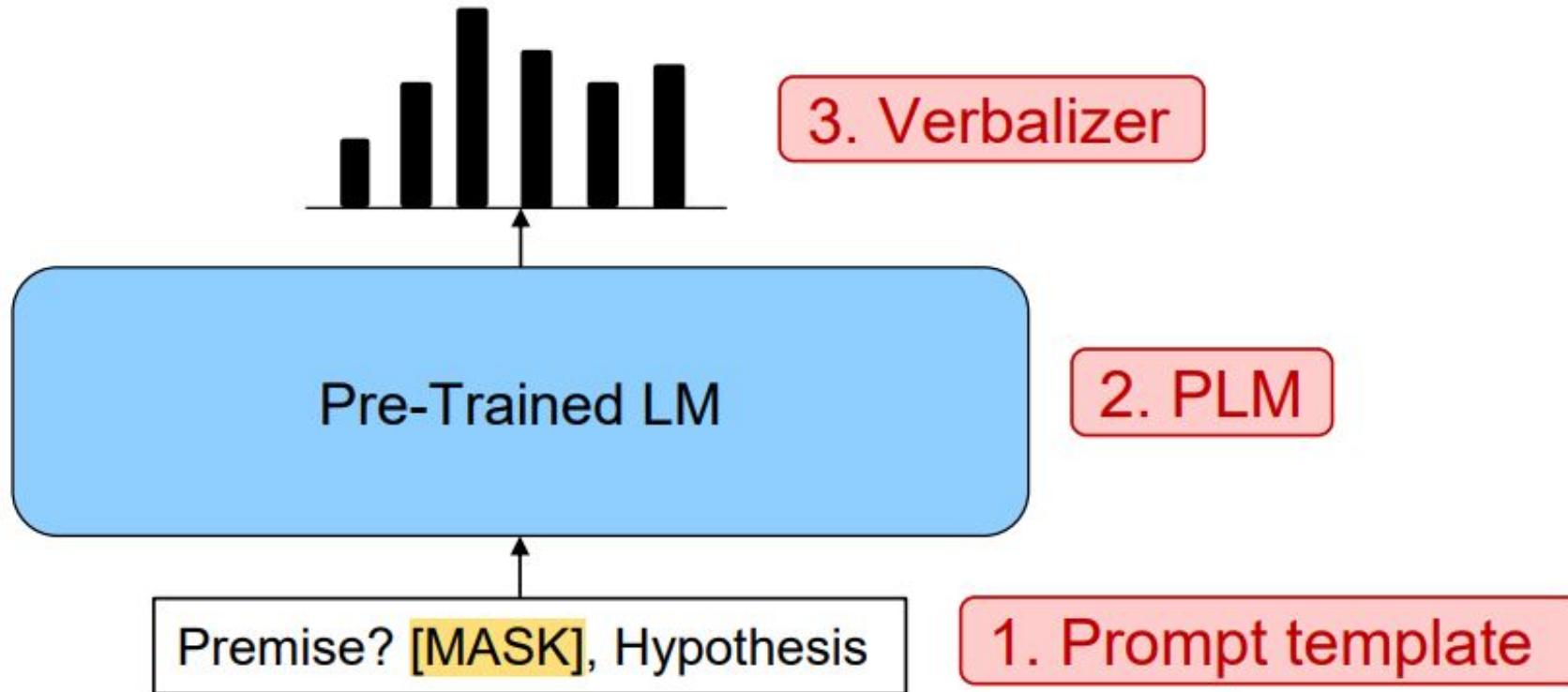
Embedding  $h$



**Low Rank Adaptation (LoRA)**

# Prompt-tuning

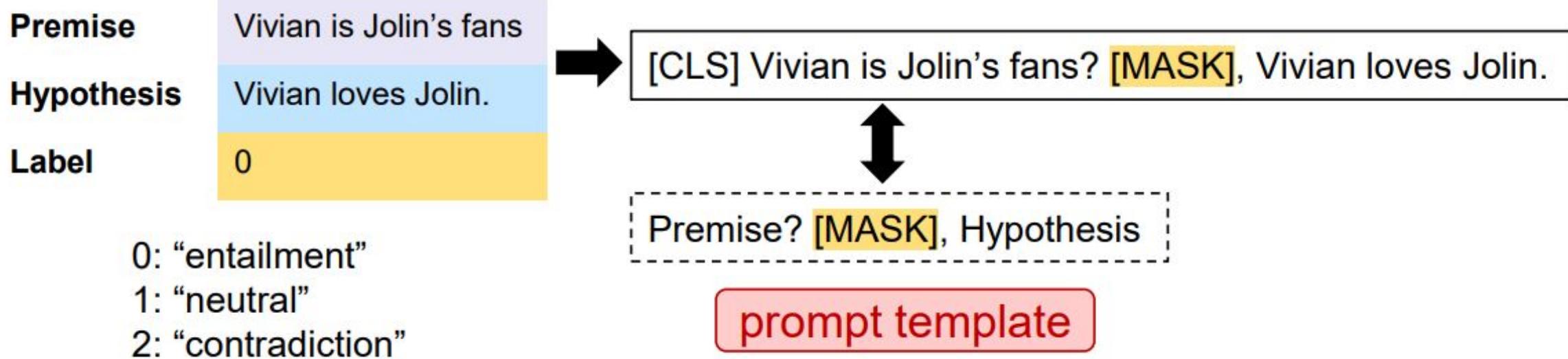
- Prompt-tuning refers to techniques that vary the input prompt to achieve better modeling results.
- Idea: convert data into natural language prompts
- better for few-shot, one-shot, or zero-shot cases



