



Large Language Models Are Aligned With The Human Language System

Martin Schrimpf

PATT, Schools of Life Sciences, and of
Computer and Communication Sciences

Neuro_X Institute

EPFL



martin.schrimpf@epfl.ch



[@martin_schrimpf](https://twitter.com/martin_schrimpf)



[mastodon.social/
@mschrimpf](https://mastodon.social/@mschrimpf)



Idan Blank



Greta Tuckute



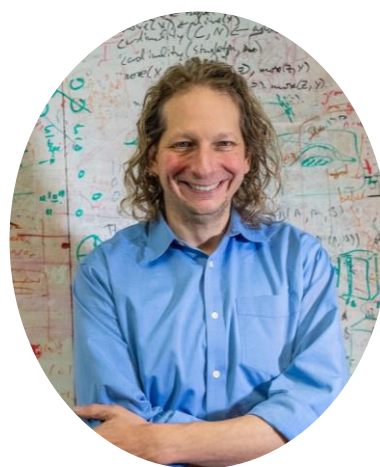
Carina Kauf



Eghbal Hosseini



Nancy Kanwisher

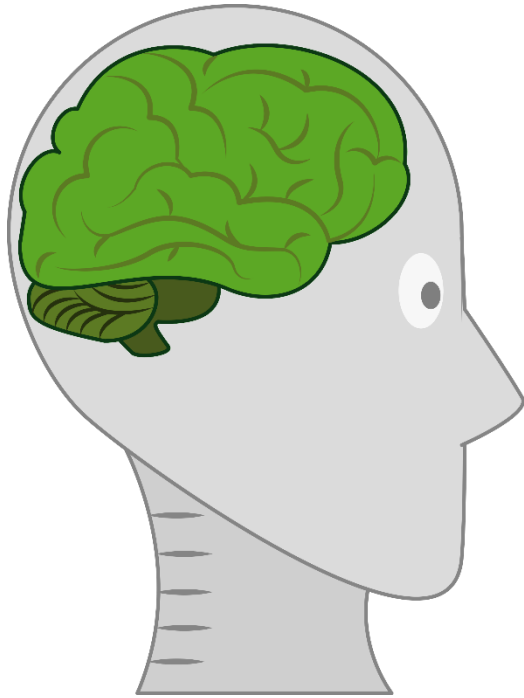


Josh Tenenbaum

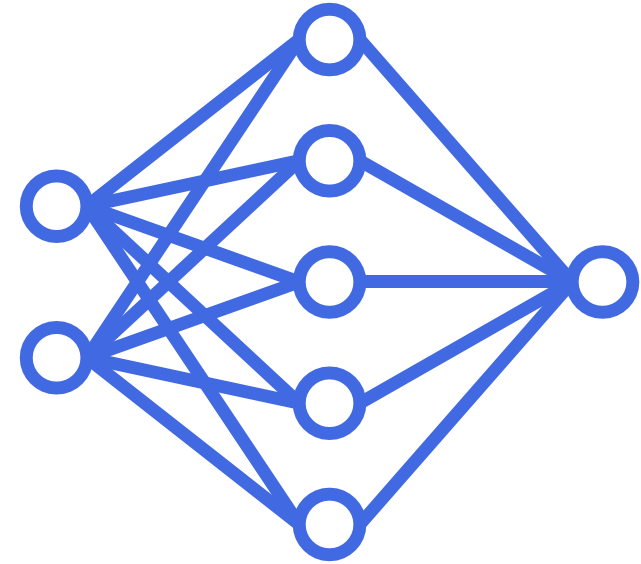


Ev Fedorenko

Goal (broadly): Model Natural (Human) Intelligence and the Underlying Neural Mechanisms



\approx

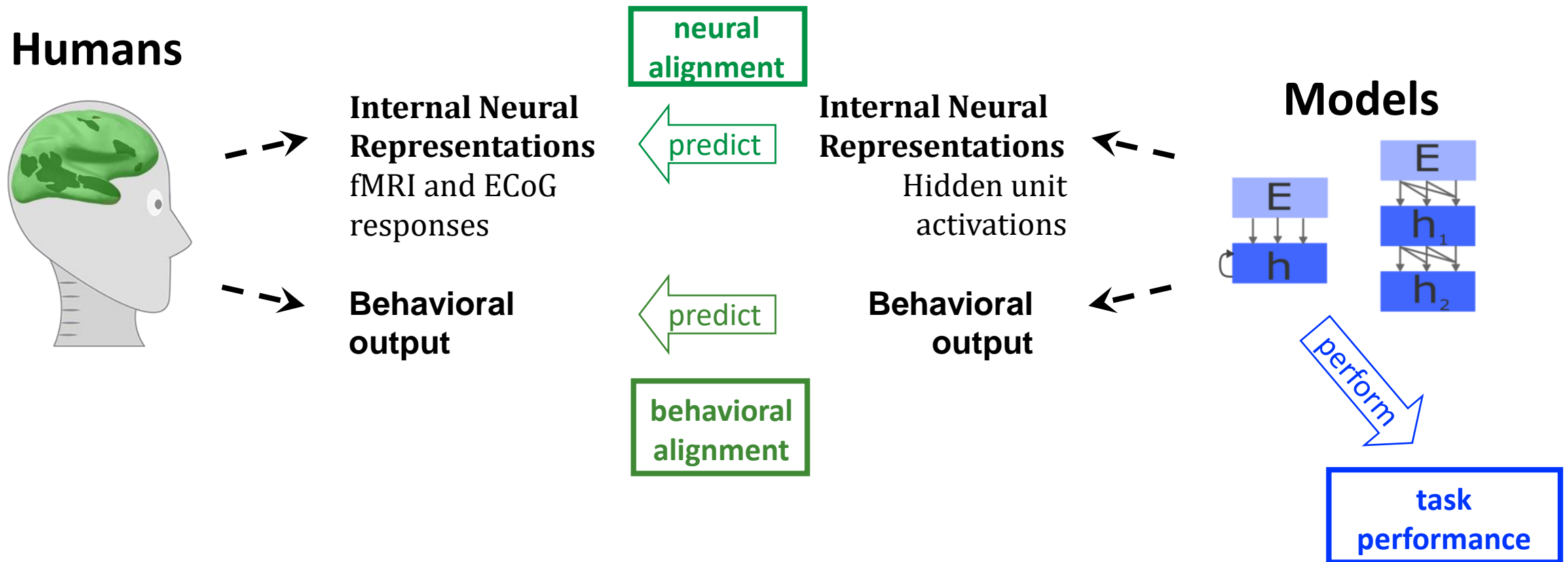


Mechanistic
understanding of
human intelligence

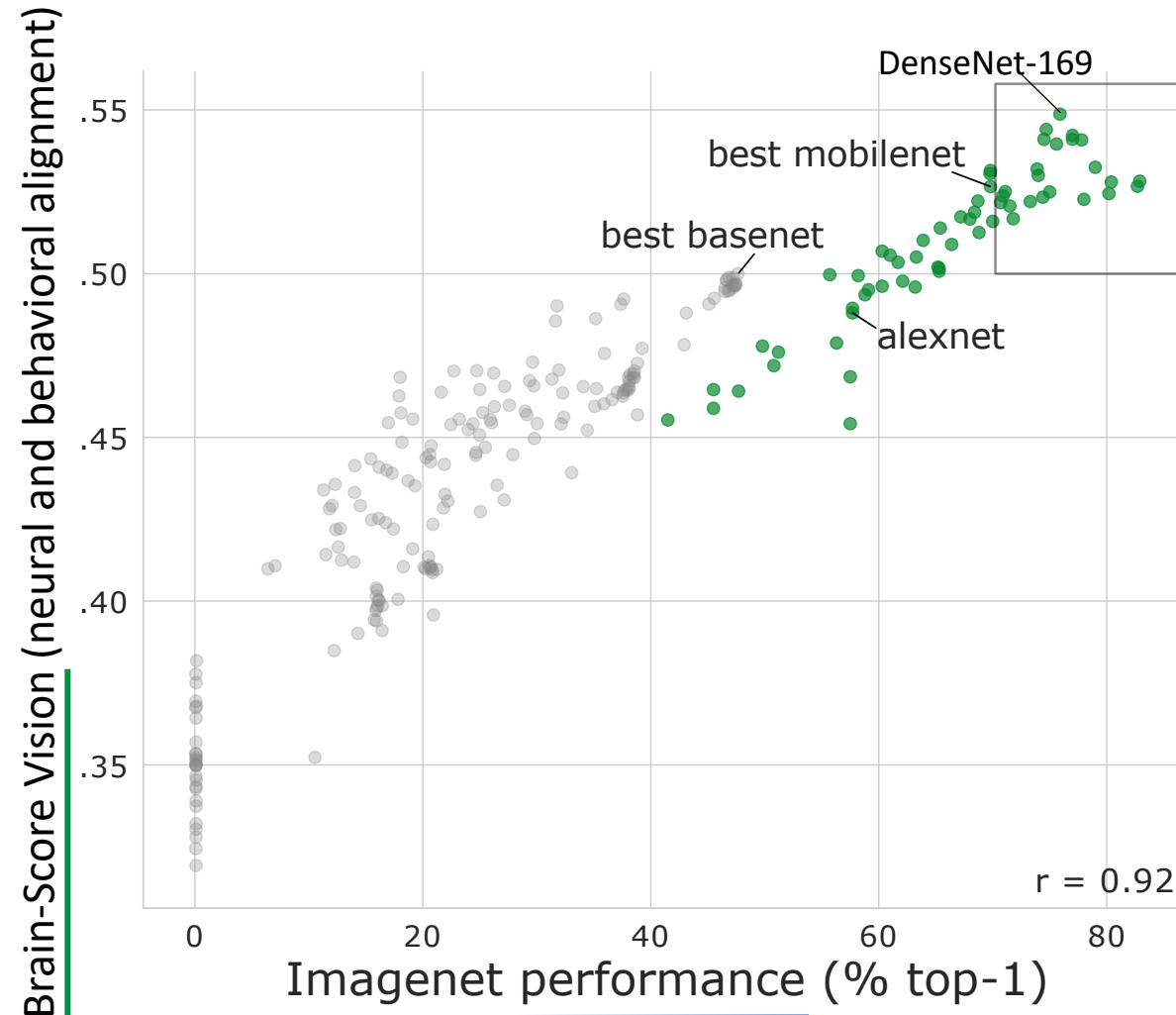
Next-generation
intelligence algorithms

Future clinical
applications

Goal (today): Model the Human Language System



What kinds of models could align with the human language system?



In sensory cortex:

- Artificial Neural Networks (ANNs) are the leading class of models for explaining brain and behavior
- ANNs make predictions for any visual input and work well for real-world stimuli
- ANNs with higher task performance generally are more aligned to brain and behavior

Schrimpf*, Kumbus*, et al. 2018 | Kumbus*, Schrimpf*, et al. 2019

see also www.brain-score.org | Yamins*, Hong*, et al. 2013, 2014 |

Khaligh-Razavi & Kriegeskorte 2014 | Zhuang et al. 2017 | Kell et al. 2018

Modeling higher cognition

Perception

Language

High-level
reasoning

Artificial neural networks have worked well in modeling sensory cortex – could they also predict higher cognition?

The human language network

working definition:

a set of **left-lateralized** regions on the lateral surfaces of **frontal** and **temporal** cortex that support **high-level** language processing.

Language

>

Perceptually
matched control

Sentences

>

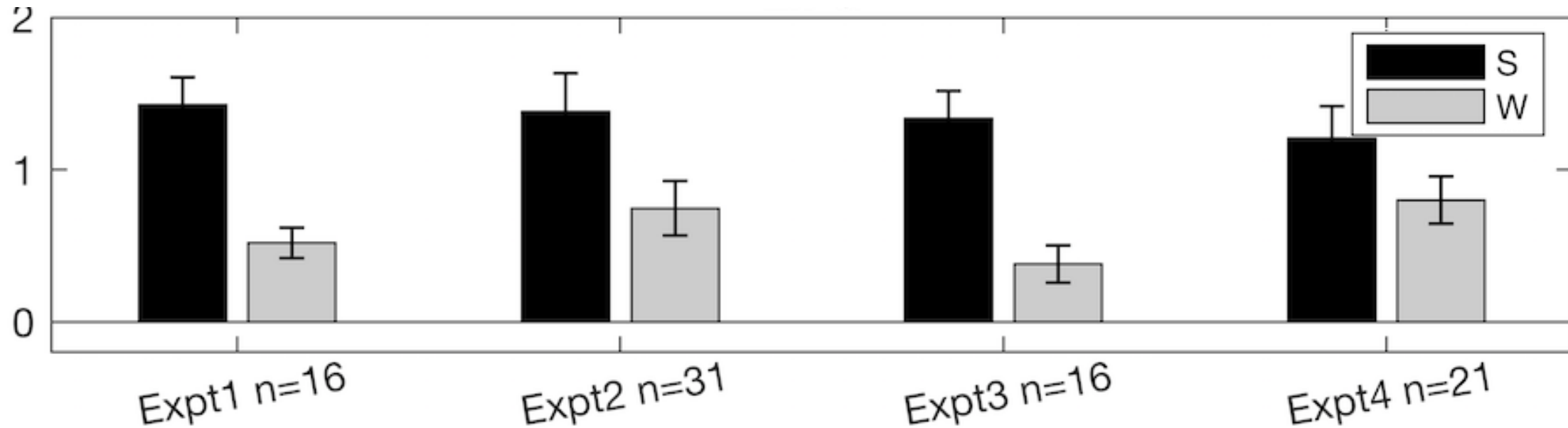
Lists of nonwords

The human language network

the dog is taking
a bath

>

dap drello smop ub
plid kav



Key signature: stronger response to sentences than lists of unconnected words

The human language network



courtesy of Idan Blank

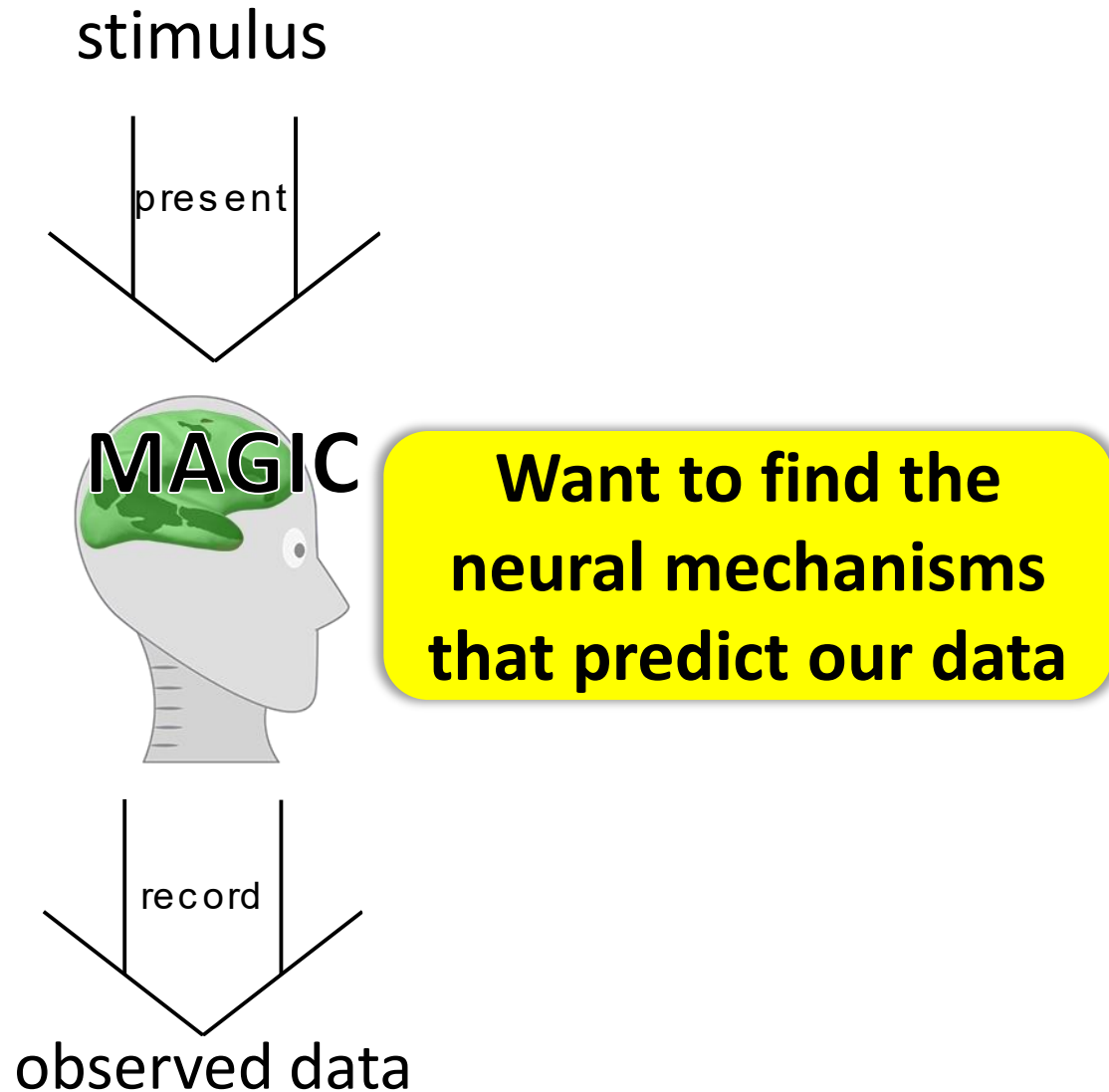
What are the mechanisms underlying human language comprehension?

