

**INTERPRETABILITY  
IN NEUROSCIENCE**

Prof. Alexander Huth

2026.1.30

# LAST TIME

- \* *Mechanistic interpretability*
- \* David Marr's 3 levels
- \* Linearity, modularity
- \* World models?

# TODAY

- \* How do we infer what the **goal** of a neural system is? (Kanwisher 2023; Schrimpf 2021; Antonello 2024)
- \* Can we build **explicitly interpretable** systems? (Benara 2024; Singh 2025)

# DISCOVERING GOALS

*David Marr's 3 levels of analysis:*

Computational – Algorithmic – Implementational

***Why?***

***What?***

***How?***

- \* How do we make inferences about the **computational** level?
- \* What do we actually **observe**?

# DISCOVERING GOALS

*David Marr's 3 levels of analysis:*

Computational – Algorithmic – Implementational

***Why?***

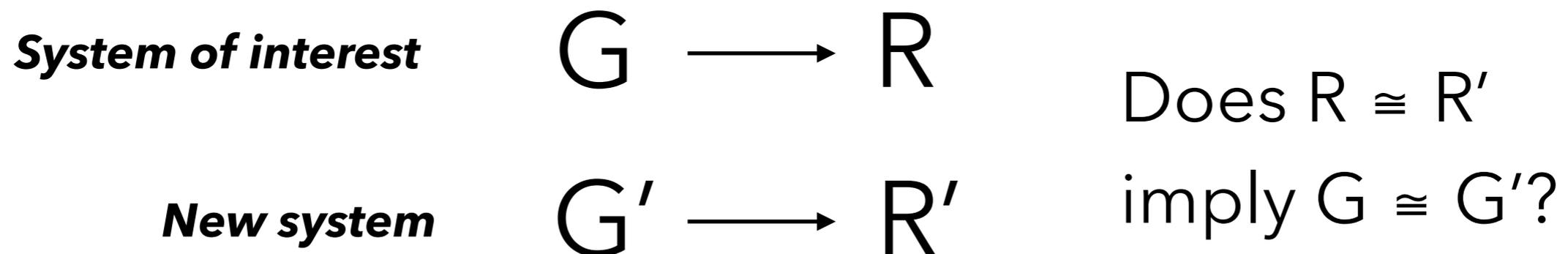
***What?***

***How?***

- \* We can observe **representations** (the algorithmic level)
- \* How do we swim upstream to the computational level?

# DISCOVERING GOALS

- \* Let's build a parallel system with a known computational goal
- \* Then we can compare its internal representations to those of our system of interest
- \* If representations are similar, we can infer that the goals are shared?