# Prompt-tuning

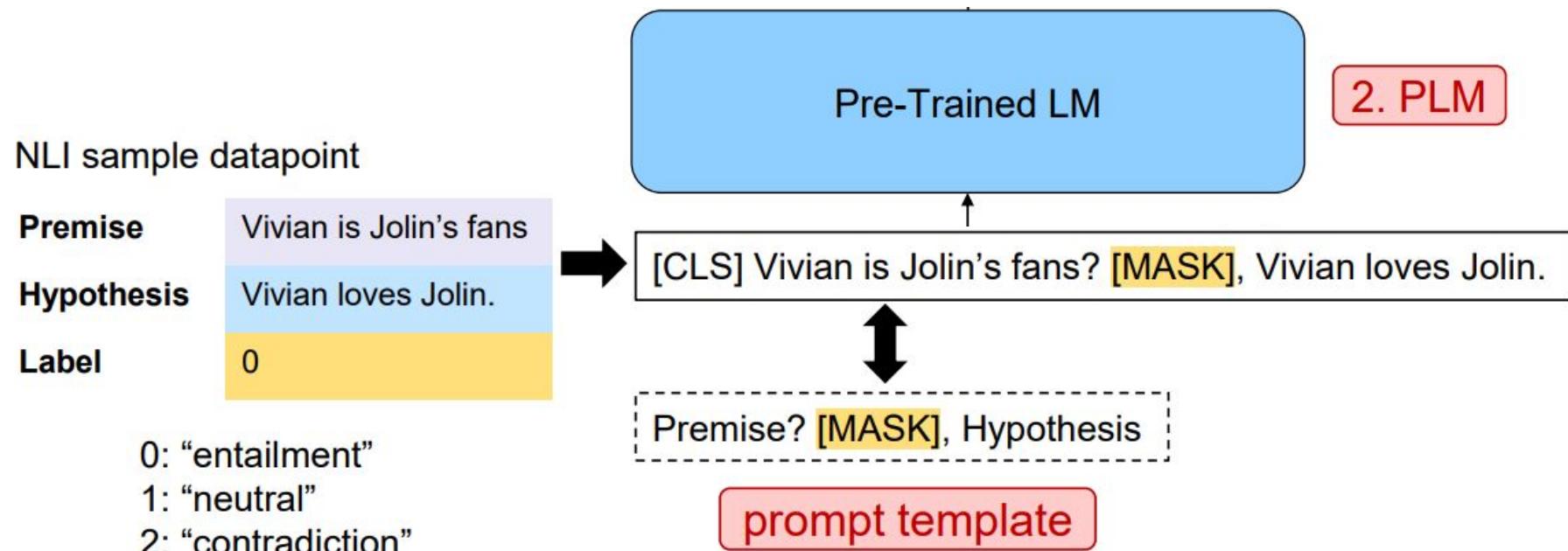
- Prompt template: manually designed natural language input for a task