the dog is taking
a bath



"meaning"

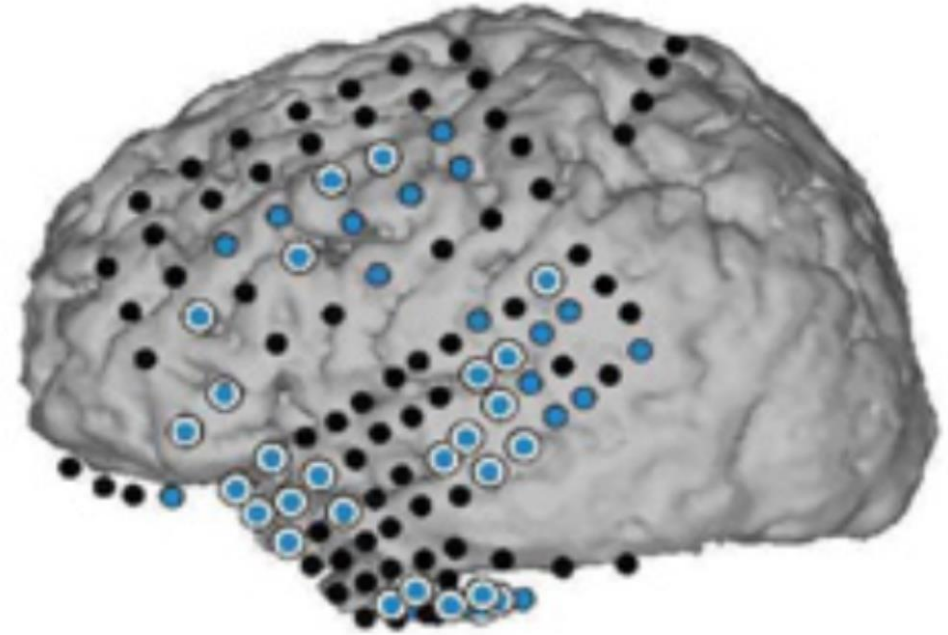
What are the mechanisms underlying human language comprehension?



Data target: human neural recordings



fMRI



ECoG

Data target: human neural recordings

Pereira et al. 2018

fMRI



627 sentences x 13,517 voxels in 10 subjects

Beekeeping encourages the conservation of local habitats. | It is in every beekeeper's interest ...

Fedorenko et al. 2016

ECoG



416 words x 97 electrodes in 5 subjects

ALEX / WAS / TIRED / SO / HE / TOOK / A / NAP

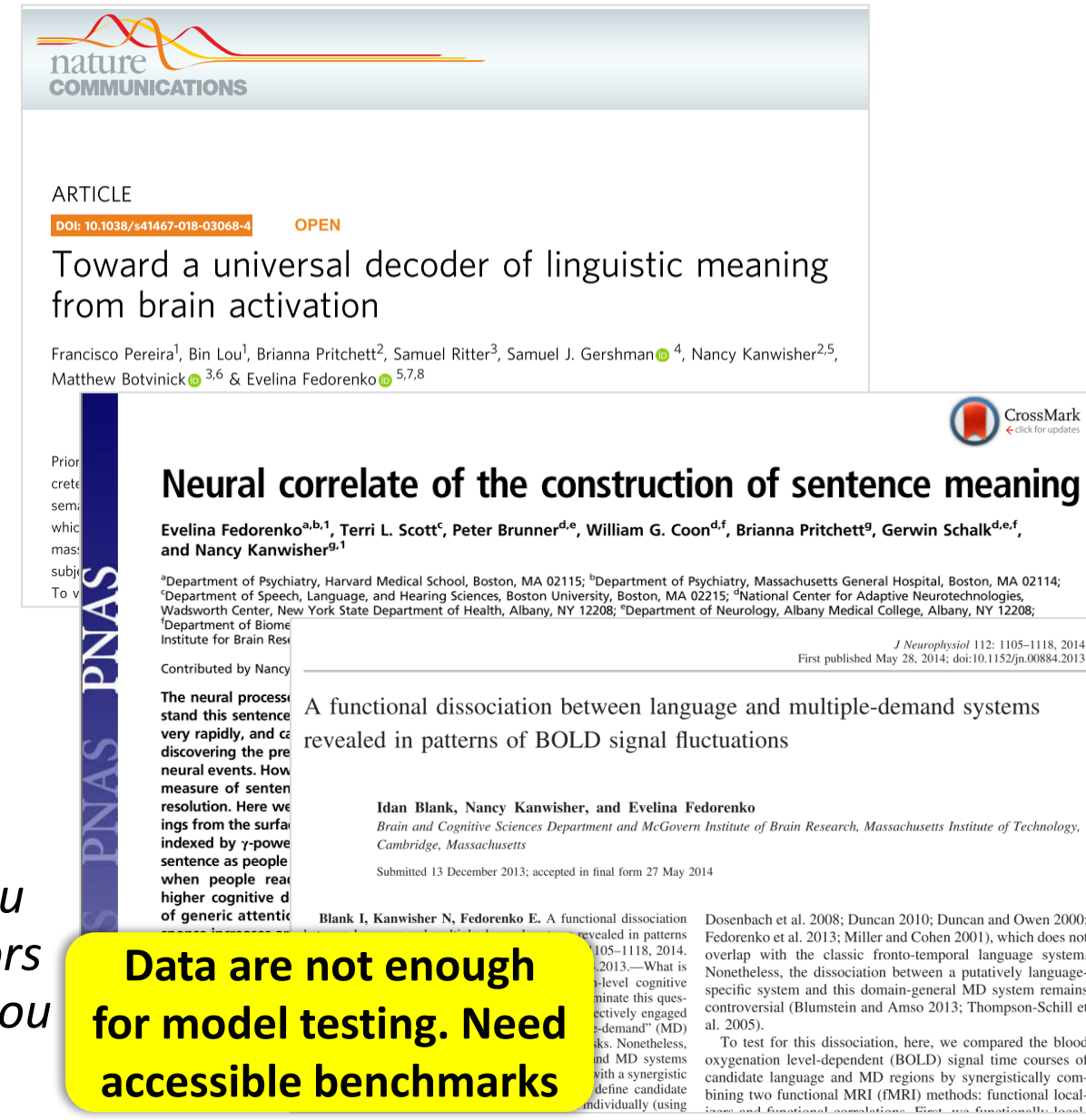
Blank et al. 2014

fMRI

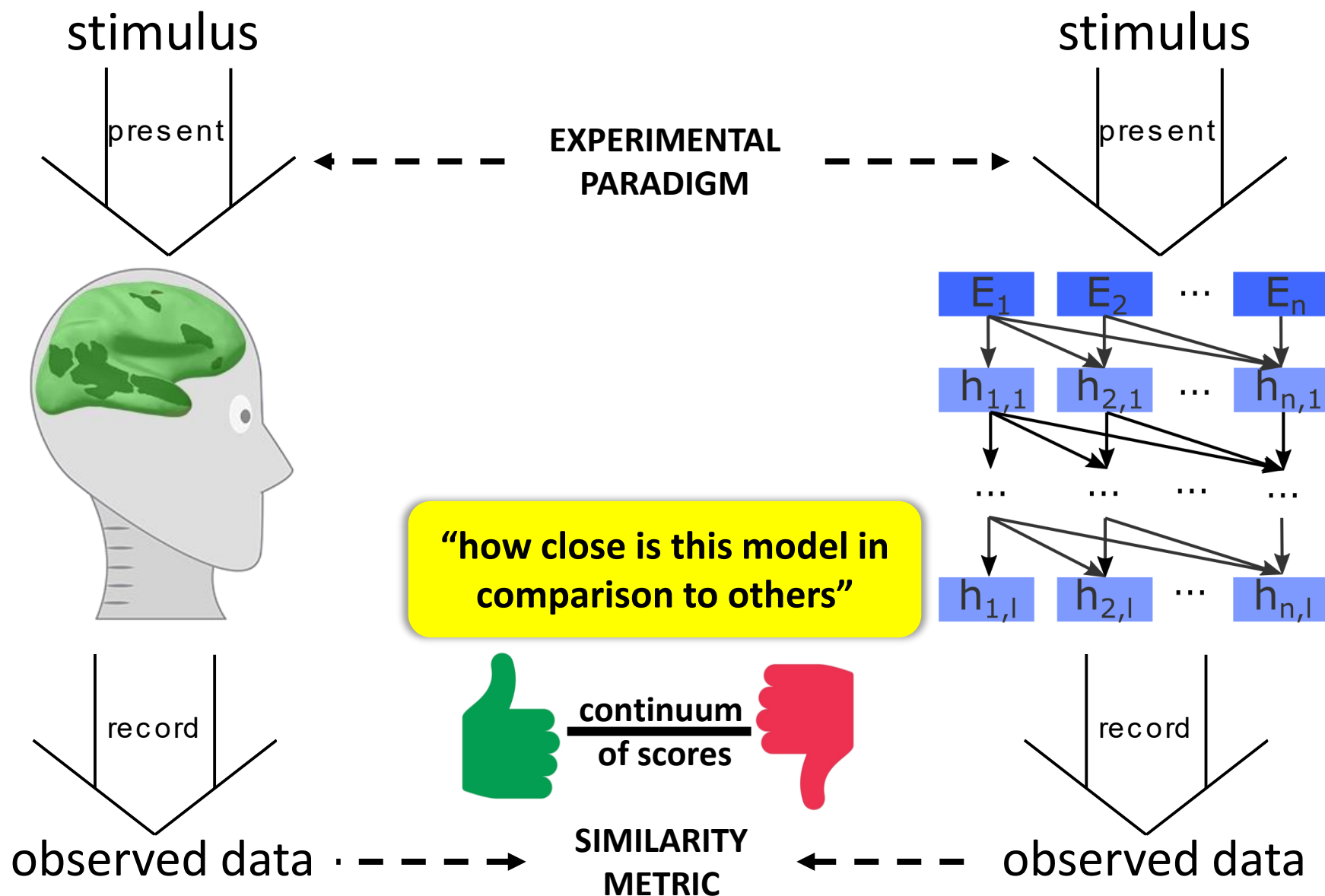


1,317 story fragments x 60 fROIs in 5 subjects

*If you were to journey to the | North of England, you
would come to a valley | that is surrounded by moors
as high as | mountains. It is in this | valley where you
would find the city of Bradford, | ...*



Quantifying match-to-brain: Benchmarking



Quantifying match-to-brain: Benchmarking

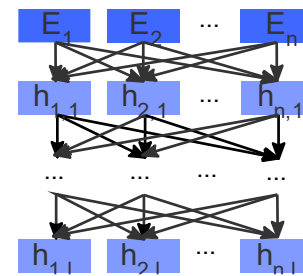
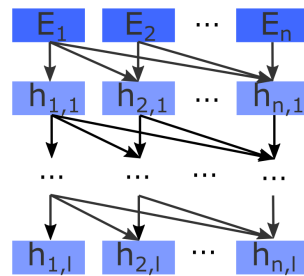
We only care about best-matching model (for now)



match-to-brain

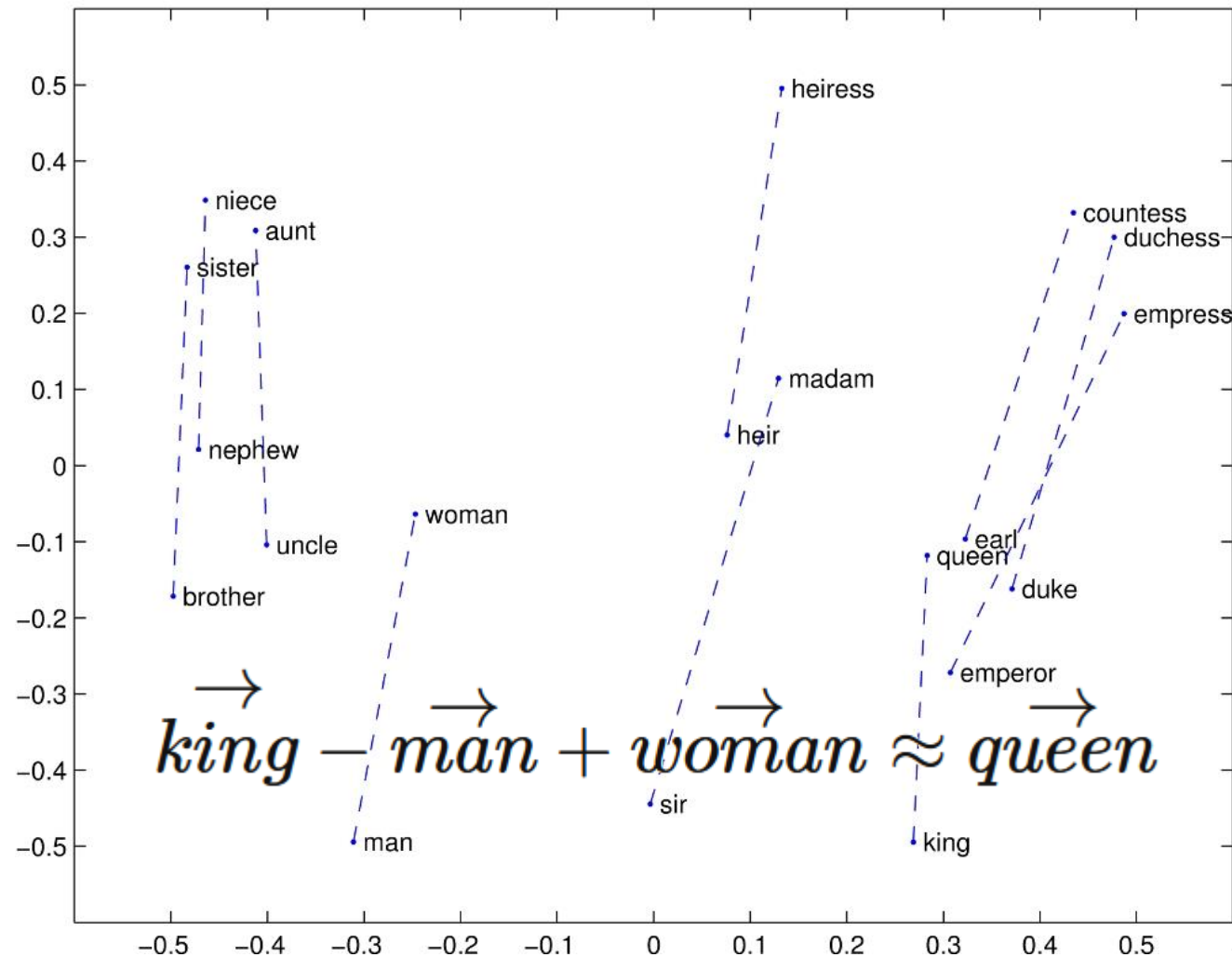
how close
are we?

differentiate
models



Models tested (n=43)

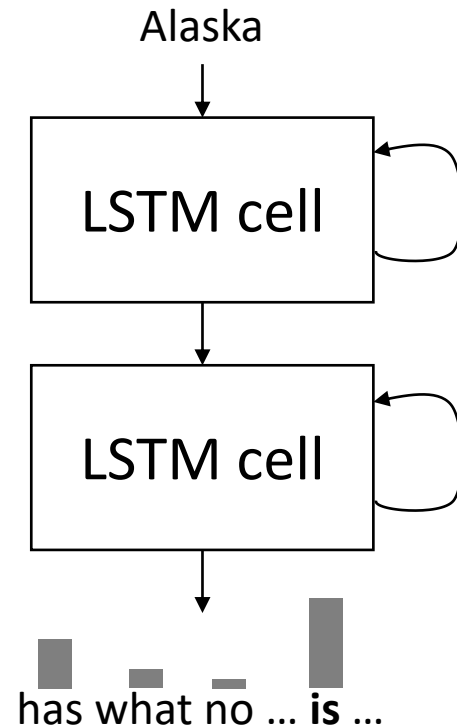
Embedding type models: GloVe, word2vec, topicETM



Models tested (n=43)