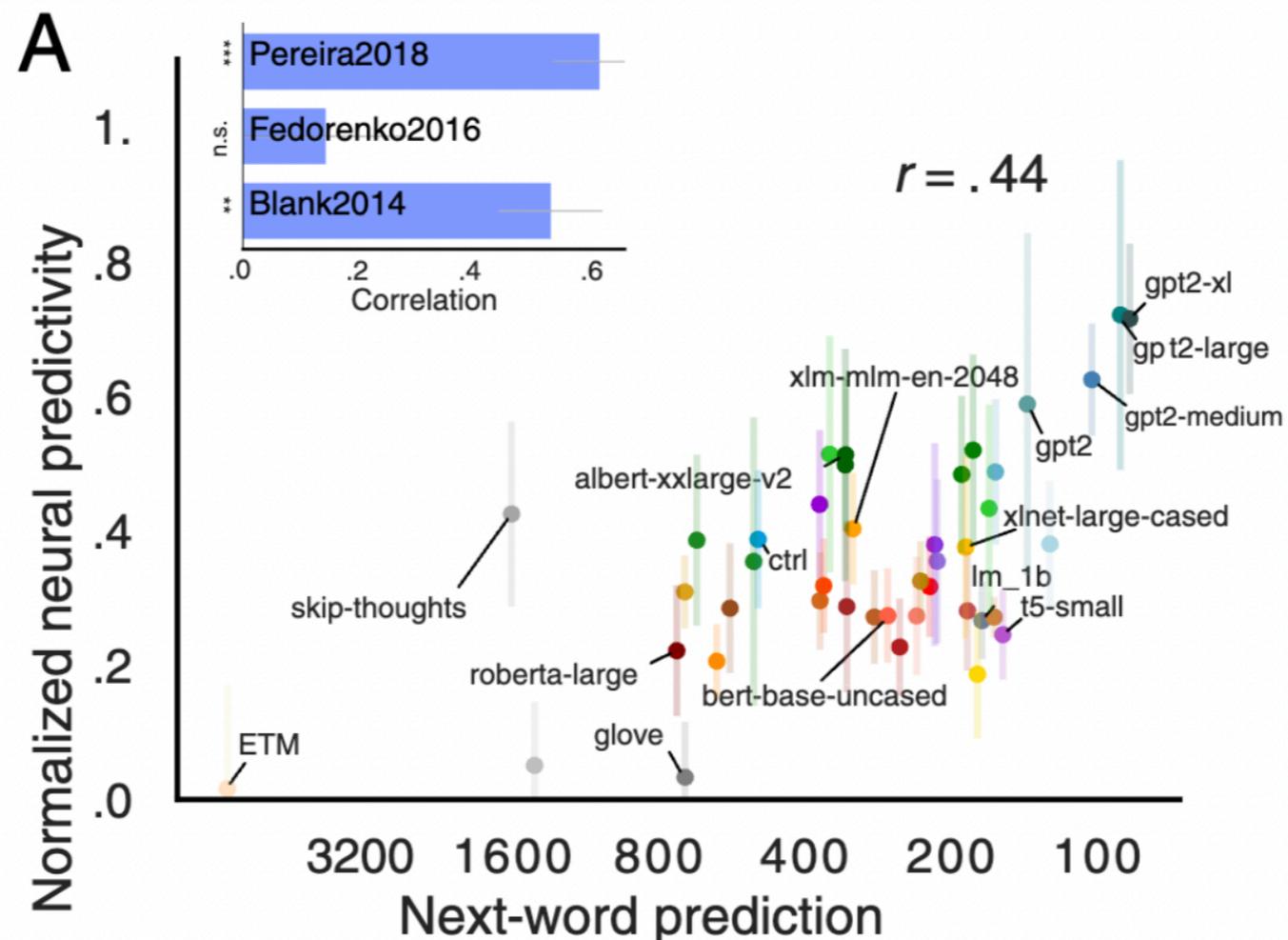


# DISCOVERING GOALS

- \* **Example:** our brains process language, but we do not know what the computational goal of our language processing system is
- \* We can build language processing systems where we know exactly what the computational goal is. **Language models** have the goal of **next-word prediction**
- \* Representations from LMs are very good at predicting representations in the brain
- \* Therefore the computational goal of our brain is next-word prediction

# DISCOVERING GOALS

- \* Or, more specifically: models that are better at next-word prediction produce representations that are better at predicting brain data



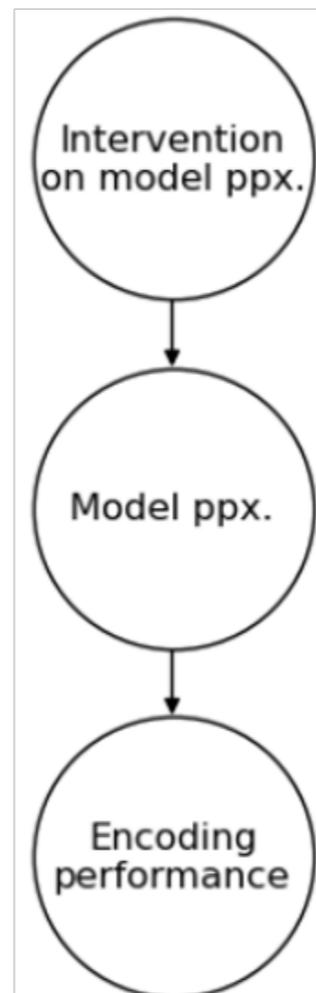
# DISCOVERING GOALS

- \* What's wrong with this type of inference?
- \* Can the **same** representations arise while pursuing **different** computational goals?