NLI sample datapoint



# Prompt-tuning

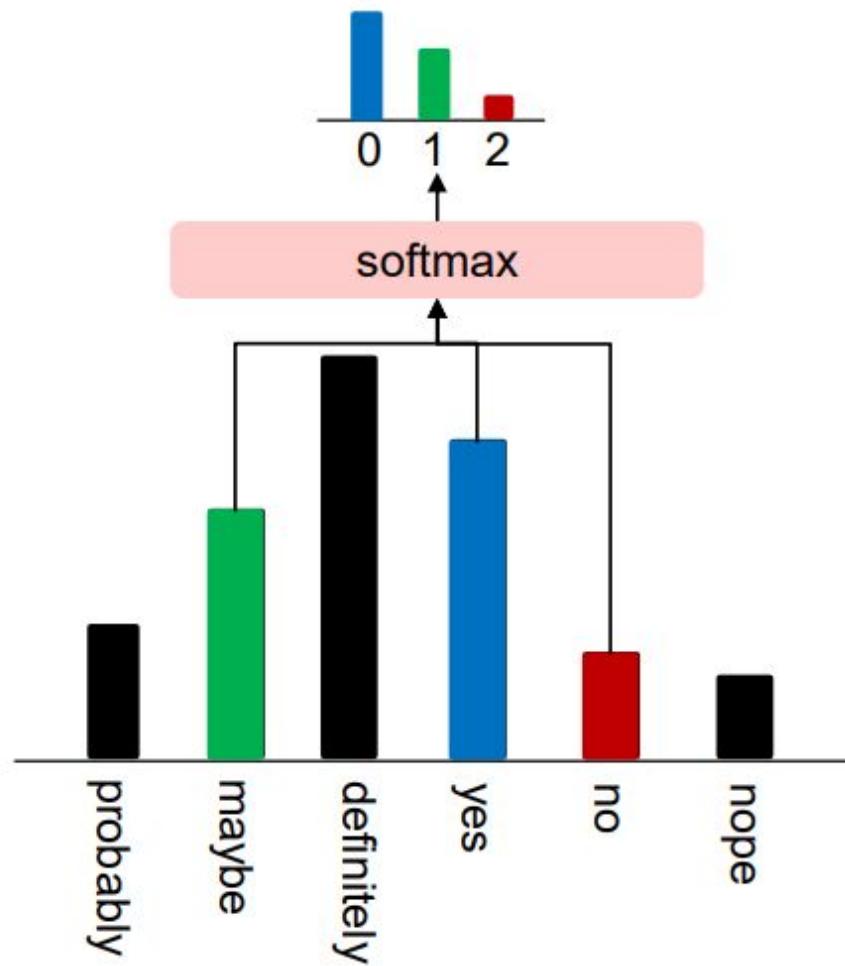
- Prompt template: manually designed natural language input for a task
- PLM: perform language modeling (masked LM or auto-regressive LM)



# Prompt-tuning

- Verbalizer: mapping from the vocabulary to labels

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



# Agenda

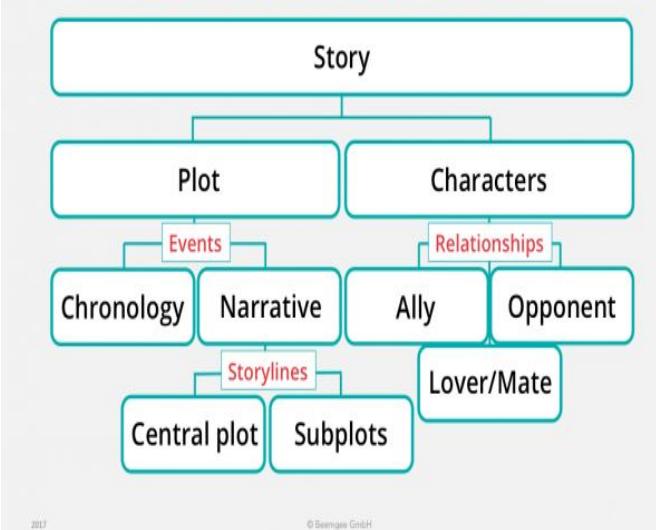
- Recap on small language models [5 min]
- Emerging abilities of language models and why they are effective? [25 min]
  - Analyzing and Interpreting language models
- Text-to-Text Transfer Transformer, **Prompting**, Instruction-tuning **[1 hour]**
  - **Behavioural interpretability**
- Lunch [45 min]
- Large language models, Chain-of-Thought reasoning, RAG [1 hour 30 min]
- Coffee break [15 min]
- Hands on session [1 hour]