Embedding type models: GloVe, word2vec, topicETM

Recurrent networks: skip-thoughts, LSTM lm_1b



Language Modeling

Alaska **is**

Alaska is **about**

Alaska is about **twelve**

Alaska is about twelve **times**

Alaska is about twelve times **larger**

Alaska is about twelve times larger **than**

Alaska is about twelve times larger than **New**

Alaska is about twelve times larger than New **York**

Image from <https://www.quora.com/What-is-a-masked-language-model-and-how-is-it-related-to-BERT>

Models tested (n=43)

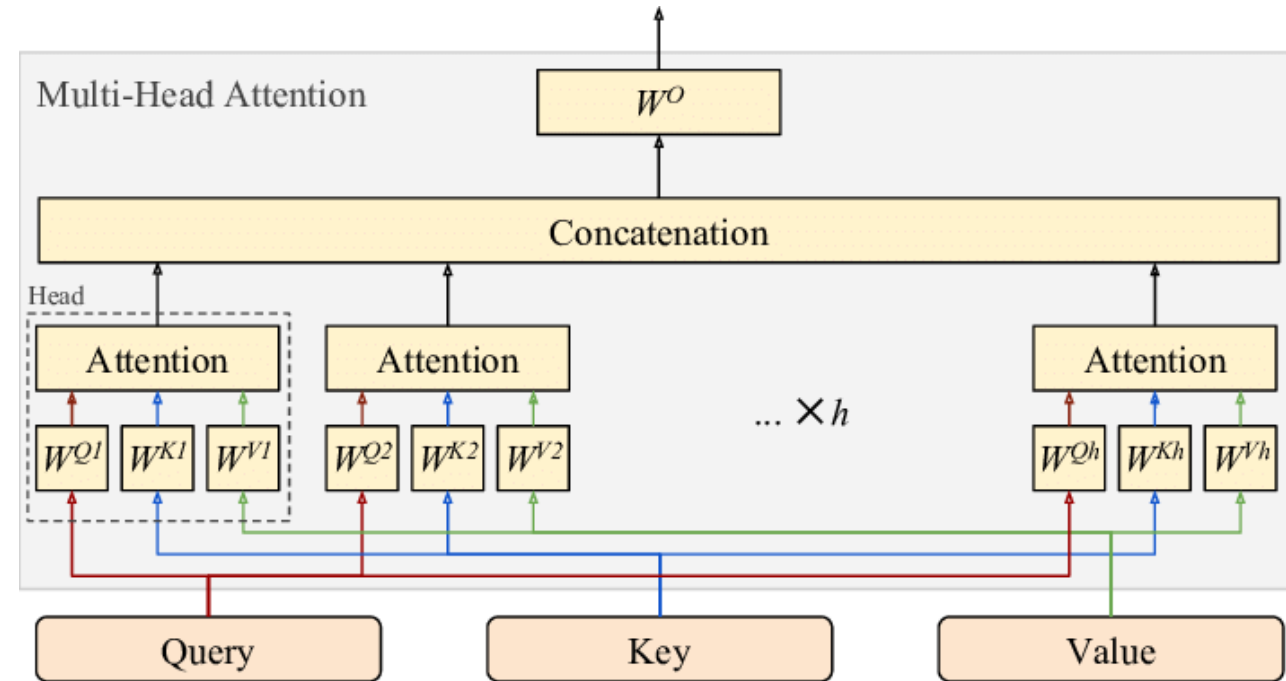
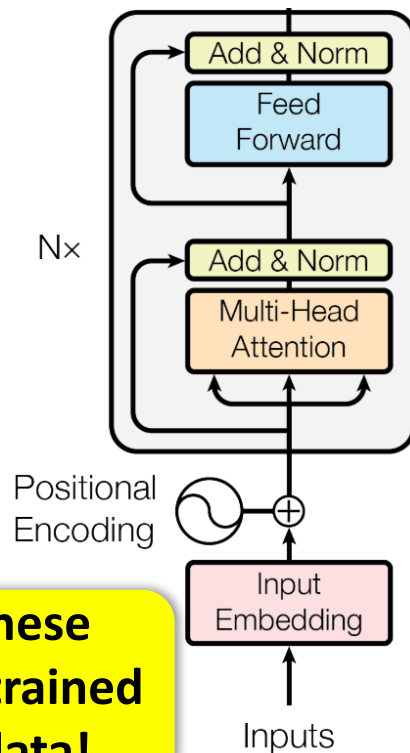
Embedding type models: GloVe, word2vec, topicETM

Recurrent networks: skip-thoughts, LSTM lm_1b

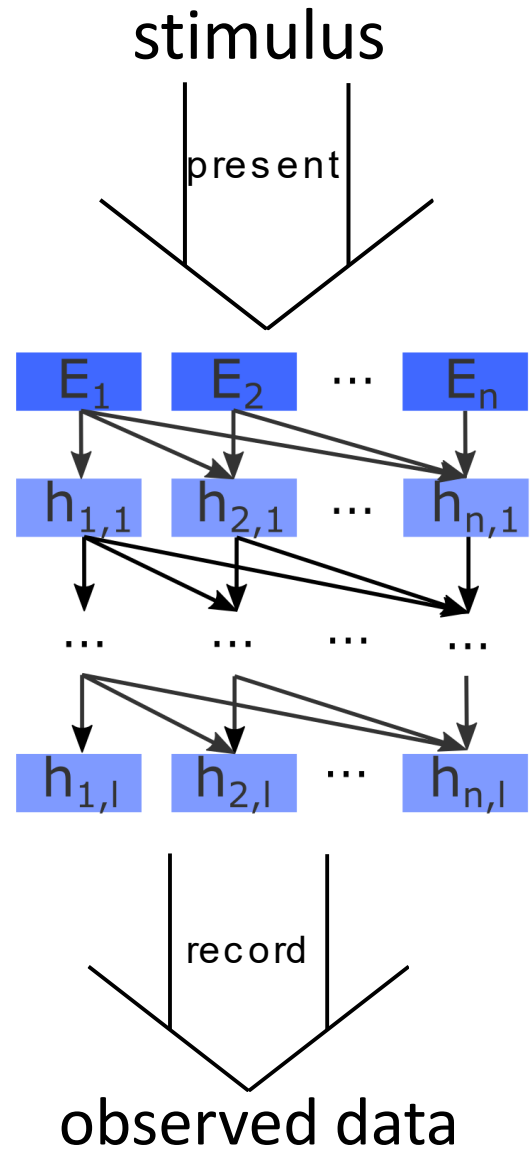
Transformers

BERTs
RoBERTas
XLMs
Transformer-XLs
XLNets
CTRL
T5s
ALBERTs
GPTs

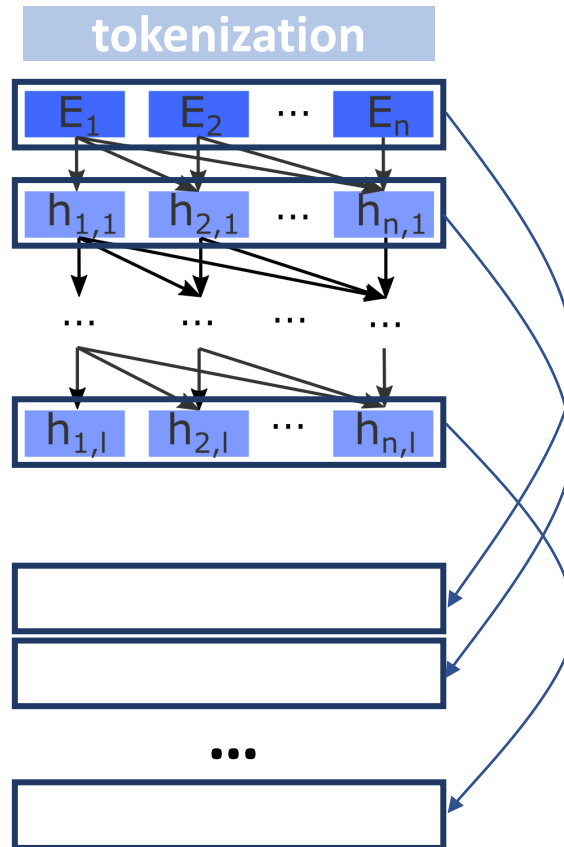
**None of these
models are trained
on brain data!**



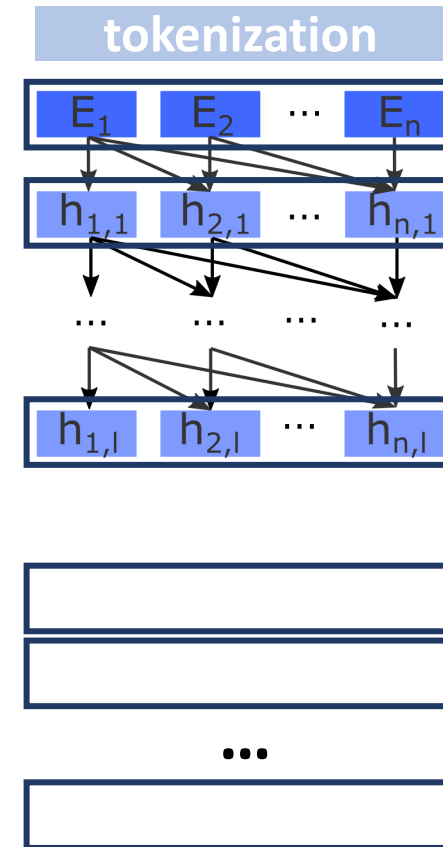
Treating models as experimental subjects



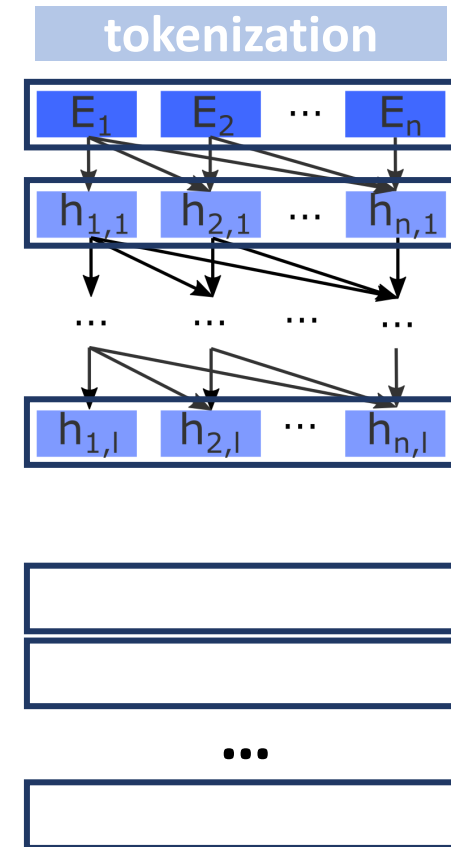
Beekeeping encourages the conservation of local habitats.



It is in every beekeeper's interest to conserve local plants that produce pollen.



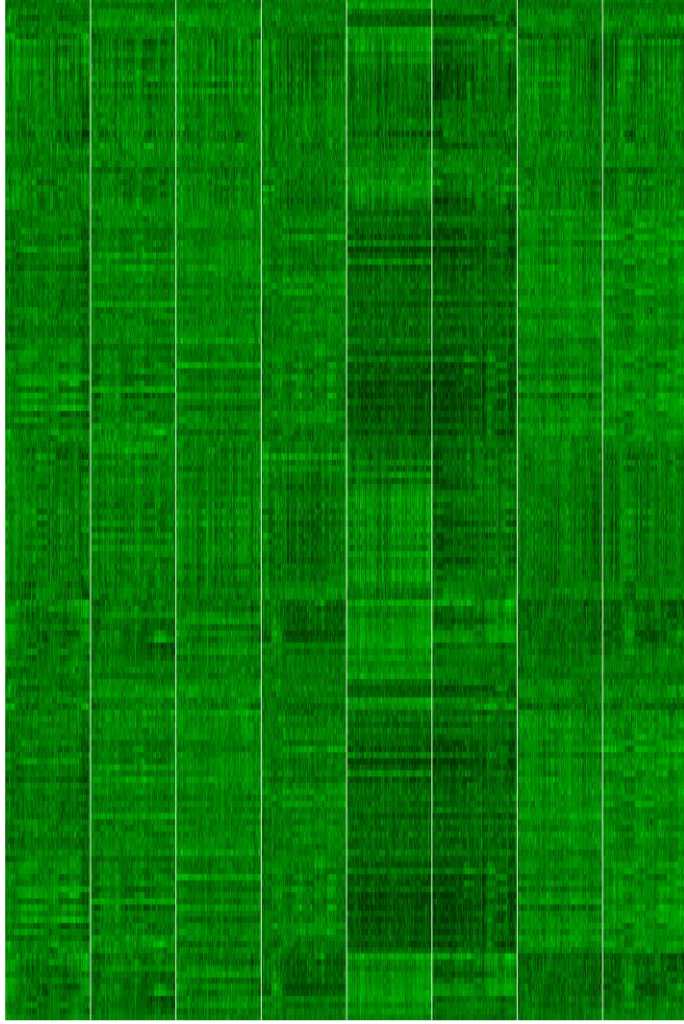
...



Neural benchmarks

← sentences →

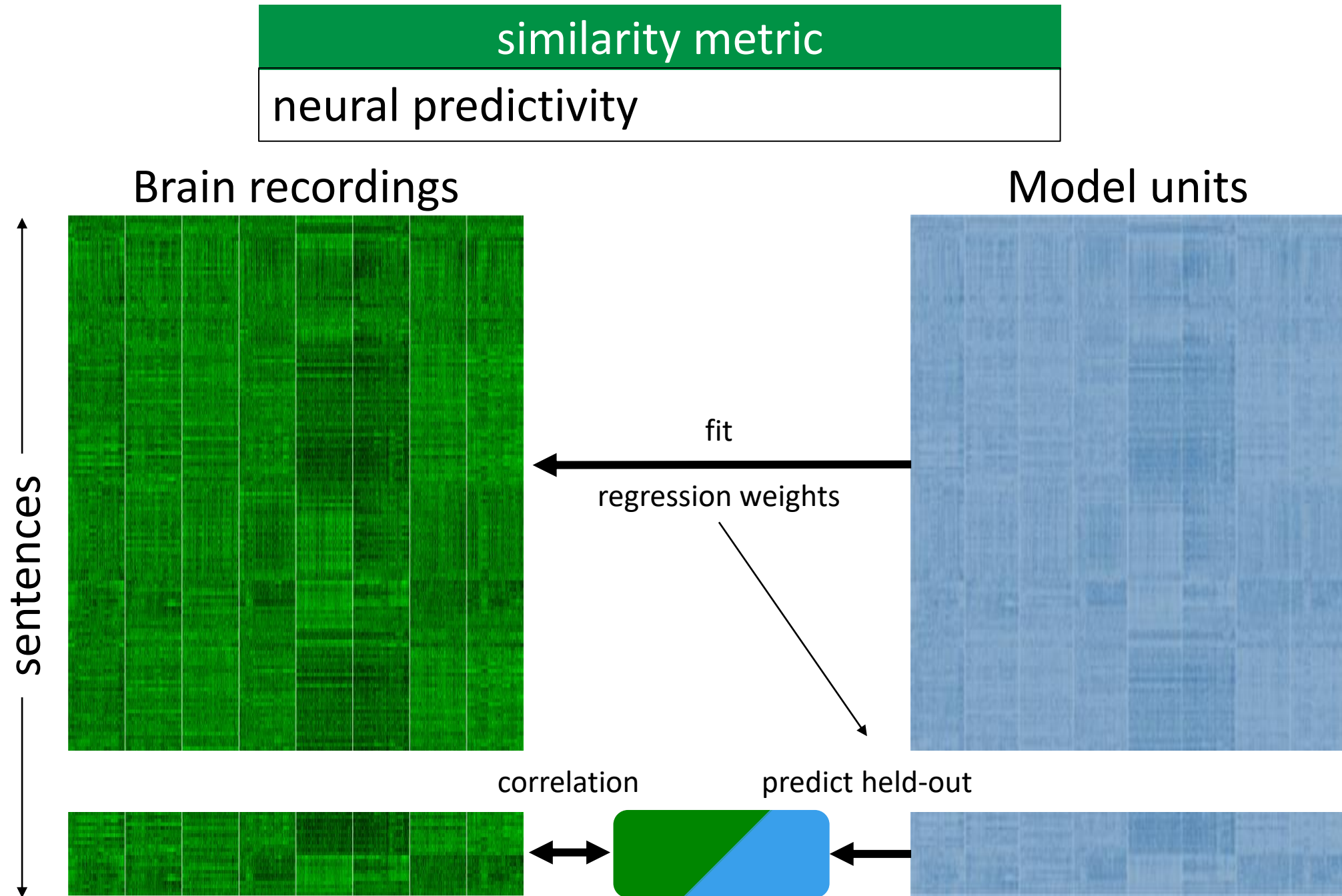
Brain recordings



Model units



Neural benchmarks



Stimuli

Pereira2018

"Beekeeping encourages the conservation of local habitats. It is in every beekeeper's interest..."