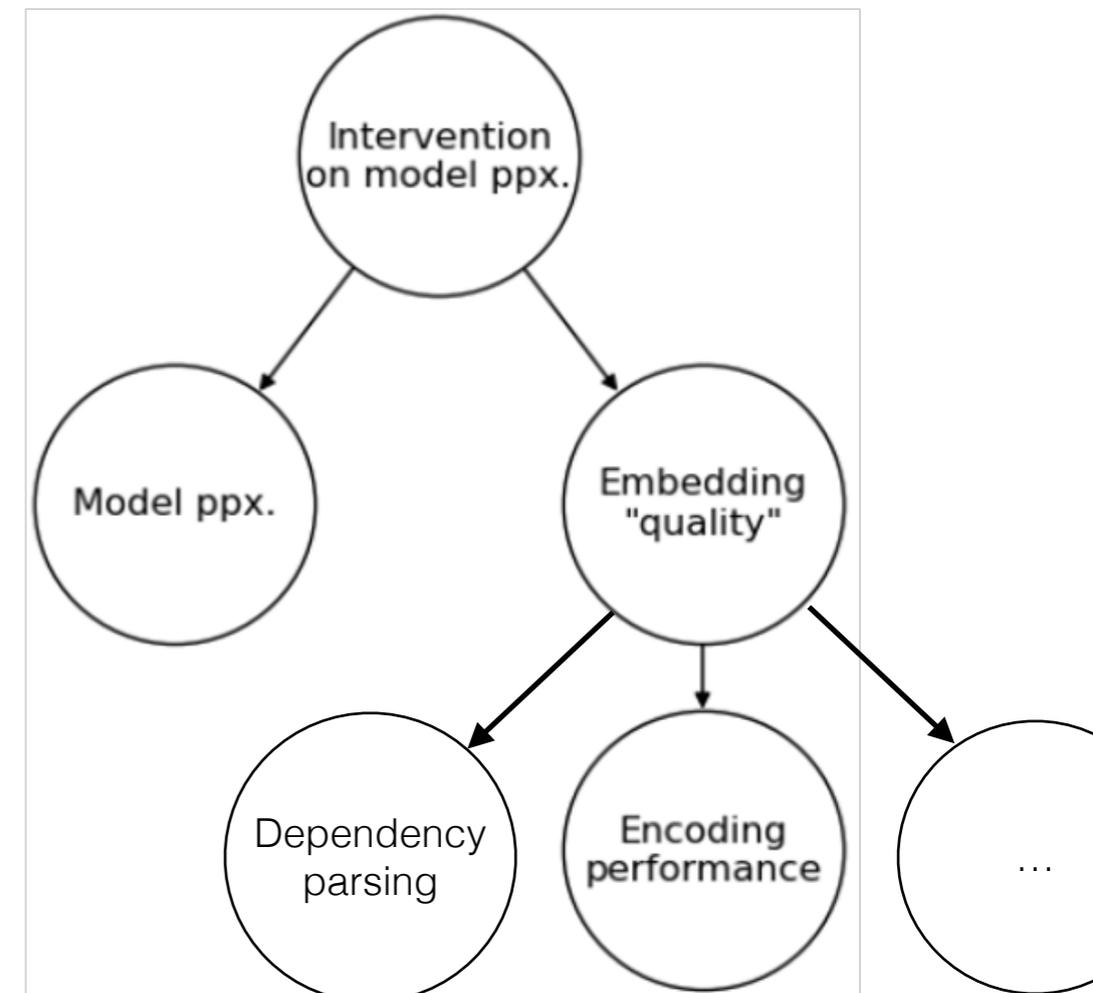
# DISCOVERING GOALS

- Two worlds:

“Predictive Coding Hypothesis”



“Useful Embedding Hypothesis”