# Large language models can segment narrative events similarly to humans

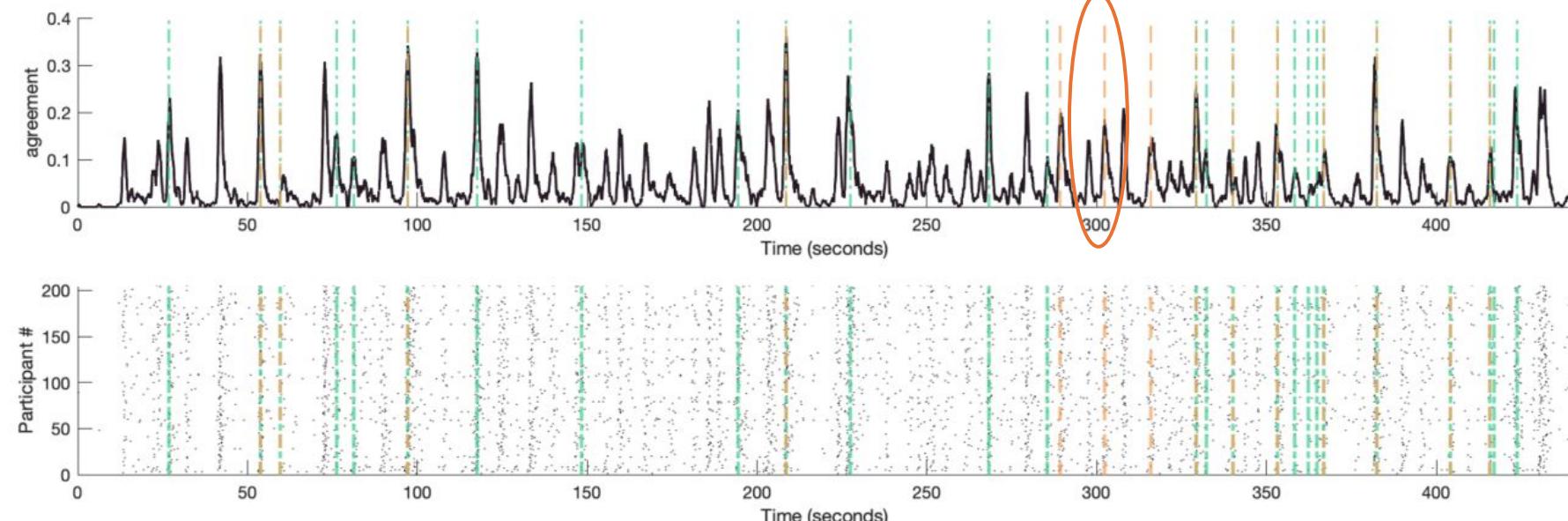
- Humans perceive discrete events such as "**restaurant visits**" and "**train rides**" in their continuous experience.
- Event cognition** with naturalistic stimuli like movies and stories typically involves **laborious hand annotation of event boundaries** that are often crowd-sourced from large behavioral samples in online experiments
- To what extent the inner workings of large language models reflect human cognition. ?  
(prompt GPT-3 to segment a continuous story into events)

- Prompting GPT-3 to segment long events results in segmentation that is even closer to the human consensus.
- Obtained a Hamming distance of 0.223 and 0.191 with the "**Pieman**" story, when comparing to human annotations from the first and second behavioral run.

## COMPONENTS OF STORY



- Vertical bars show event boundaries detected by GPT-3 with (orange)



# 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 langauge models and why they are effective? [25 min]
  - Analyzing and Interpreting LLMs
- Text-to-Text Transfer Transformer, Prompting, **Instruction-tuning** [1 hour]
- Lunch [45 min]
- Large language models, Chain-of-Thought reasoning, RAG [1 hour 30 min]
- Coffee break [15 min]
- Hands on session [1 hour]

# 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.

<b>Natural language inference</b> (7 datasets)	<b>Commonsense</b> (4 datasets)	<b>Sentiment</b> (4 datasets)	<b>Paraphrase</b> (4 datasets)	<b>Closed-book QA</b> (3 datasets)	<b>Struct to text</b> (4 datasets)	<b>Translation</b> (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
<b>Reading comp.</b> (5 datasets)	<b>Read. comp. w/ commonsense</b> (2 datasets)	<b>Coreference</b> (3 datasets)	<b>Misc.</b> (7 datasets)	<b>Summarization</b> (11 datasets)		WMT-16 EN/DE
BoolQ	OBQA	DPR	CoQA	AESLC	Multi-News	WMT-16 EN/FI
DROP	SQuAD	CosmosQA	QuAC	CoLA	SamSum	WMT-16 EN/RO
MultiRC		ReCoRD	WIC	Math	AG News	WMT-16 EN/RU
		WSC273	Fix Punctuation (NLG)	Gigaword	Newsroom	WMT-16 EN/TR
					Opin-Abs: iDebate	XSum
					Opin-Abs: Movie	