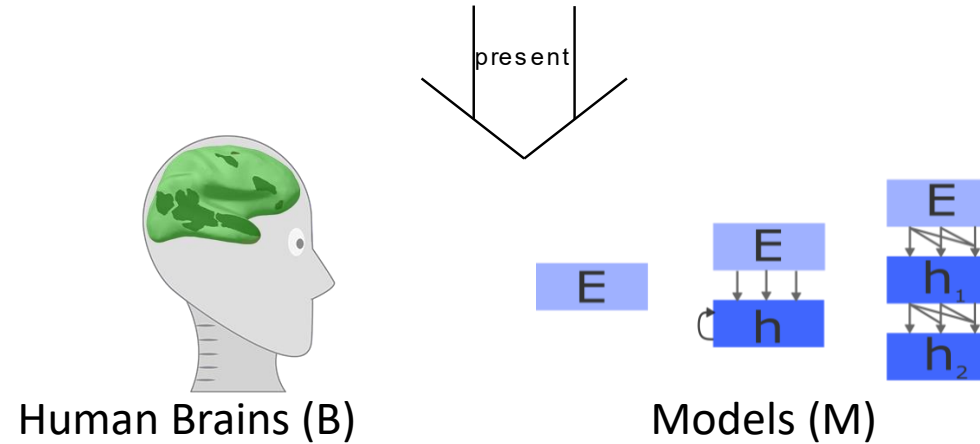
Fedorenko2016

“Alex was tired so he took a nap.”

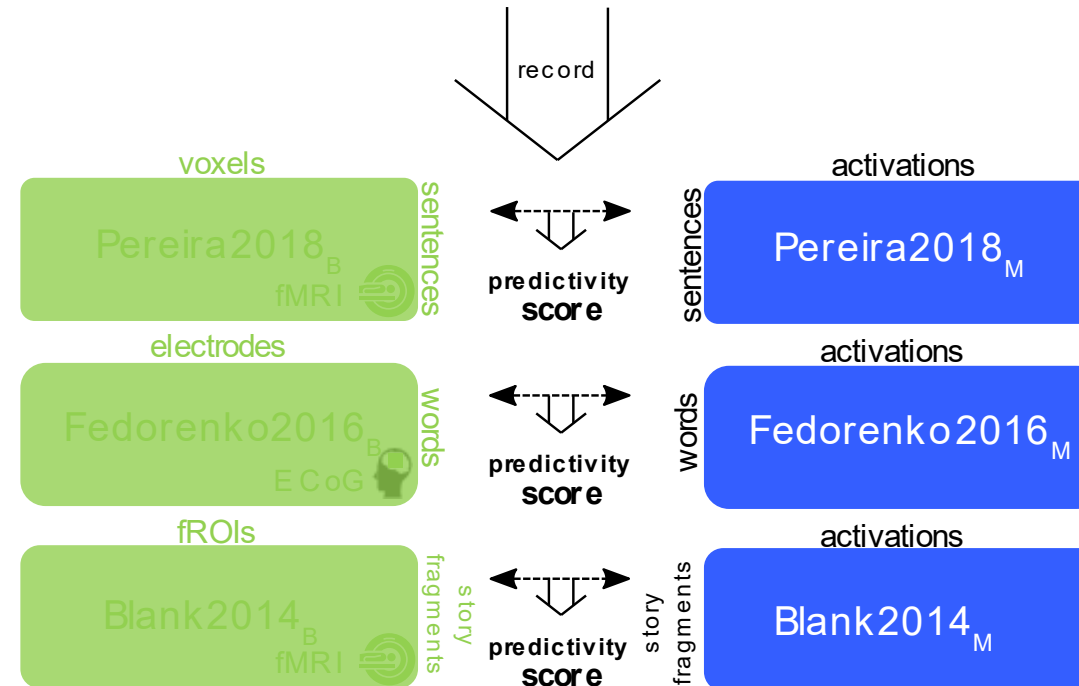
Blank2014

"If you were to journey to the North of England, you would come to a valley that is surrounded by moors as high as mountains. It is in this valley where you..."

Experimental Participants



Comparative Measurements

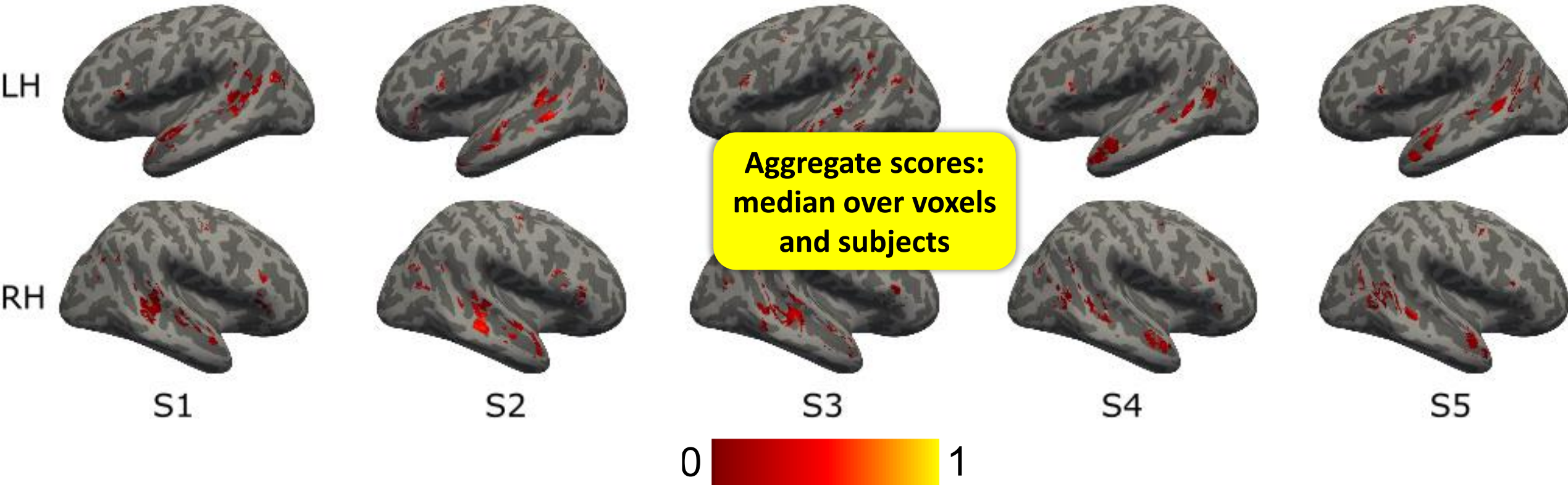


**We want one model
to predict *all* data**

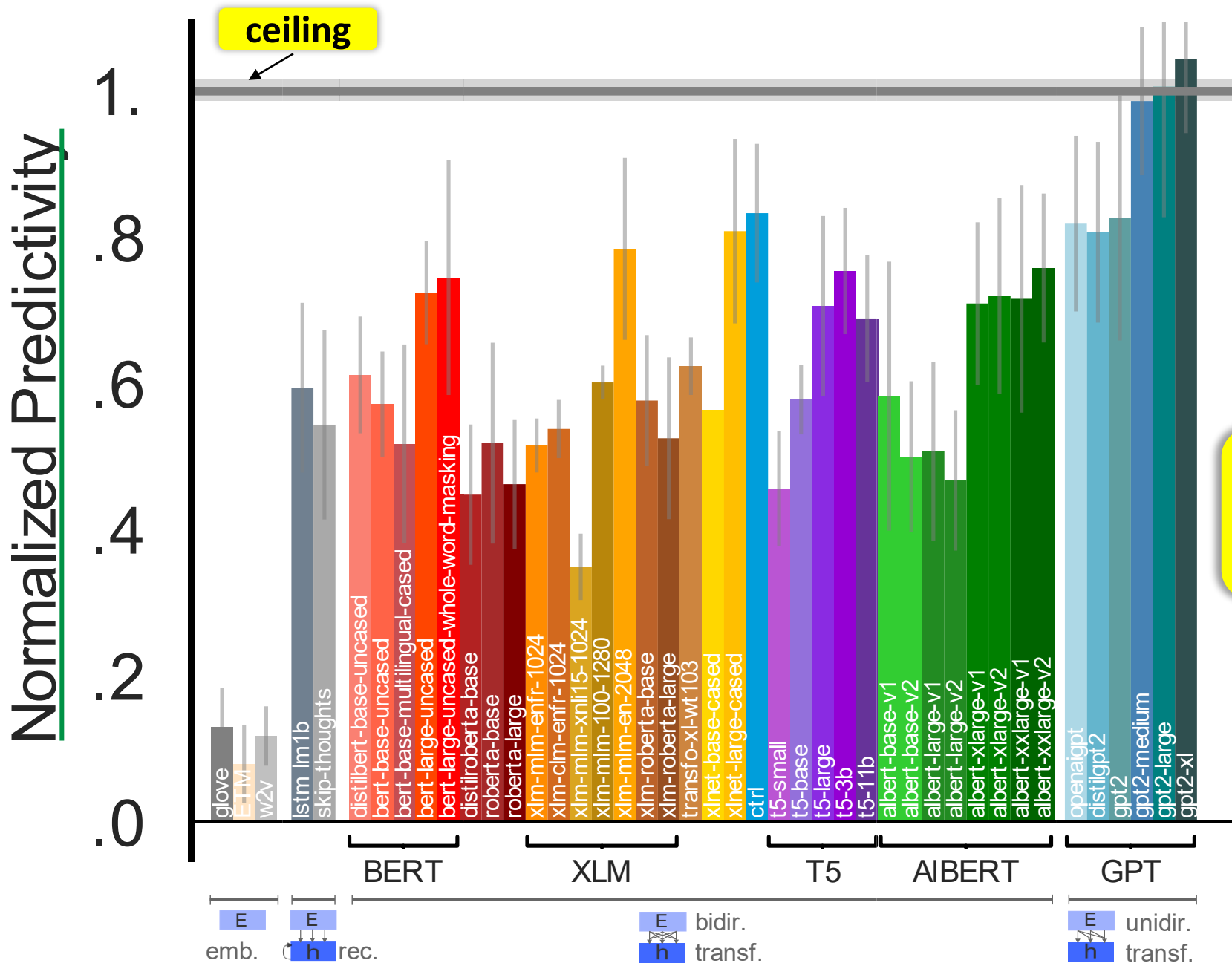
GloVe voxel-wise predictivity scores



← GloVe



Certain language models predict human language recordings



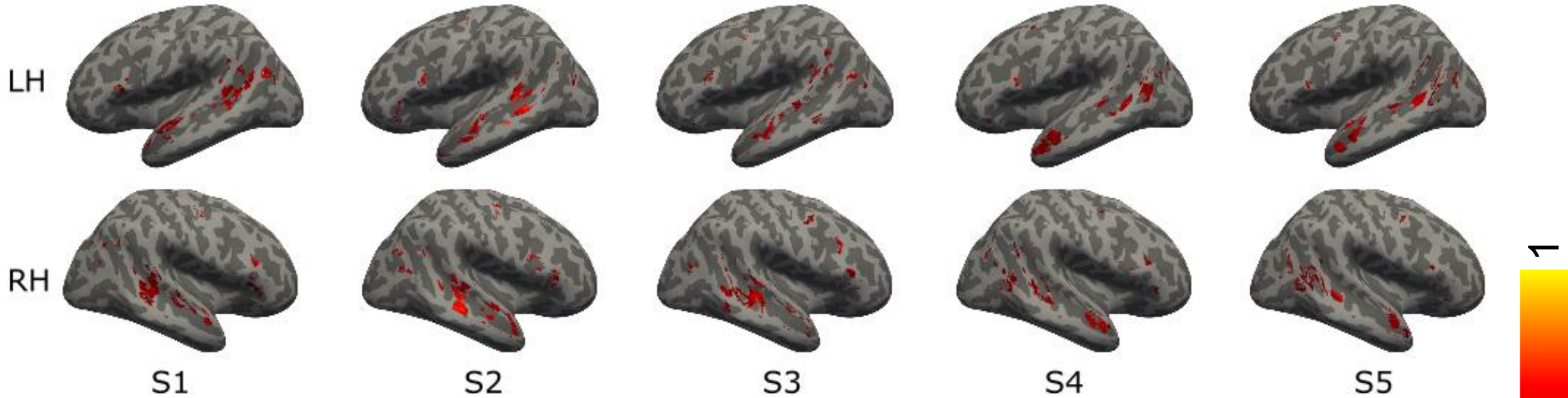
gpt2-xl hits our estimated ceiling for this benchmark

Small differences can lead to very different brain predictivities, warranting a full survey

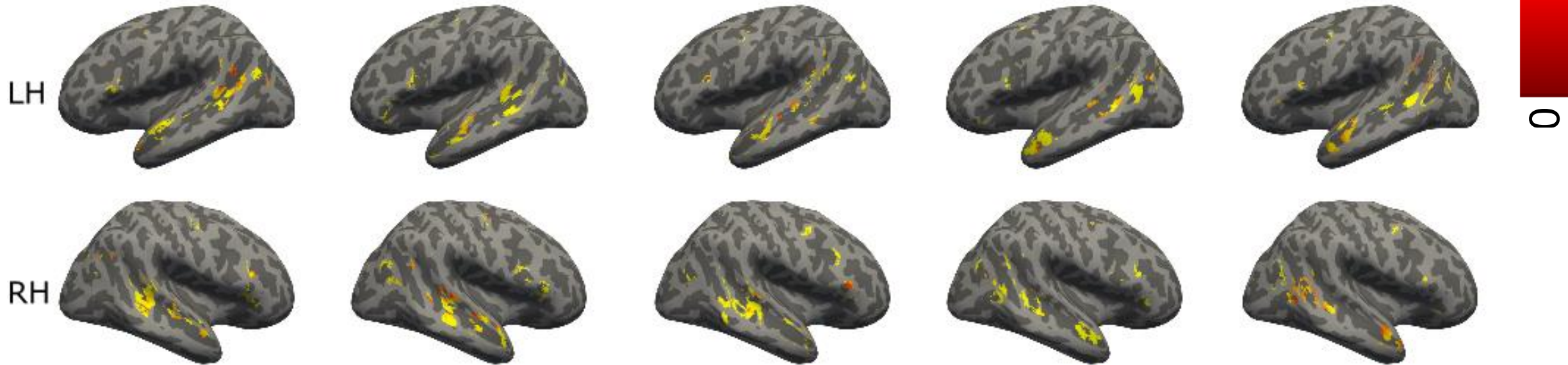
Jain & Huth 2018
Gauthier & Ivanova 2018
Jat et al. 2019
Toneva & Wehbe 2019
Gauthier & Levy 2020
Wang et al. 2020