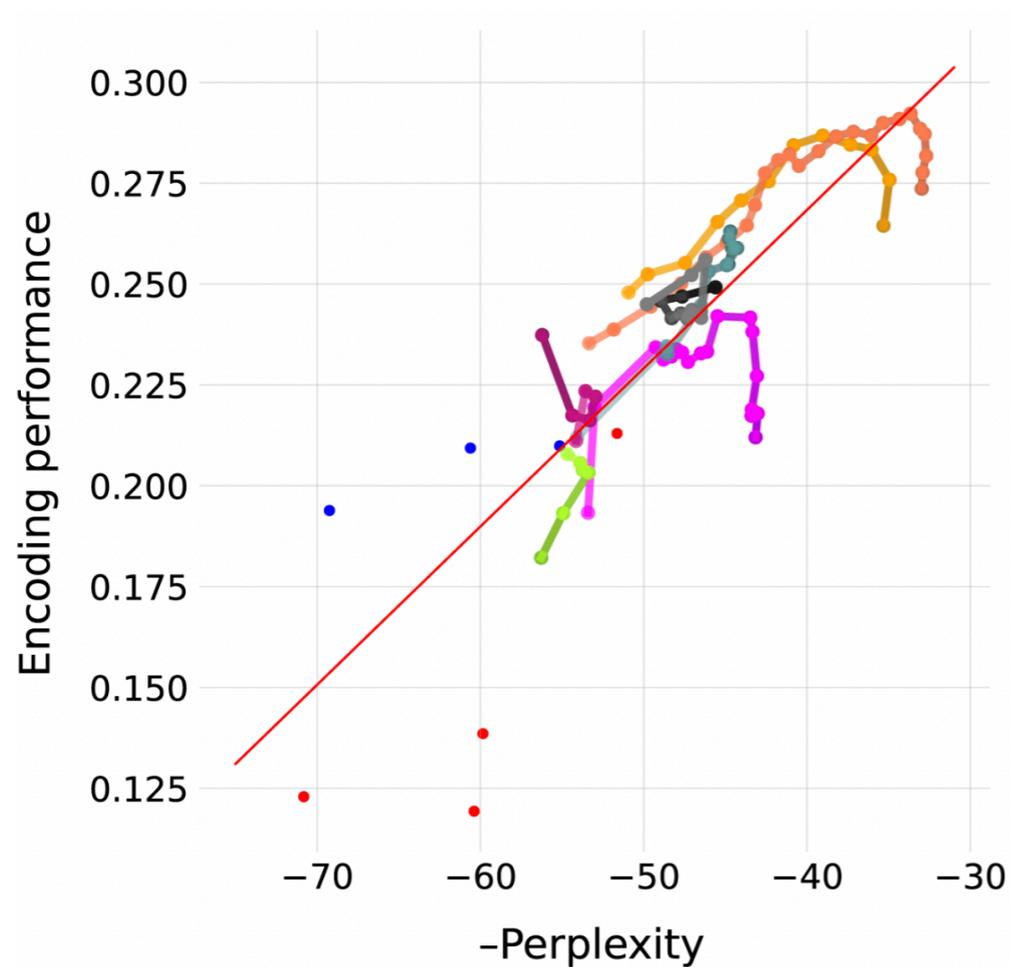
# DISCOVERING GOALS

- \* How do we measure “embedding quality” or “generality”?
- \* Create a set of embeddings (representations), then measure how well each one predicts each other one (*Antonello et al., NeurIPS 2021*)