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

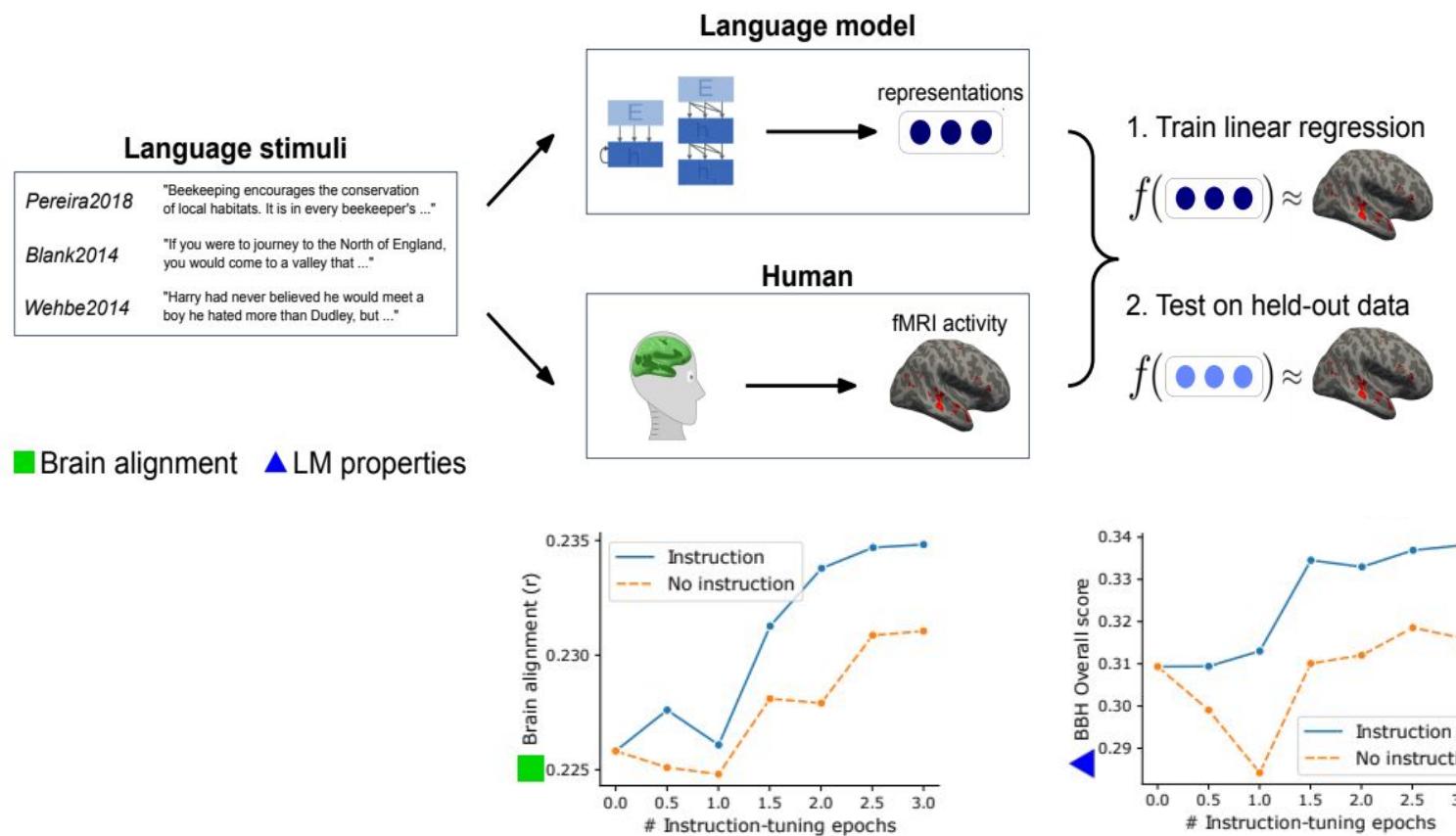
## Cross-lingual Reasoning

- 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.