GPT2-xl accurately predicts a large portion of voxels

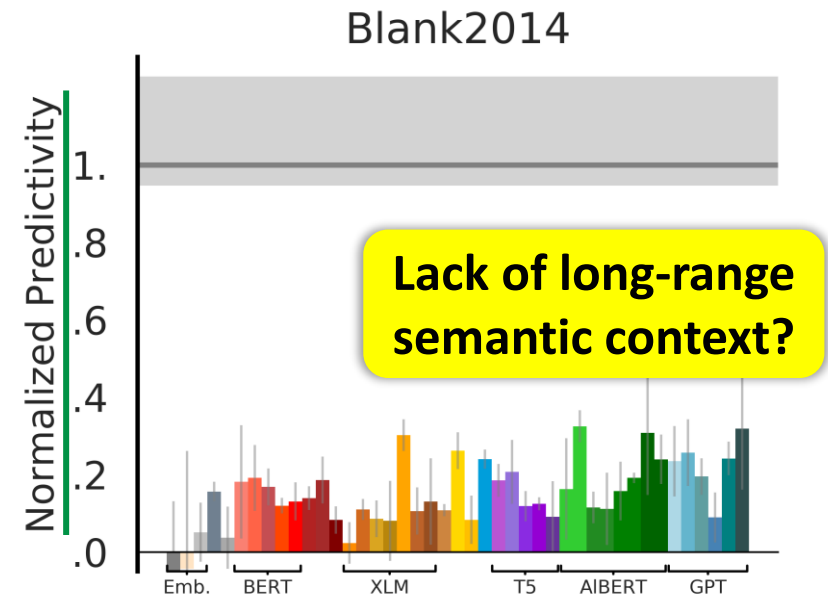
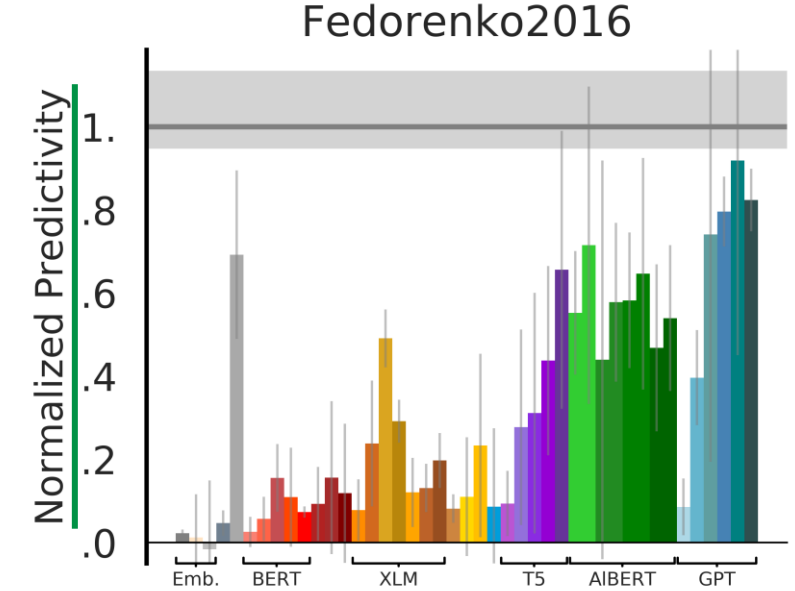
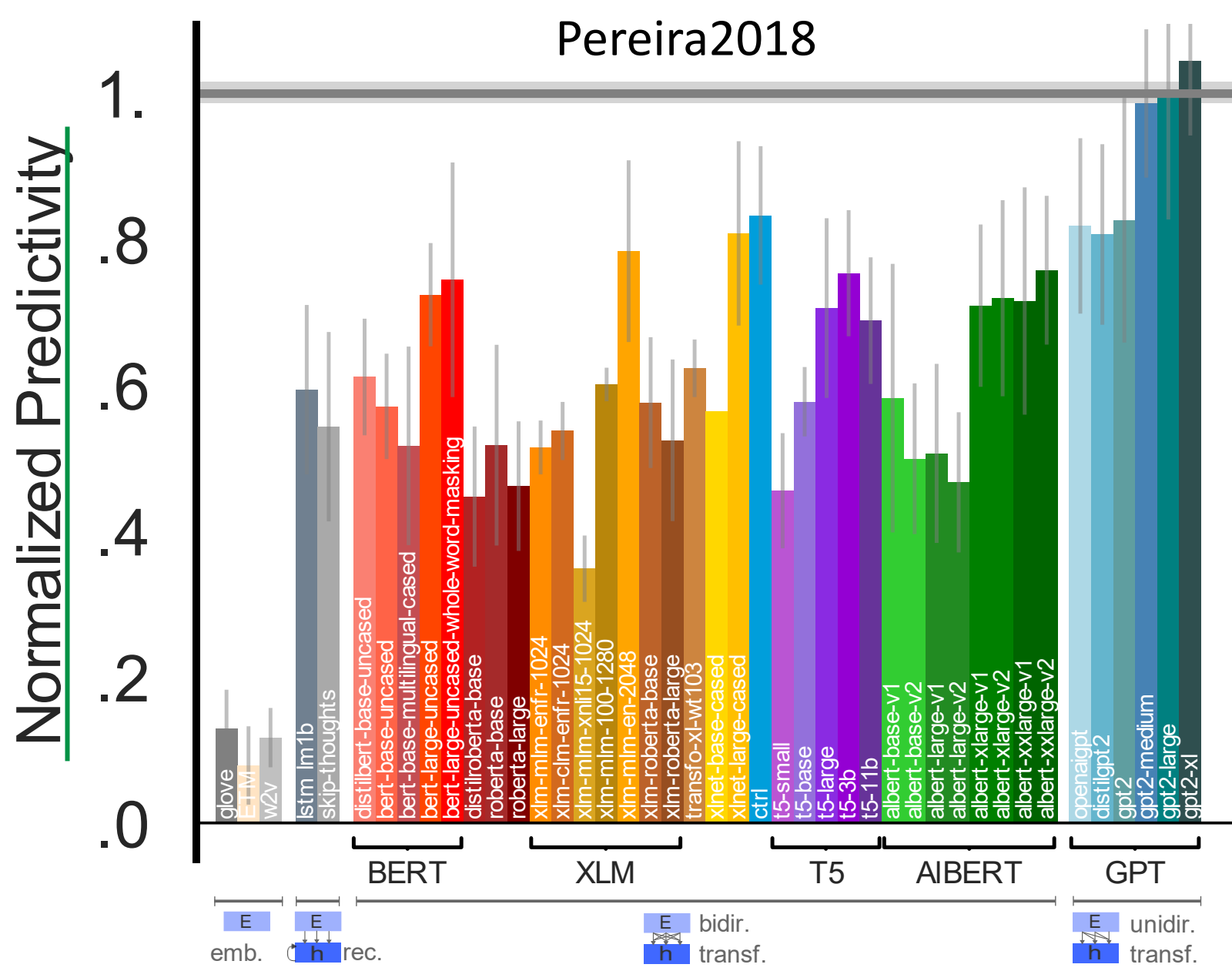
Glove



GPT2-xl

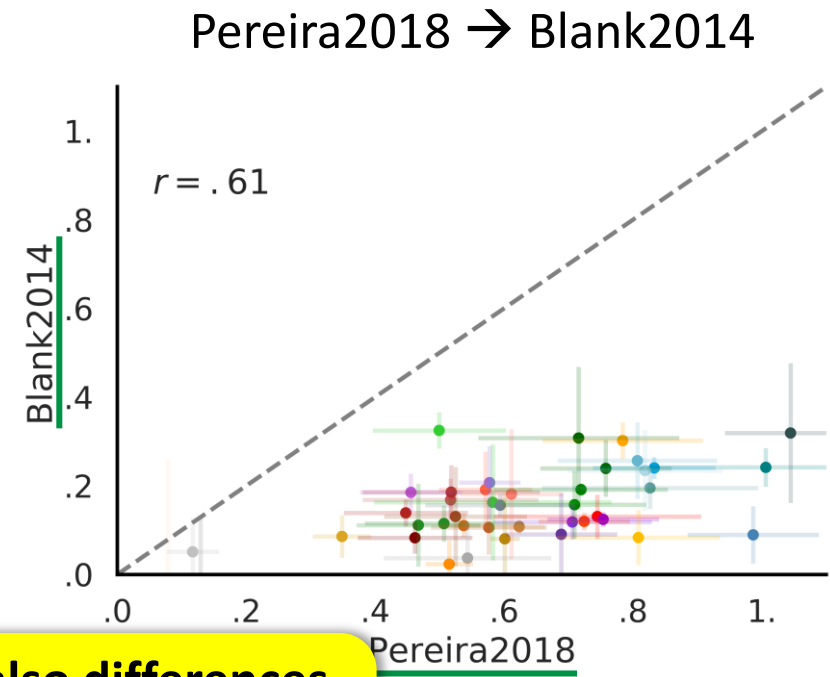
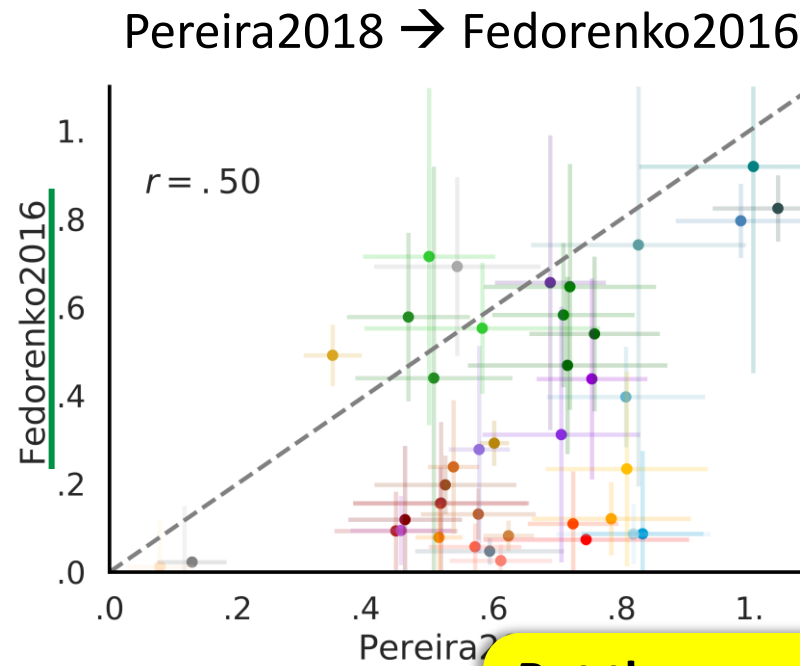
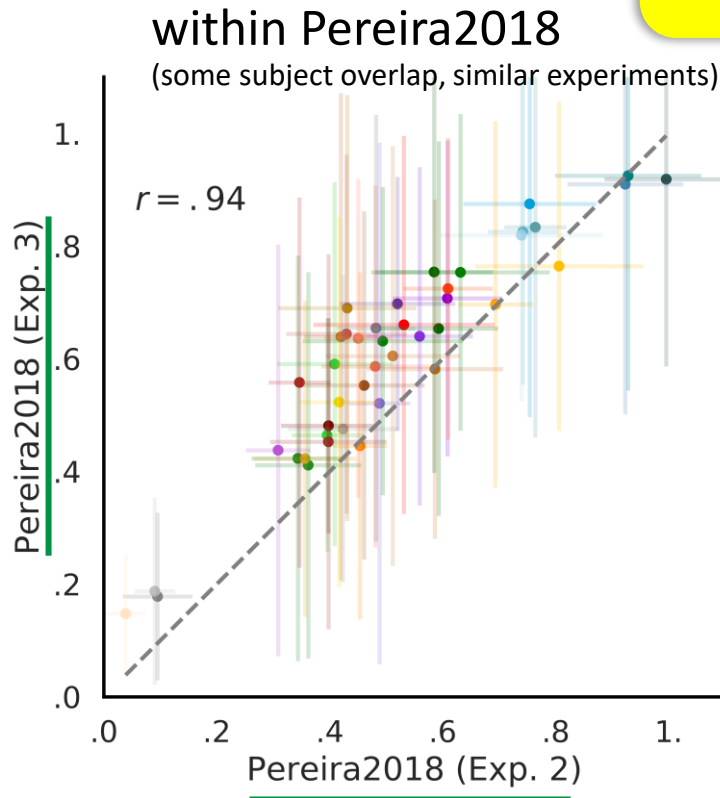


Language Models predict human language recordings



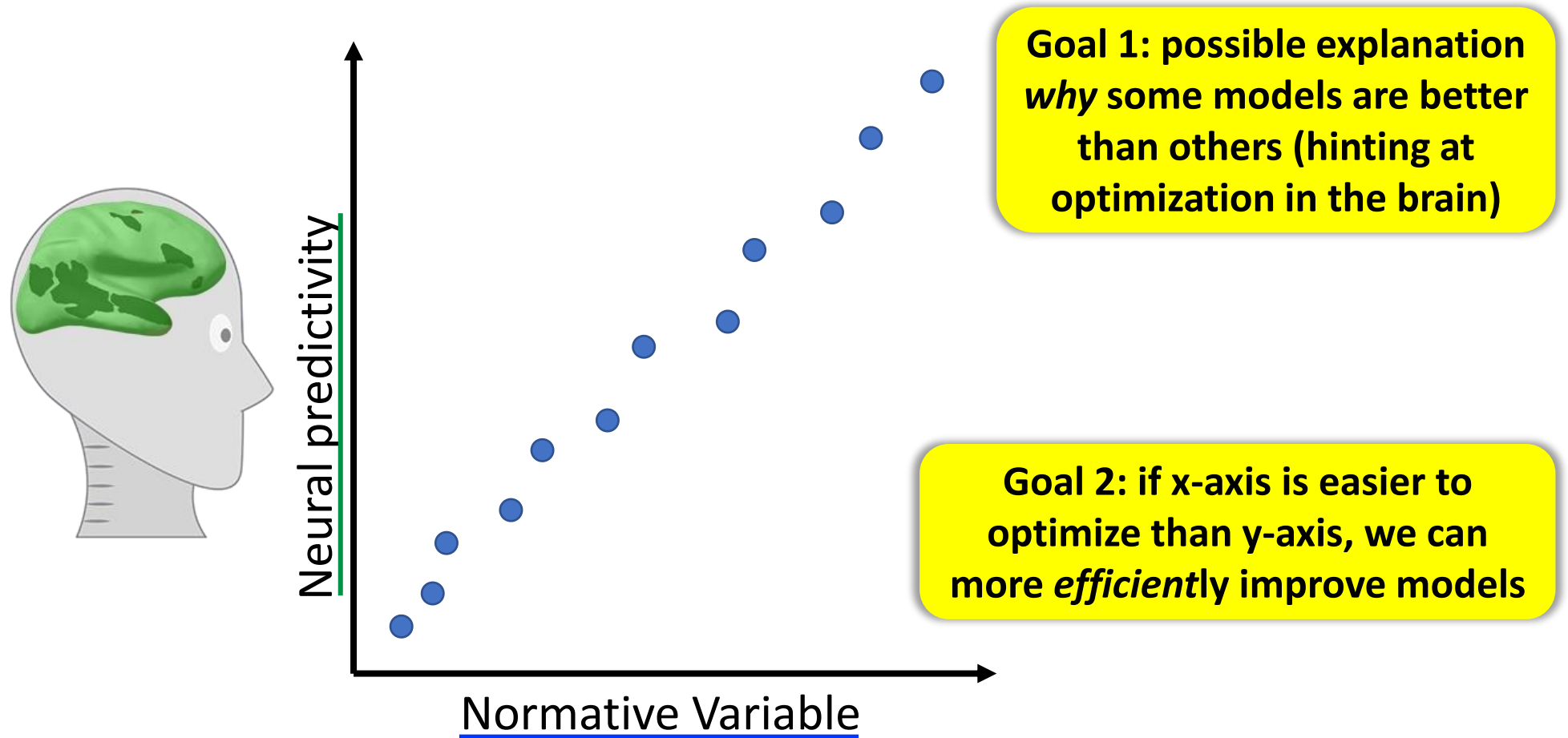
Control: model scores across benchmarks are correlated, although differences exist

Scores generalize to a good extent



But there are also differences, making each individual benchmark valuable

What explains the model differences?



Next-Word Prediction on WikiText-2

= Gold dollar =

The gold dollar or gold one @-@ dollar piece was a coin struck as a regular issue by the United States Bureau of the Mint from 1849 to 1889 . The coin had three types over its lifetime , all designed by Mint Chief Engraver James B. Longacre . The Type 1 issue had

...