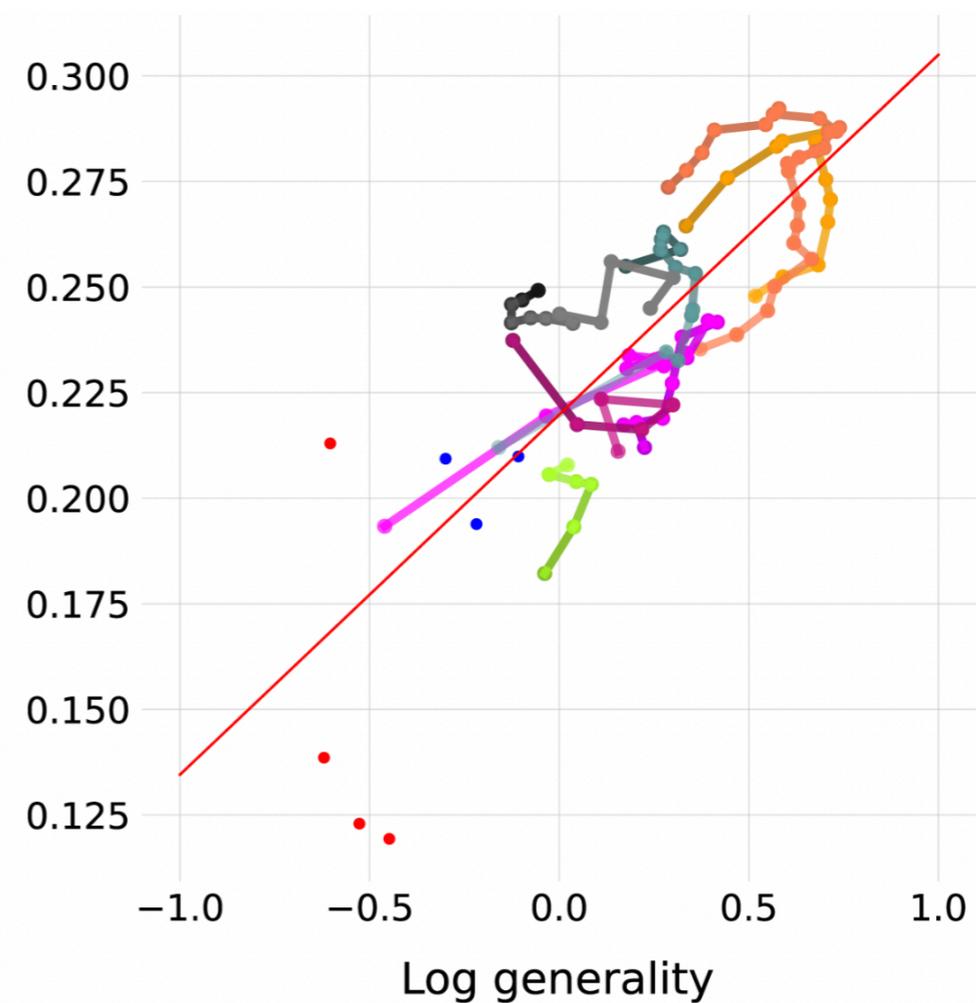
# DISCOVERING GOALS

- GPT-2 Small
- GPT-2 Medium
- Transformer-XL
- BERT
- ALBERT
- Eng → De
- Eng → Zh
- Word embeddings
- Interpretable