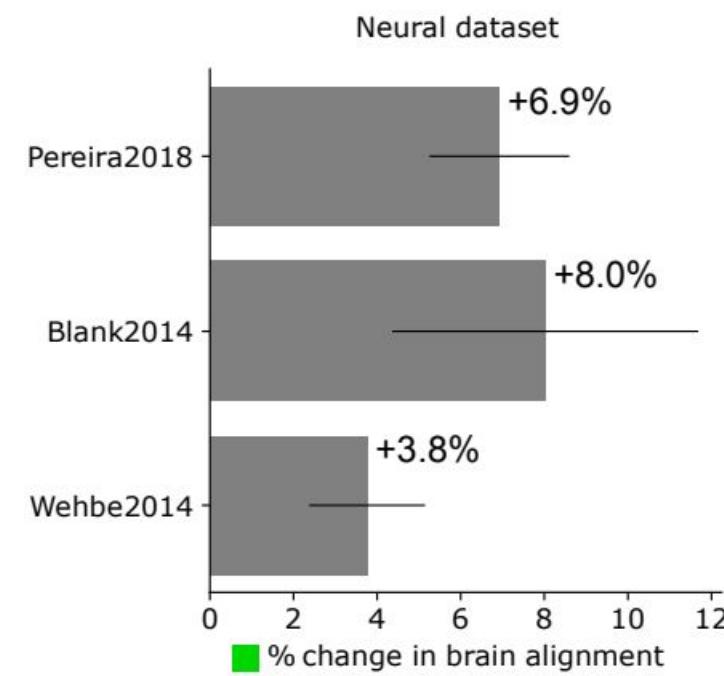
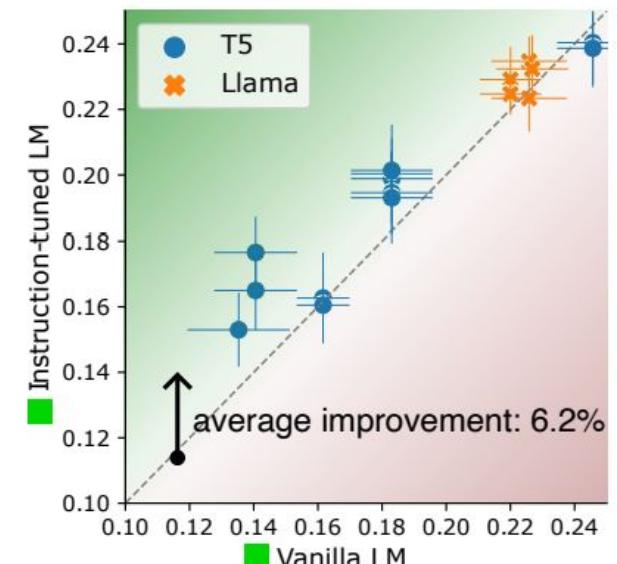
Task Type	Query & Options
Arithmetic-Int	<b>Query:</b> $2331 + 2693 = ?$ <b>Options:</b> 5024 (✓); 5018; 5005; 5025 <b>Query:</b> $109848 \div 199 = ?$ <b>Options:</b> 552.0 (✓); 516.0; 558.0; 567.0
Arithmetic-Flo	<b>Query:</b> $7.682 + 28.894 = ?$ <b>Options:</b> 36.576 (✓); 28.576; 40.909; 38.076 <b>Query:</b> $25.204 \times 88.29 \div 12.133 = ?$ <b>Options:</b> 183.406 (✓); 183.739; 185.406; 181.962
MPS-Cal	<b>Query:</b> 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? <b>Options:</b> 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	<b>Query:</b> 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? <b>Options:</b> 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

# INSTRUCTION-TUNING ALIGNS LLMS TO THE HUMAN BRAIN

- Stimuli: passages, narratives
- Stimulus representation: Instruction-tuned NLP models
- Brain recording & modality: fMRI, reading and listening of different stimuli



Average Brain alignment (Pearson corr.)



# 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. ✓

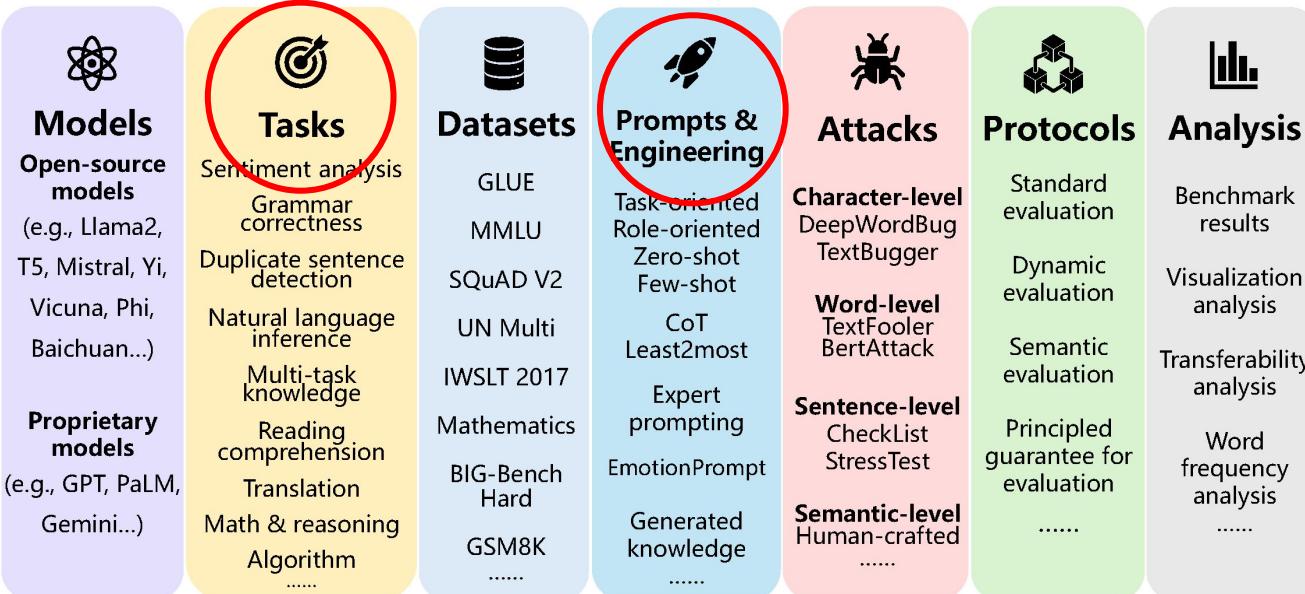
**Prompt**  
As a mathematics instreector, calculate the ansxer to the following problem related to if a number is a prime:  
  