	WikiText-2		
	Train	Valid	Test
Articles	600	60	60
Tokens	2,088,628	217,646	245,569
Vocab	33,278		
OoV	2.6%		

Merity et al. 2016

Alaska

Alaska is

Alaska is about

Alaska is about twelve

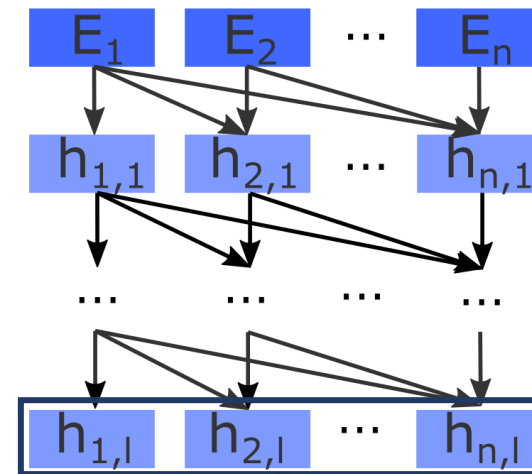
Alaska is about twelve times

Alaska is about twelve times larger

Alaska is about twelve times larger than

Alaska is about twelve times larger than New

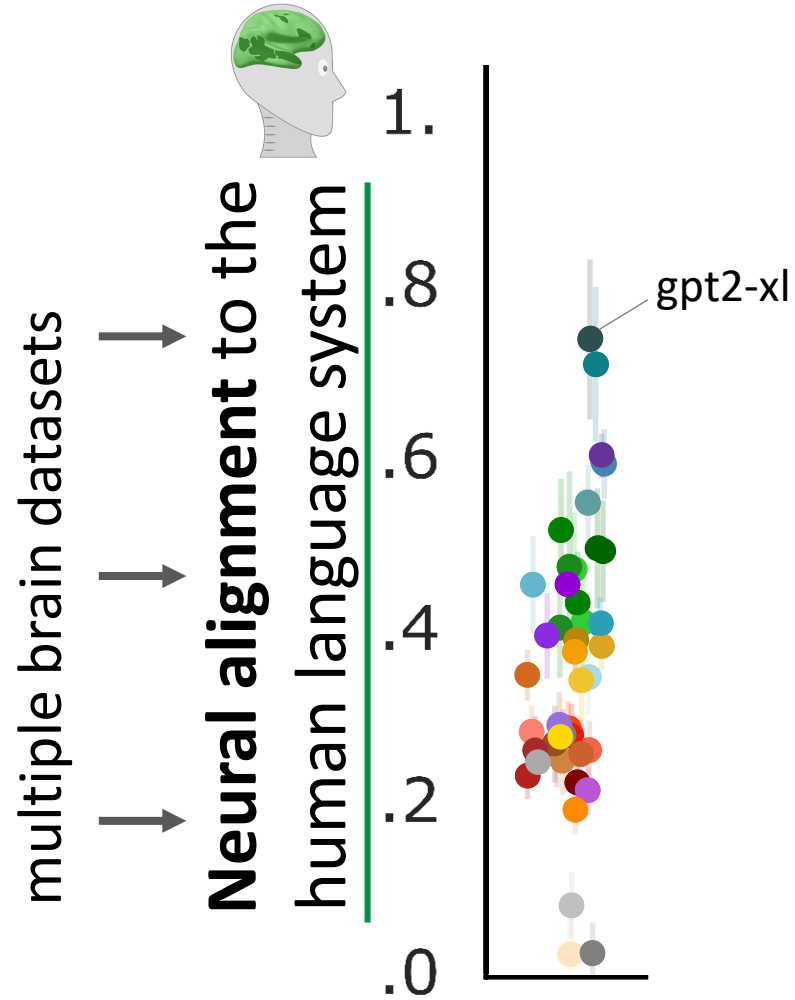
Alaska is about twelve times larger than New York



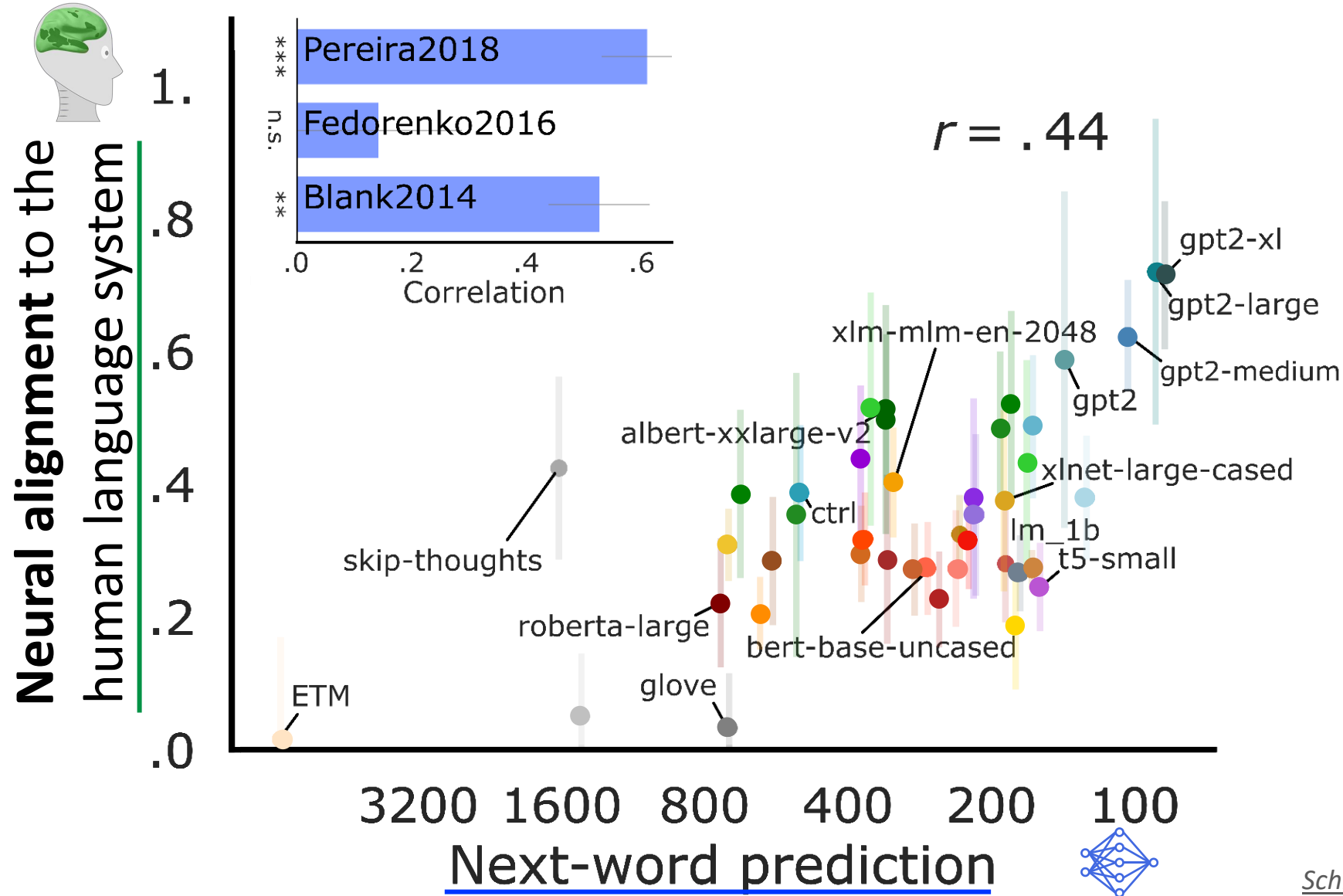
Surprisal of seeing
actual next word:
perplexity =
 $\exp(\text{NLL Loss})$

afternoon | alaska | animation | article | ...

The better models can predict the next word, the more brain-like they are



The better models can predict the next word, the more brain-like they are



What about other language tasks?



9 “General Language Understanding Evaluation” tasks:

Sentence grammaticality (CoLa)

Sentence sentiment (SST-2)

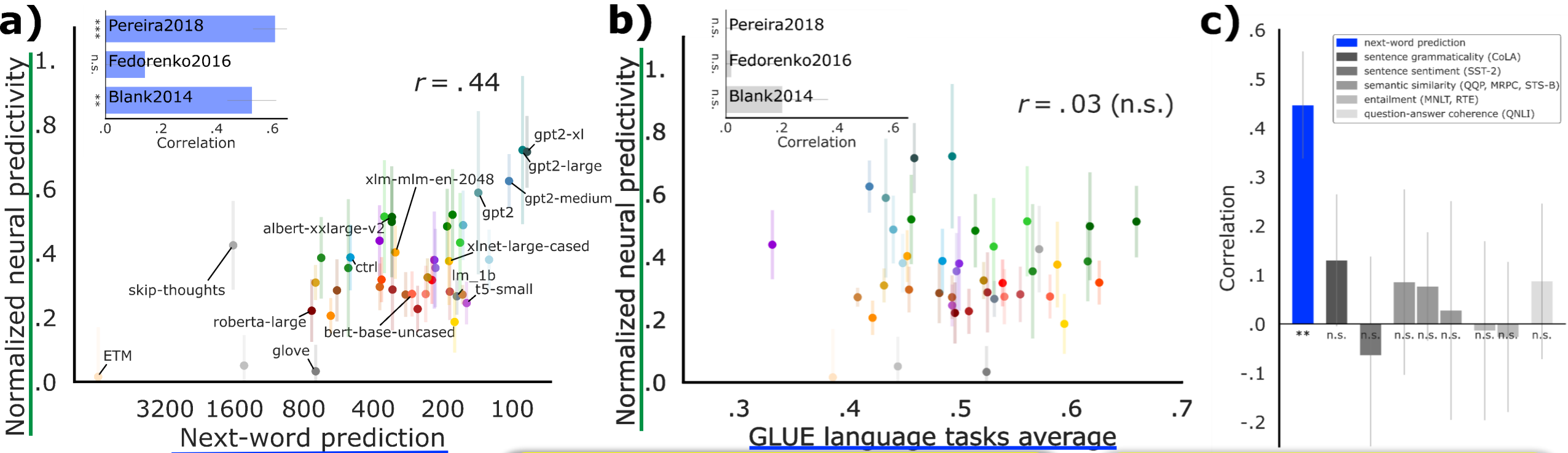
Semantic similarity (QQP, MRPC, STS-B)

Entailment (MNLT, RTE)

Question-answer coherence (QNLI)

Winograd (WNLI; ignored due to known issues)

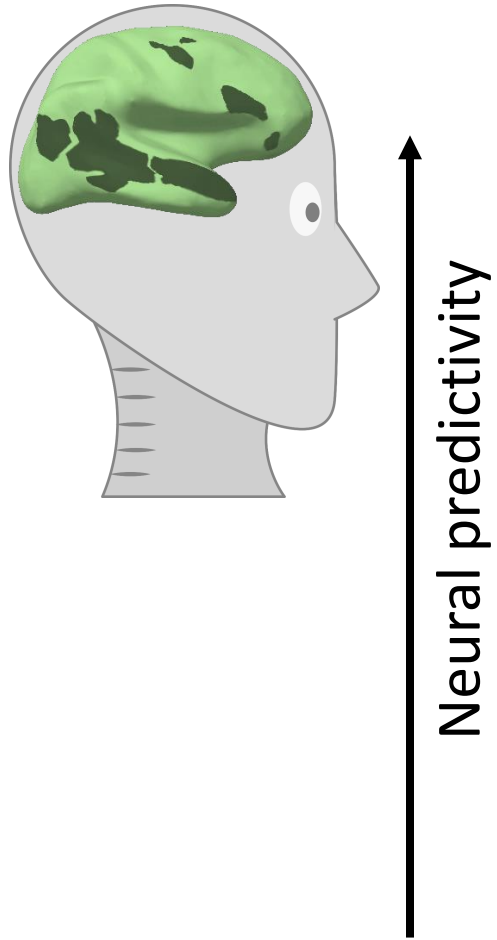
Next-Word Prediction performance **selectively** correlates with neural predictivity



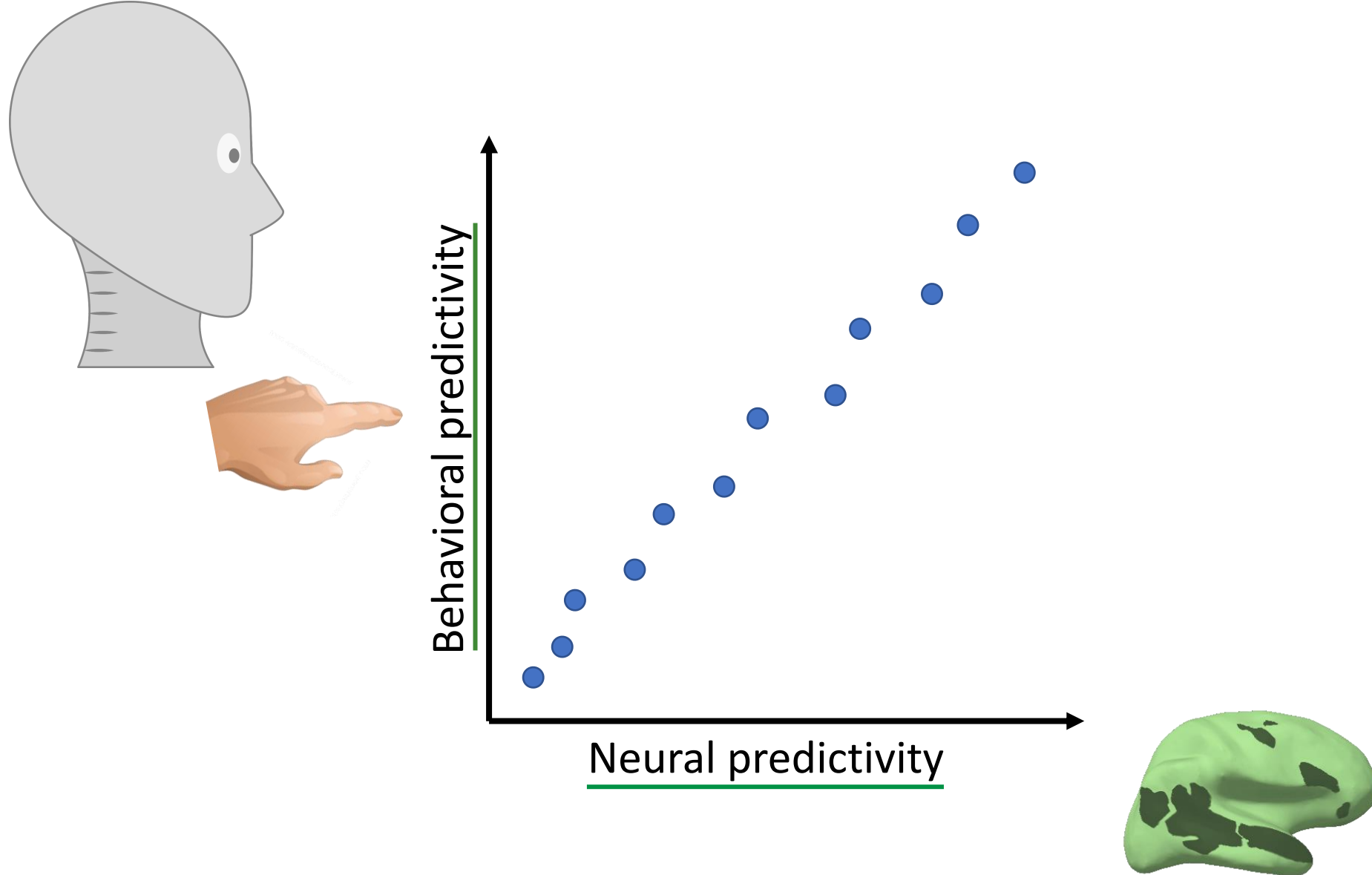
Online prediction may fundamentally shape language processing in the brain

Ongoing work: for high-performing models, reasoning capabilities (MMLU/BBH) also drive brain alignment (*Aw et al. in prep*)

Is any of this behaviorally relevant?



Is any of this behaviorally relevant?



Behavioral target: human reading times

Futrell et al. 2018

10256 words x 179 subjects

*If | you | were | to | journey | to | the |
North | of | England, | you | would | come
| to | a | valley | that | is | surrounded |
by | moors | as | high | as | mountains. |
It | is | in | this | valley | where | you |
would | find | the | city | of | Bradford, |
where | once | a | thousand | spinning | ...*

Treat reading times as representation target

The Natural Stories Corpus

Richard Futrell¹, Edward Gibson¹, Harry J. Tily², Idan Blank¹,
Anastasia Vishnevetsky¹, Steven T. Piantadosi³, and Evelina Fedorenko^{4,5}

¹MIT Department of Brain and Cognitive Sciences ²Netflix, Inc.