*Replication of perplexity effect:*



*Relationship with “generality”:*



# DISCOVERING GOALS

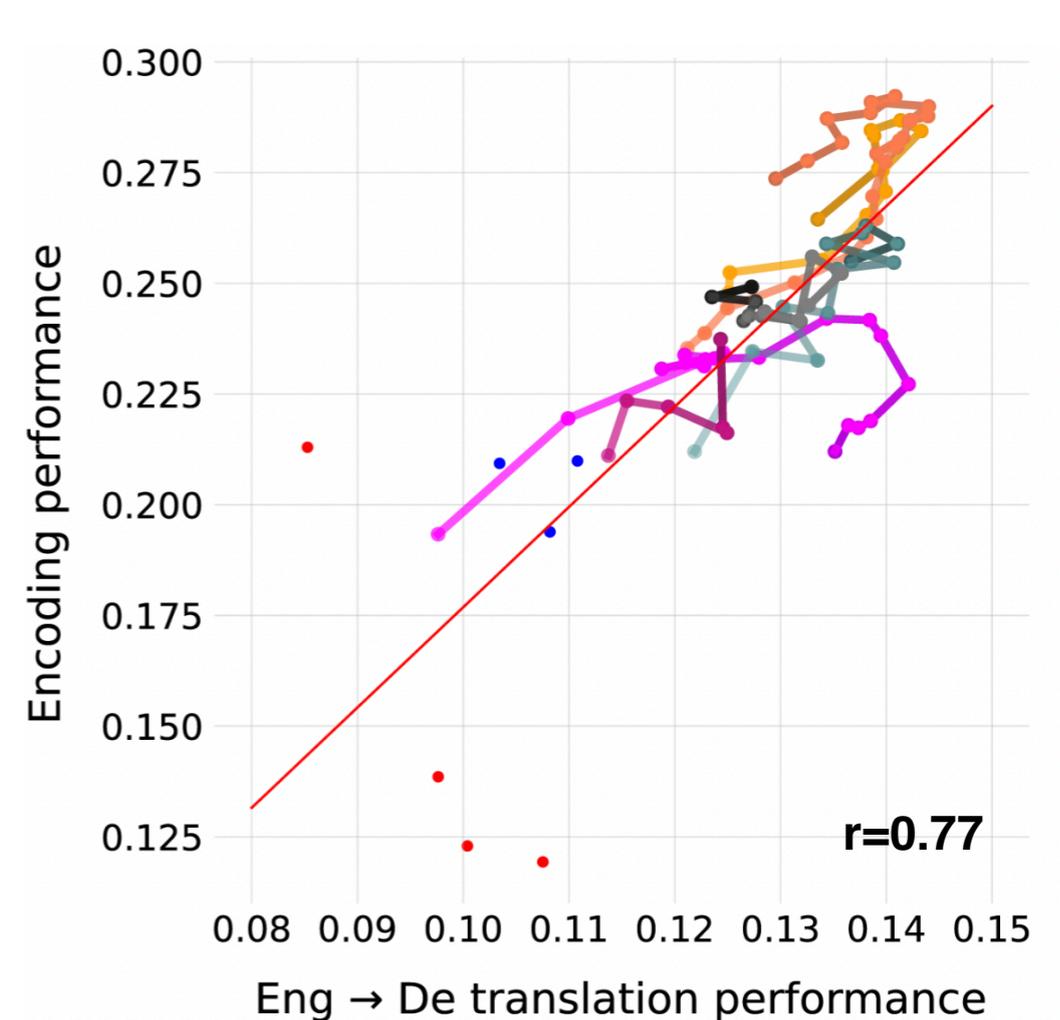
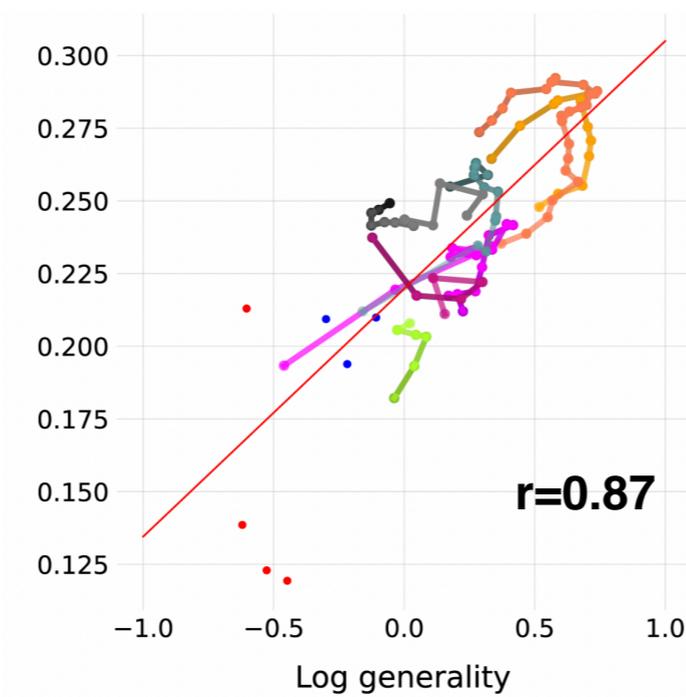
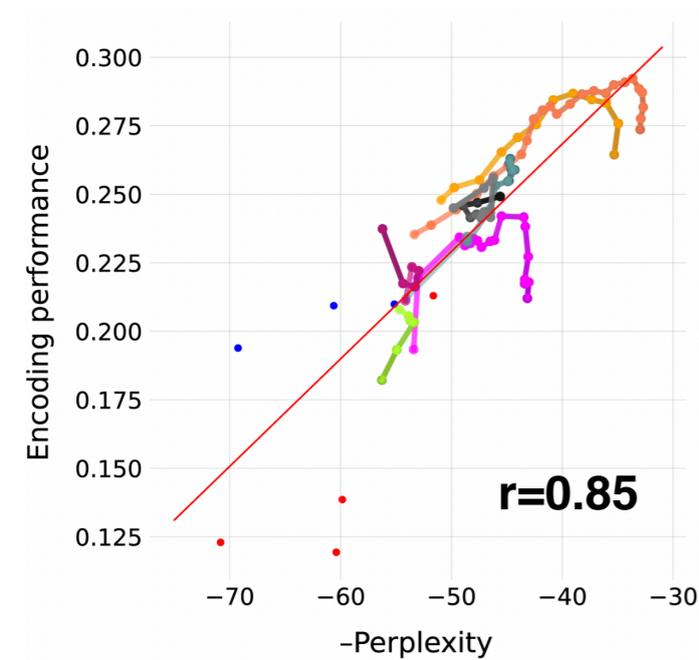
- \* Reductio ad absurdum: what if the computational goal of the brain was to translate incoming language to German?
- \* Then similarity to a translation model should be correlated with similarity to the brain!