**Sample**  
Question: Let  $z(a) = -871*a + 415$ . Is  $z(-16)$  a composite number? Answer:  
  
**User 2**  
No. ✗

(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.

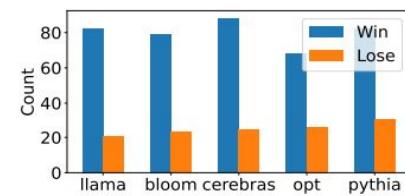
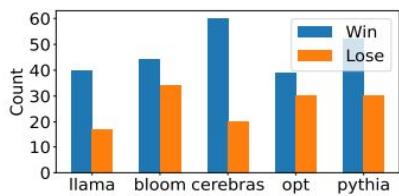
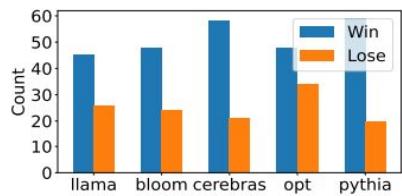
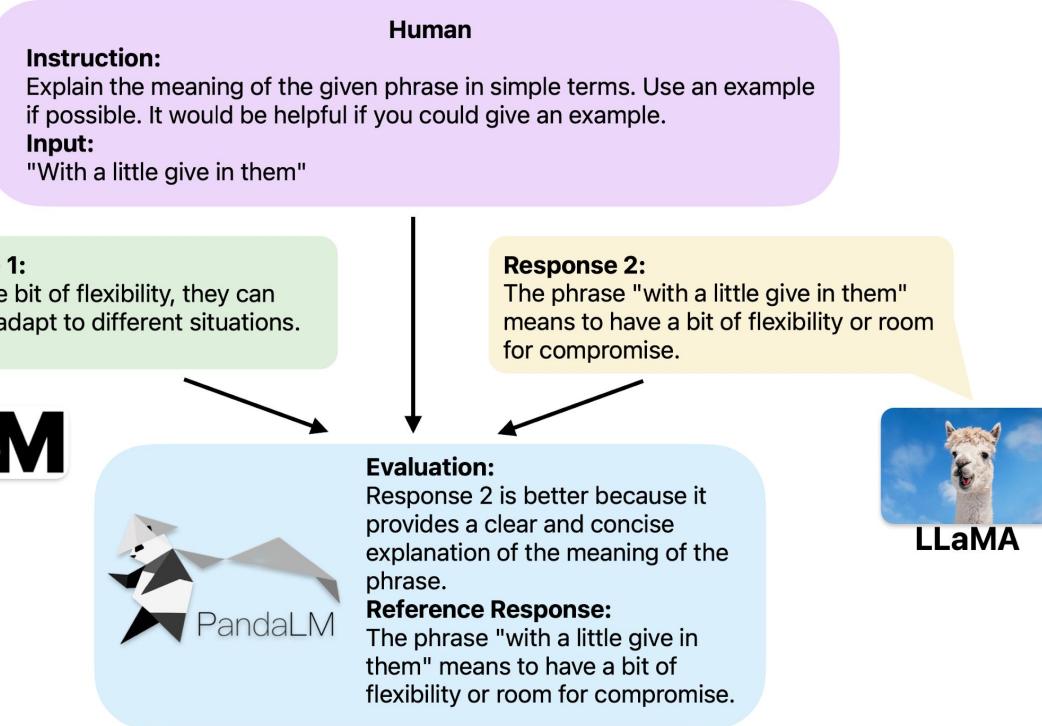


The APDR on different LLMs.

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

- 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

# PandaLM: Judge language model



(a) Comparison Results of GPT-3.5. (b) Comparison Results of GPT-4.

(c) Comparison Results of Human.