³University of Rochester Department of Brain and Cognitive Sciences

⁴Massachusetts General Hospital Department of Psychiatry

⁵Harvard Medical School Department of Psychiatry

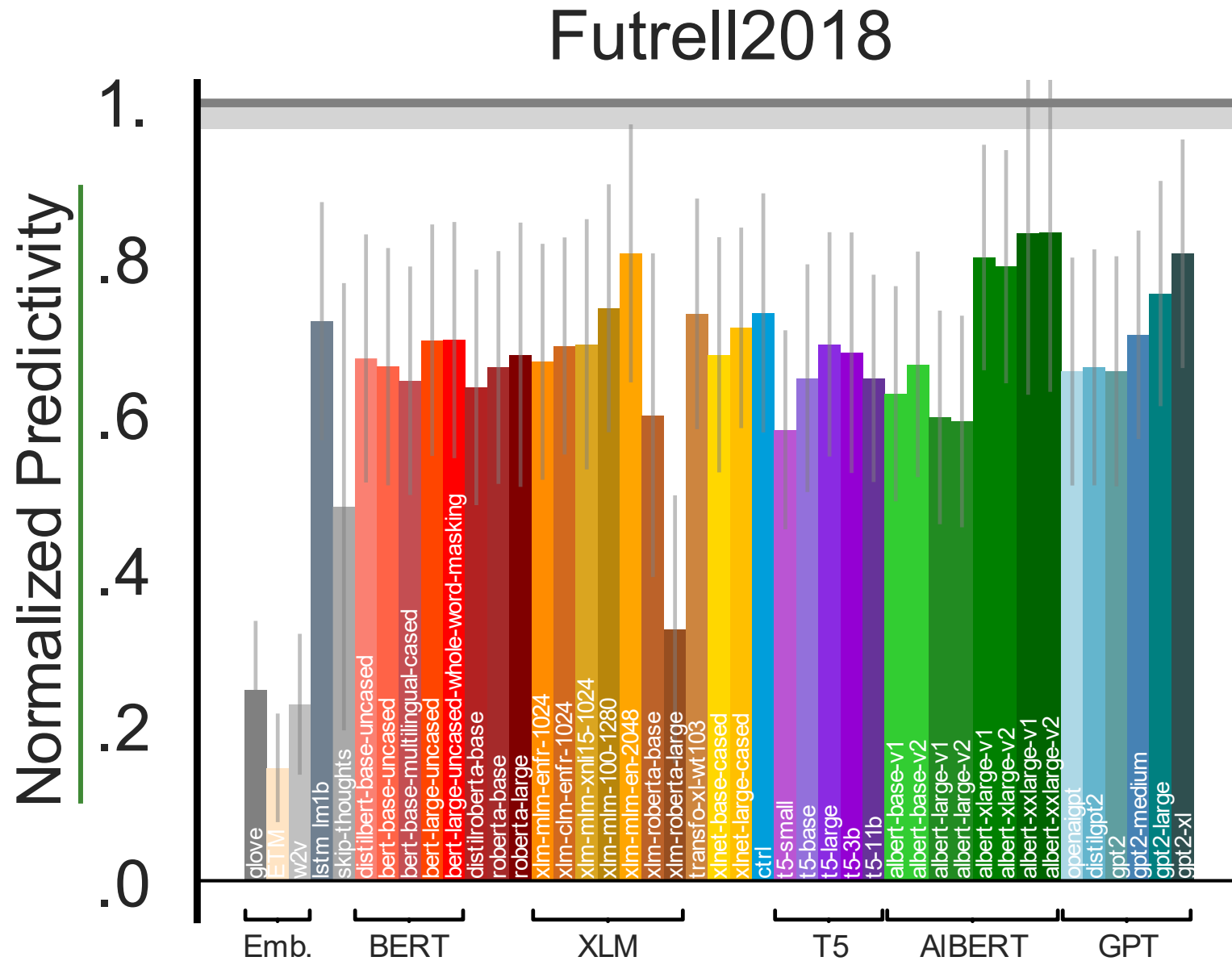
{futrell, egibson, iblack, evelina9}@mit.edu,
hal.tily@gmail.com, staseyvi@mail.med.upenn.edu

Abstract

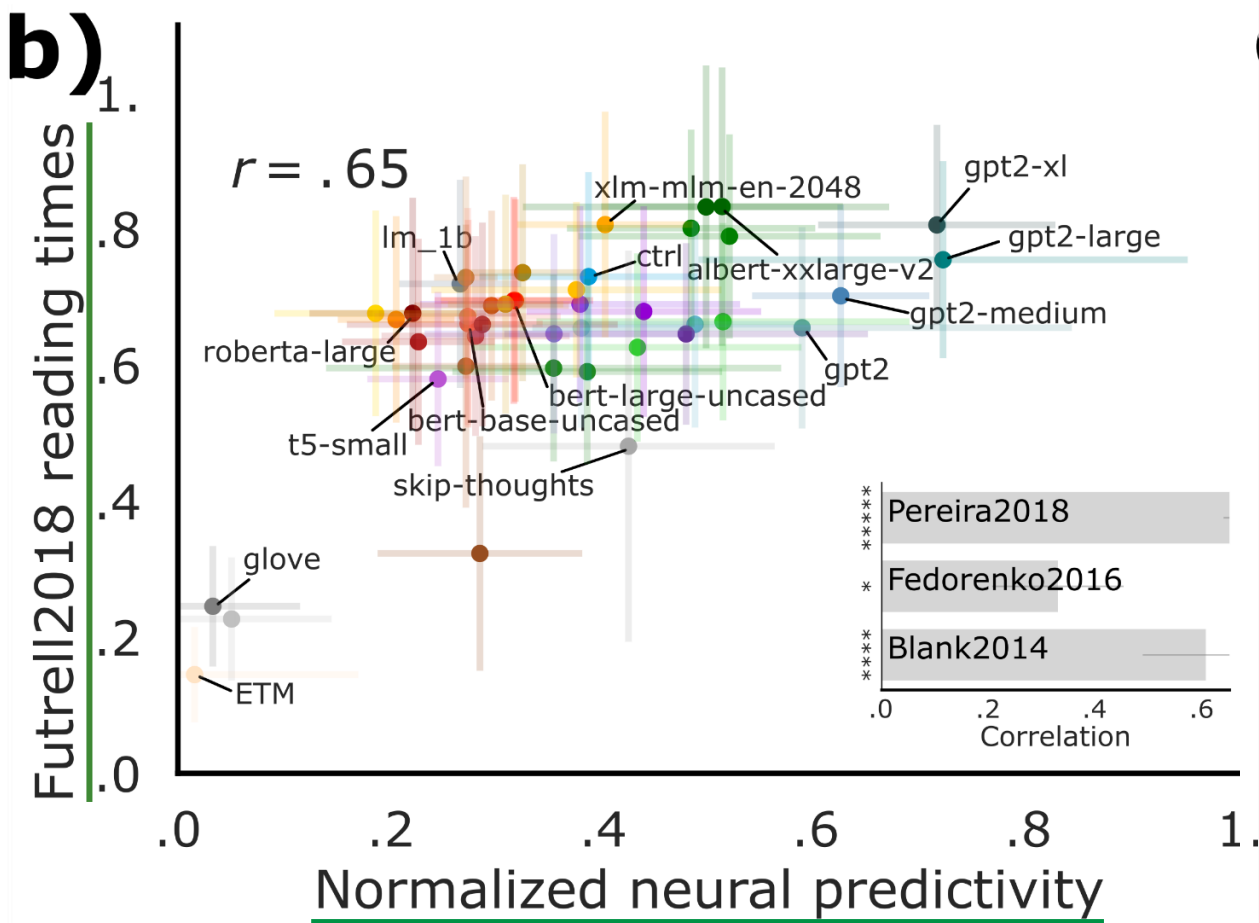
It is now a common practice to compare models of human language processing by comparing how well they predict behavioral and neural measures of processing difficulty, such as reading times, on corpora of rich naturalistic linguistic materials. However, many of these corpora, which are based on naturally-occurring text, do not contain many of the low-frequency syntactic constructions that are often required to distinguish between processing theories. Here we describe a new corpus consisting of English texts edited to contain many low-frequency syntactic constructions while still sounding fluent to native speakers. The corpus is annotated with hand-corrected Penn Treebank-style parse trees and includes self-paced reading time data and aligned audio recordings. Here we give an overview of the content of the corpus and release the data.