# DISCOVERING GOALS

- GPT-2 Small
- GPT-2 Medium
- Transformer-XL
- BERT
- ALBERT
- Eng → De
- Word embeddings
- Interpretable
- Eng → Zh

Replication of perplexity effect: Relationship with “generality”:

Relationship with Eng-De translation:



# DISCOVERING GOALS

- \* Upshot: we should be skeptical about claims that representational similarity entails mechanistic similarity
- \* But does this also tell us something deeper? What makes some representations “general”?
- \* Platonic representation hypothesis (*Huh, Cheung, Wang & Isola (2024) ICML*)

# INTERPRETABLE SYSTEMS

- \* What does it mean for a representation to be **interpretable**?

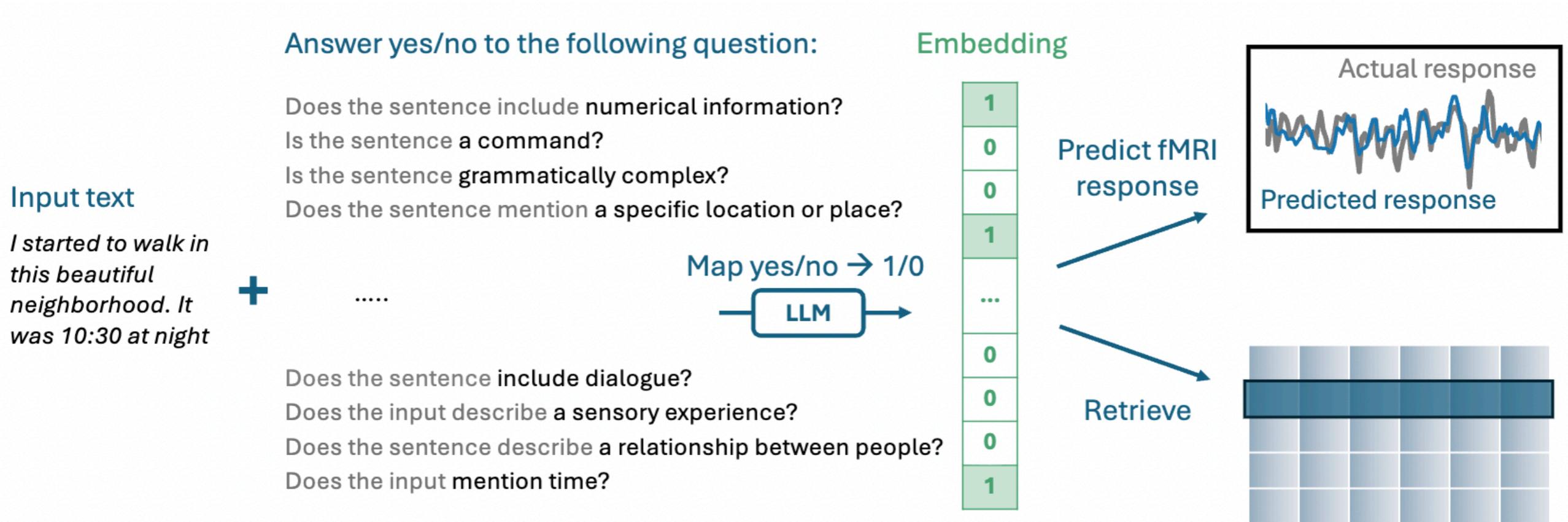
# INTERPRETABLE SYSTEMS

- \* What does it mean for a representation to be **interpretable**?
- \* One possibility: *it can be expressed simply and completely in words*

# INTERPRETABLE SYSTEMS

- \* What does it mean for a representation to be **interpretable**?
- \* One possibility: *it can be expressed simply and completely in words*
- \* What if we start by expressing features in words, then turn them into a representation?