- Achieving 93.75% of GPT-3.5's evaluation ability and 88.28% of GPT4's in terms of F1-score on our diverse human annotated test dataset.

```
{  
    "index": "749",  
    "motivation_app": "CNN News",  
    "task_id": "user_oriented_task_165",  
    "cmp_key": "opt-7b_pythia-6.9b", ## It means response 1 is from opt-7B and response 2 is from pythia-6.9b  

```

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

Weixin Liang et al.

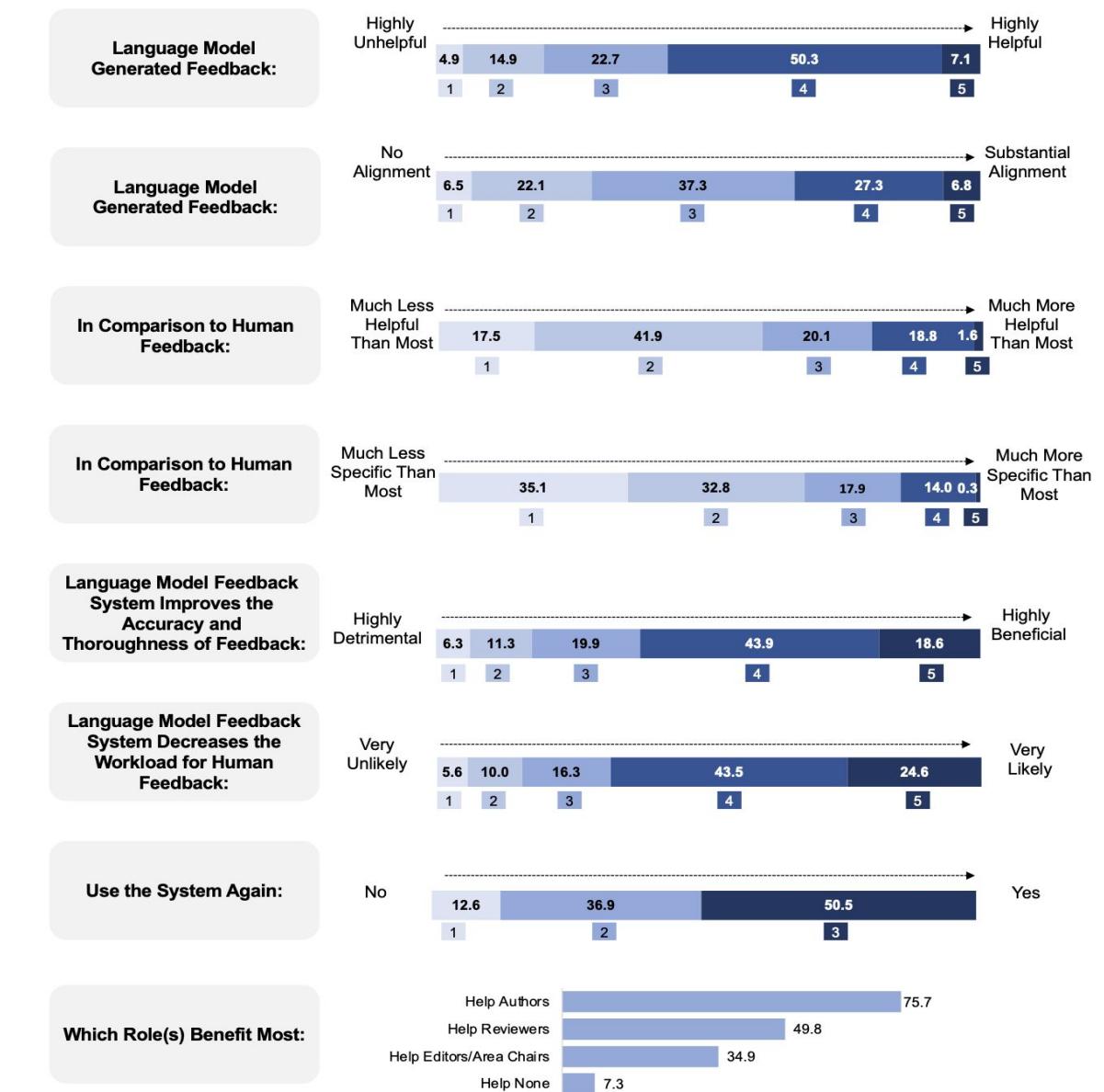
<https://arxiv.org/abs/2310.01783>

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.



Lunch Break – 12:25 to 1:00 pm

# QR Code for Github Link!

<https://github.com/mounikamarreddy/Lamarr-Tutorial---13.03.2024/tree/main>