Keywords: Cognitive modeling, reading time, psycholinguistics

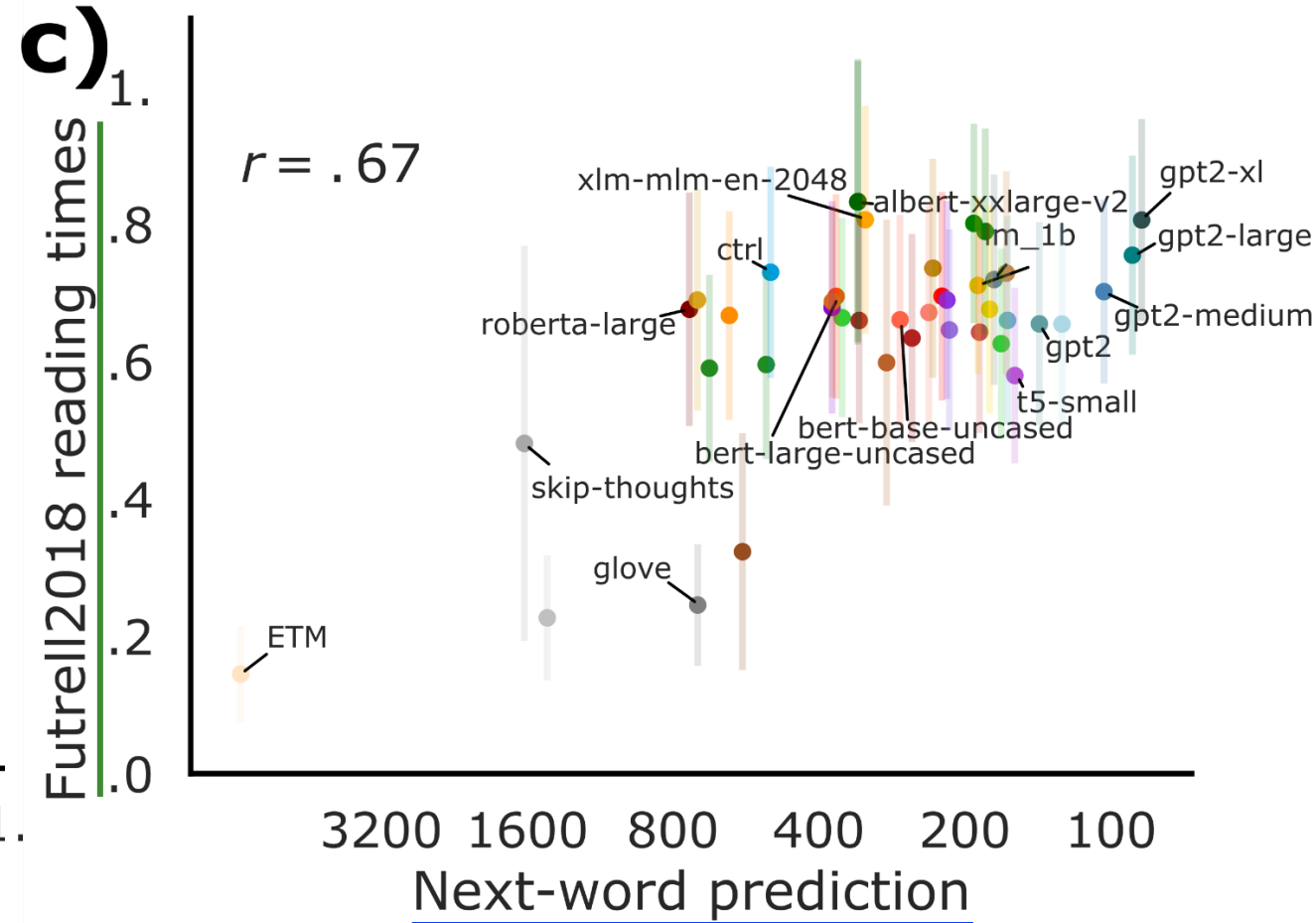
Behavioral scores

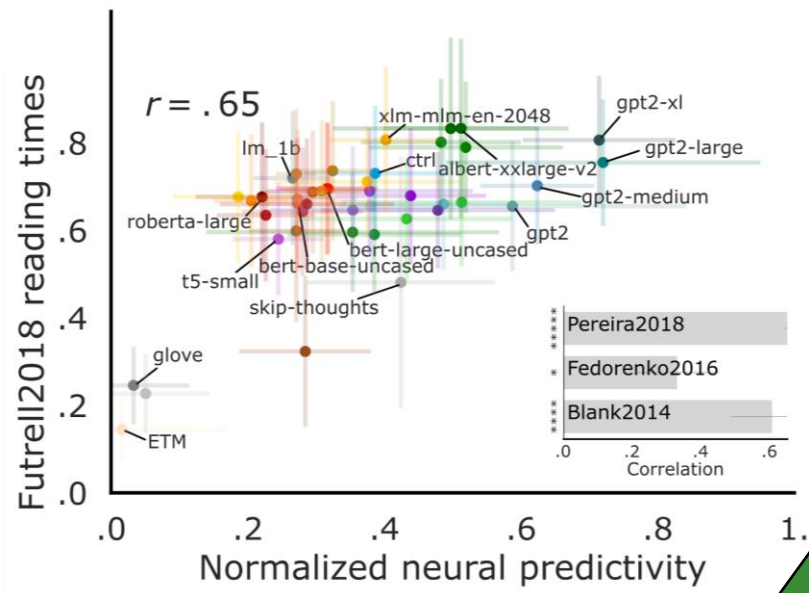


Neural scores correlate with Behavioral scores

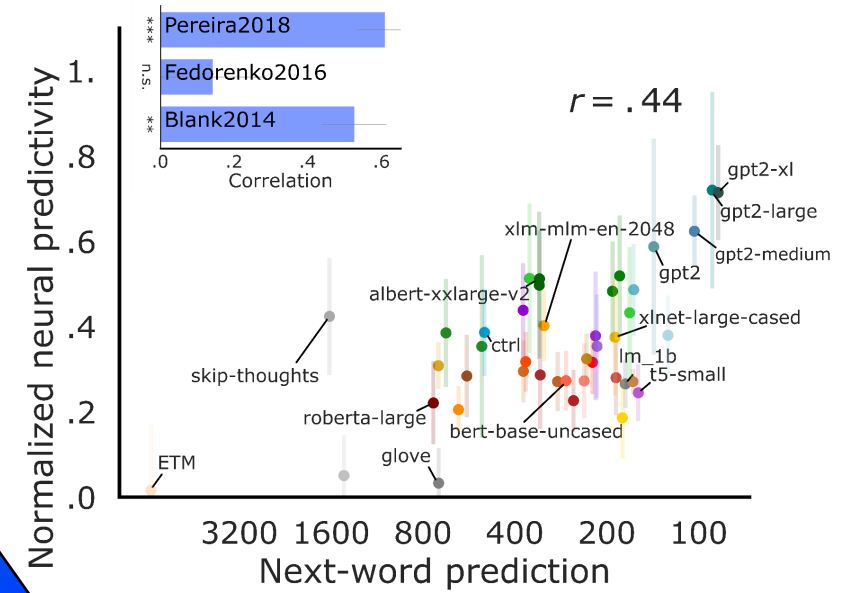


Task scores correlate with Behavioral scores





Neural

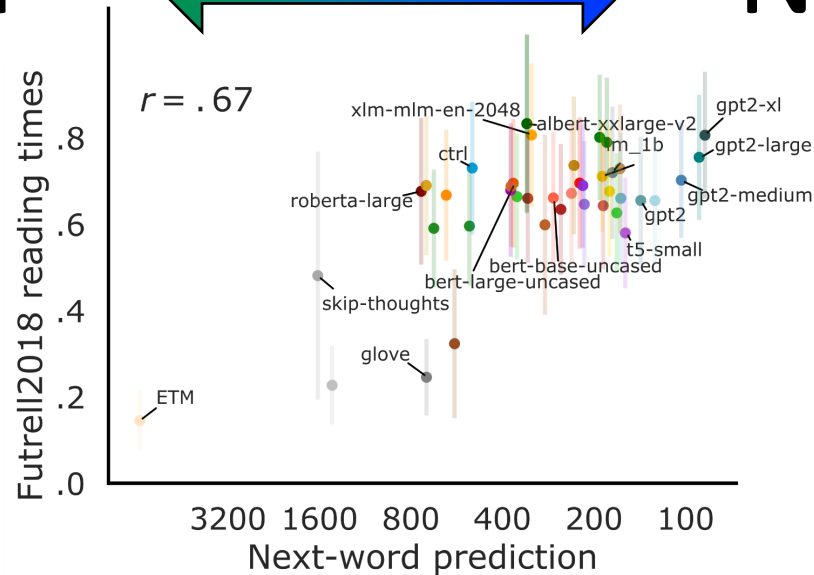


Integrative Modeling:
link neural mechanisms,
behavior, and computation

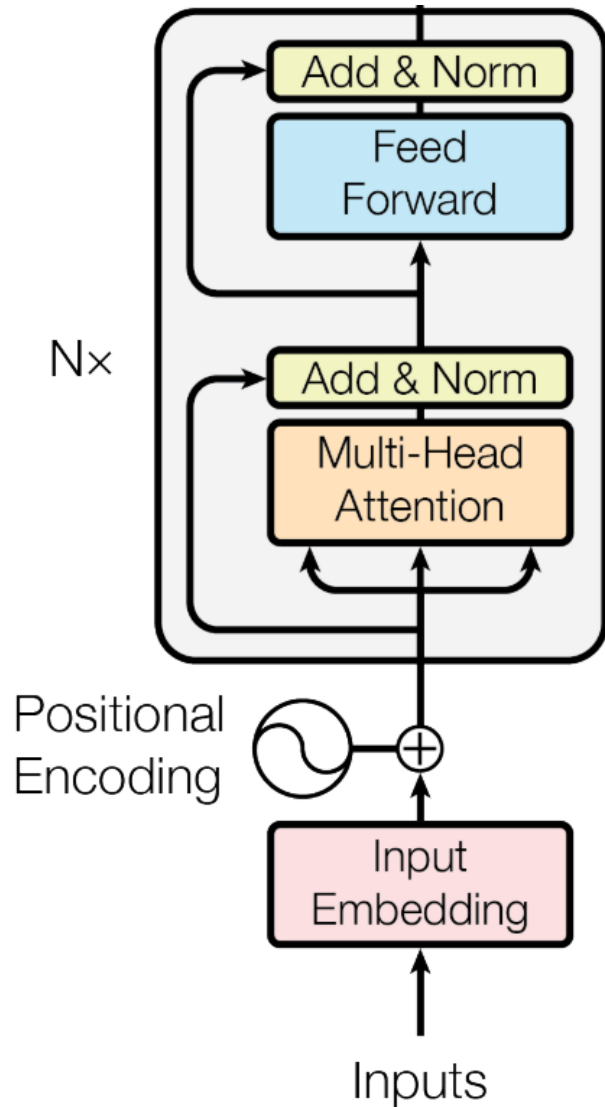
Schrimpf et al. Neuron 2020

Behavioral

Normative Task



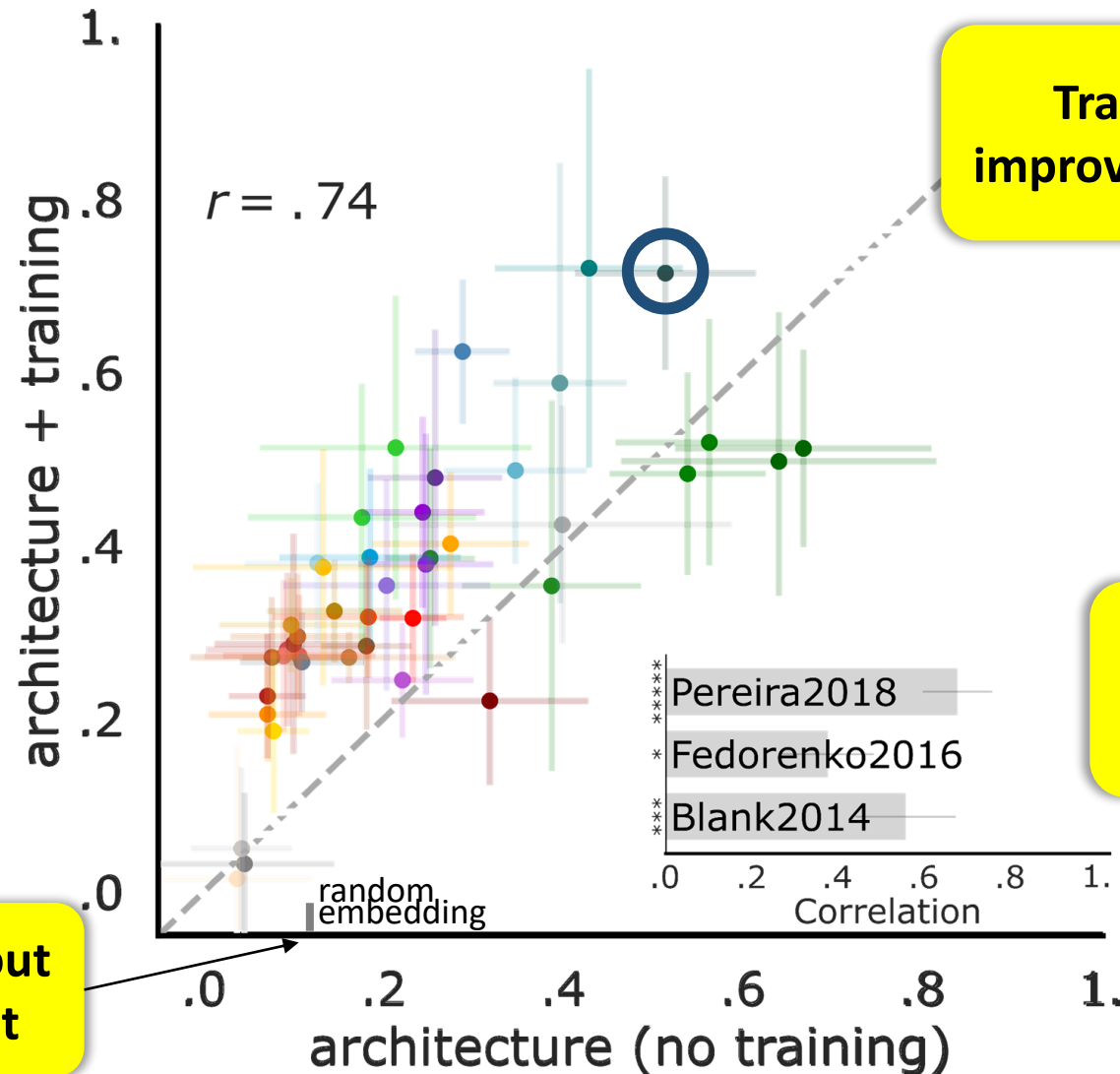
What is the relative importance of evolutionary and learning-based optimization?



Evolution \simeq community optimization over architectural properties

Experience-dependent learning \simeq updating of weights over training

Architecture substantially contributes to models' brain predictivity

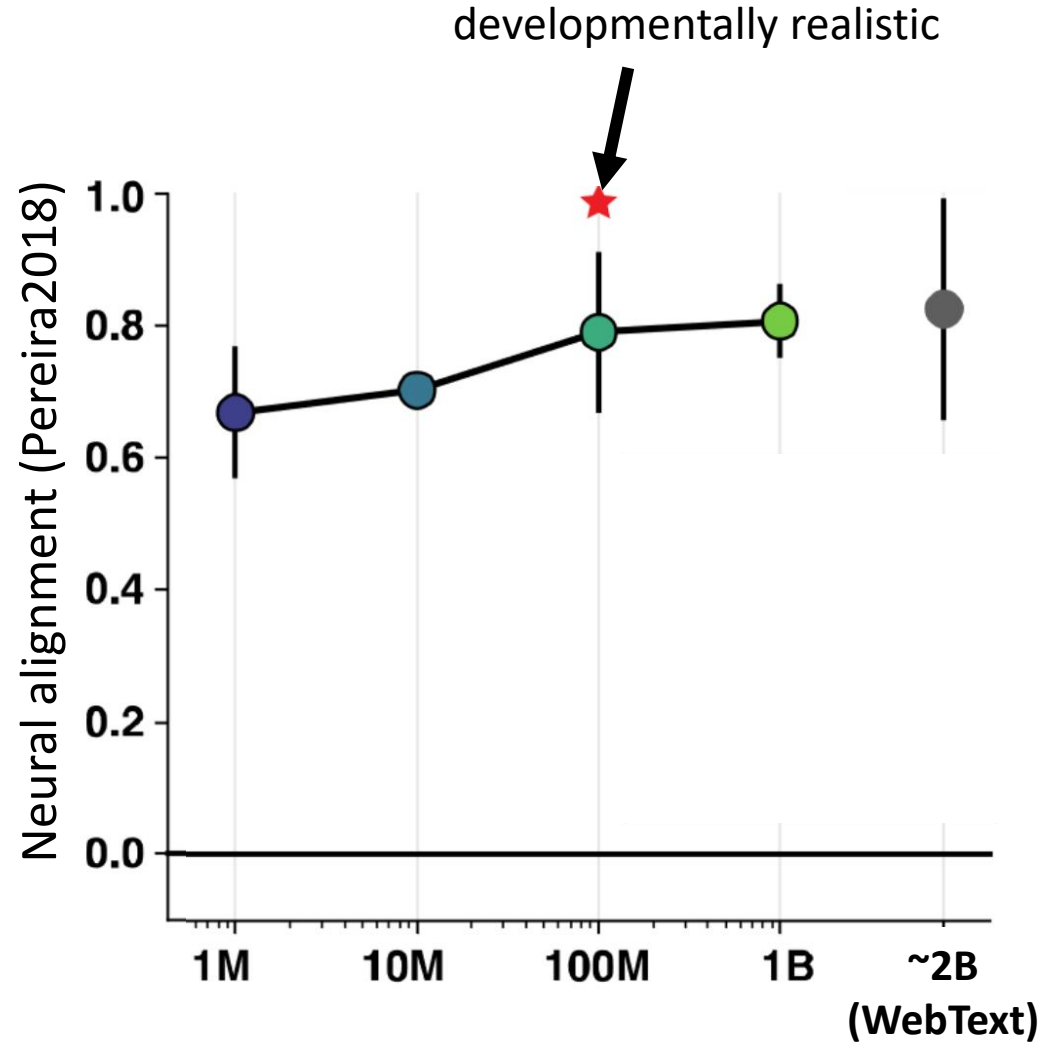


Training generally improves scores by ~53%

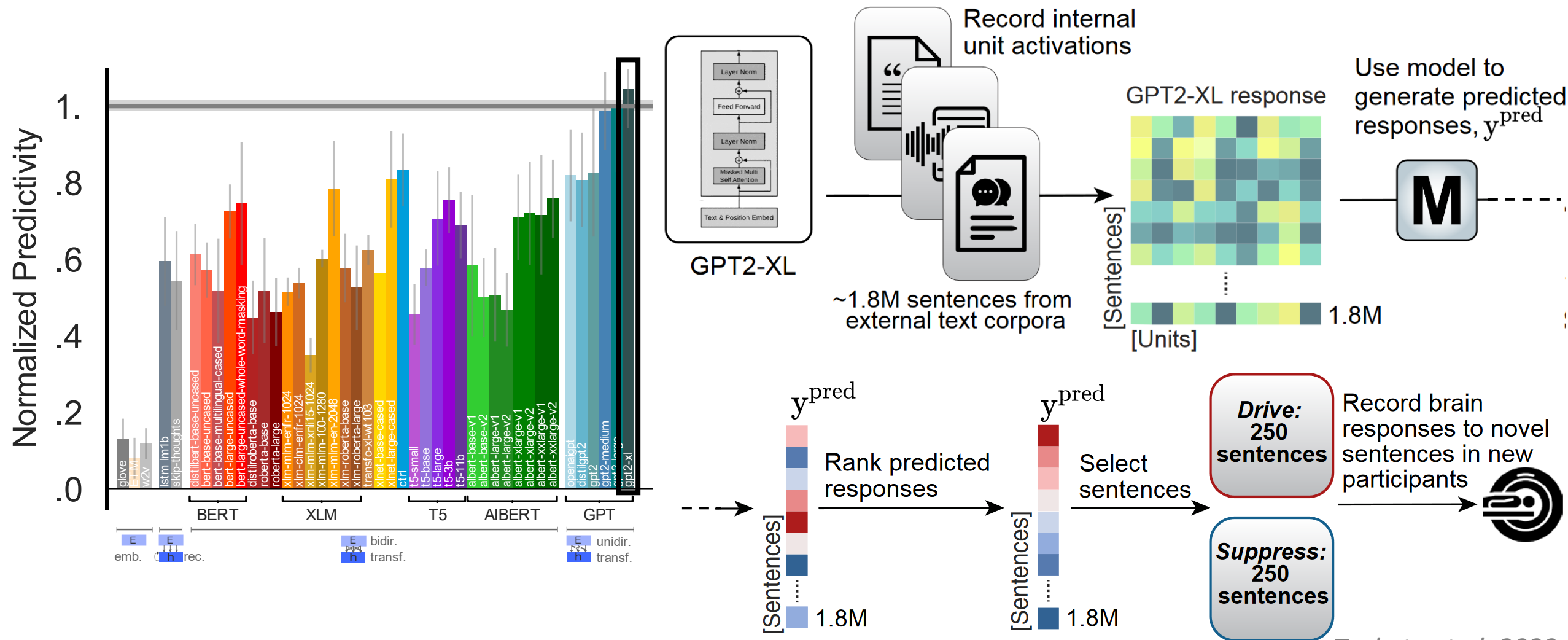
Inherent structure might be a key driver of brain-like language representations

Large feature size without structure is insufficient

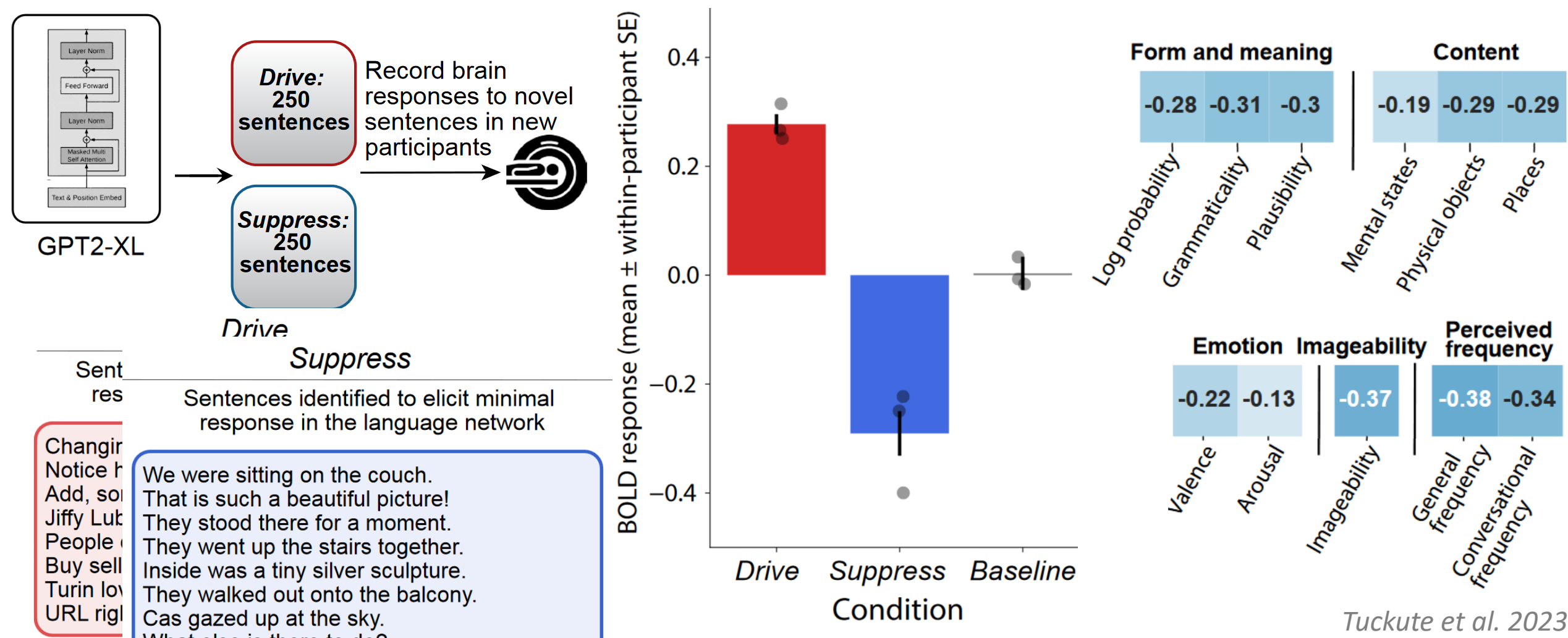
LLMs align to the brain's language system after developmentally realistic amounts of training



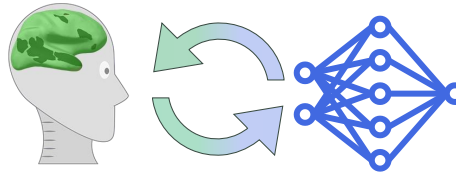
We can use brain-aligned LLMs to noninvasively control neural activity



We can use brain-aligned LLMs to noninvasively control neural activity



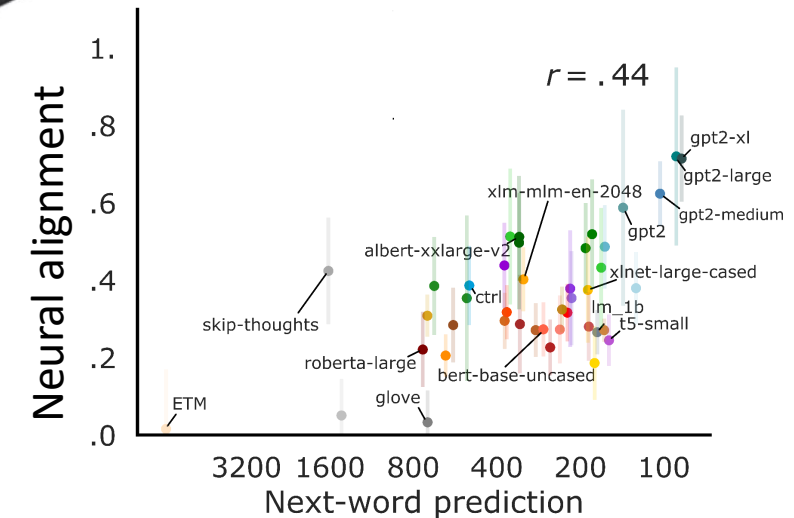
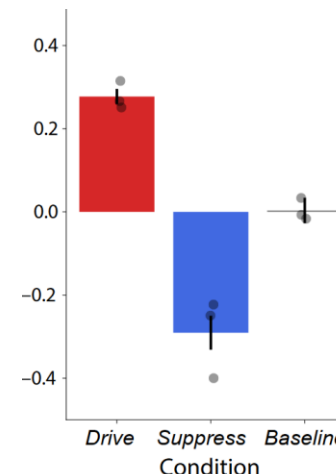
Contributions



1 Particular LLMs are strong models of the human language system

2 Next-word prediction performance relates to brain and behavioral alignment

3 The best models can be used to noninvasively control neural activity



www.Brain-Score.org/language