# INTERPRETABLE SYSTEMS



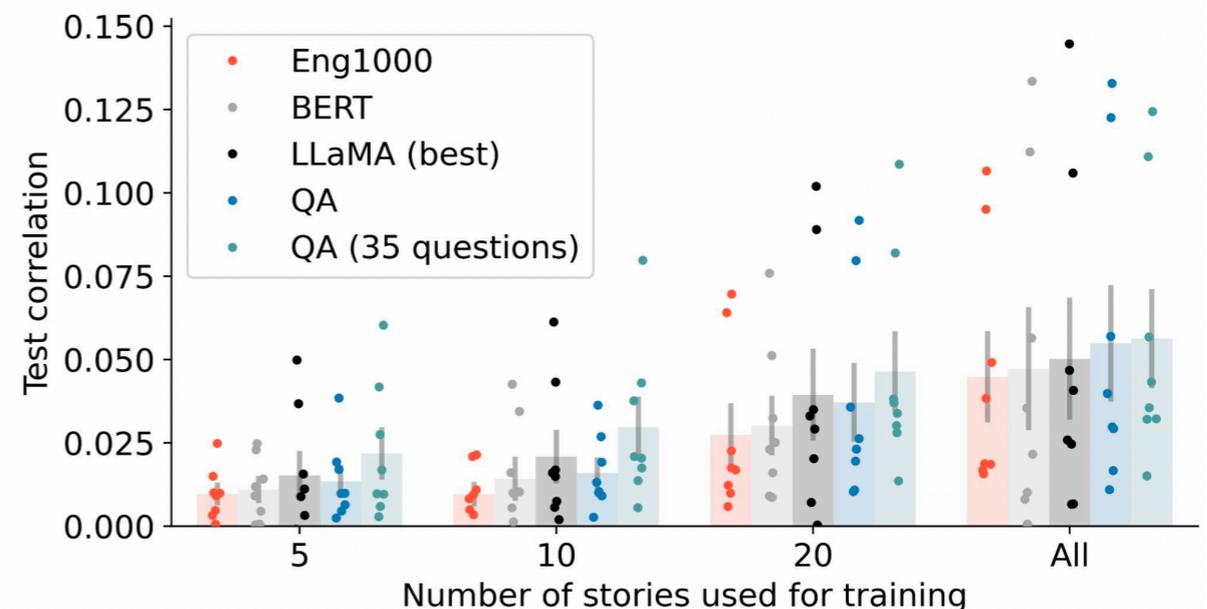
# INTERPRETABLE SYSTEMS

\* Are these explicitly interpretable embeddings **good**? (*Yeah they're pretty good!*)

## Information retrieval task

	Mean reciprocal rank	Recall@1	Recall@5	Size
Bag of words	$0.37 \pm 0.01$	$0.28 \pm 0.02$	$0.42 \pm 0.02$	27,677
Bag of bigrams	$0.39 \pm 0.01$	$0.30 \pm 0.02$	$0.44 \pm 0.02$	197,924
Bag of trigrams	$0.39 \pm 0.02$	$0.30 \pm 0.02$	$0.44 \pm 0.02$	444,403
QA-Emb	$0.45 \pm 0.01$	$0.34 \pm 0.01$	$0.50 \pm 0.01$	† <b>2,000</b>
BM-25	$0.77 \pm 0.01$	$0.69 \pm 0.01$	$0.82 \pm 0.01$	27,677
<b>BM-25 + QA-Emb</b>	<b><math>0.80 \pm 0.01</math></b>	<b><math>0.71 \pm 0.01</math></b>	<b><math>0.84 \pm 0.01</math></b>	29,677

## Brain prediction task



# INTERPRETABLE SYSTEMS

- \* Using natural language as a representation creates explicit interpretability. Is this broadly applicable?
- \* Related: **chain-of-thought** prompting in LLMs/ chatbots (*Wei et al., NeurIPS 2022*)
- \* Does this offer interpretability? (*Barez et al. 2025*)

# TODAY

- \* How do we infer what the **goal** of a neural system is? (Kanwisher 2023; Schrimpf 2021; Antonello 2024)
- \* Can we build **explicitly interpretable** systems? (Benara 2024; Singh 2025)