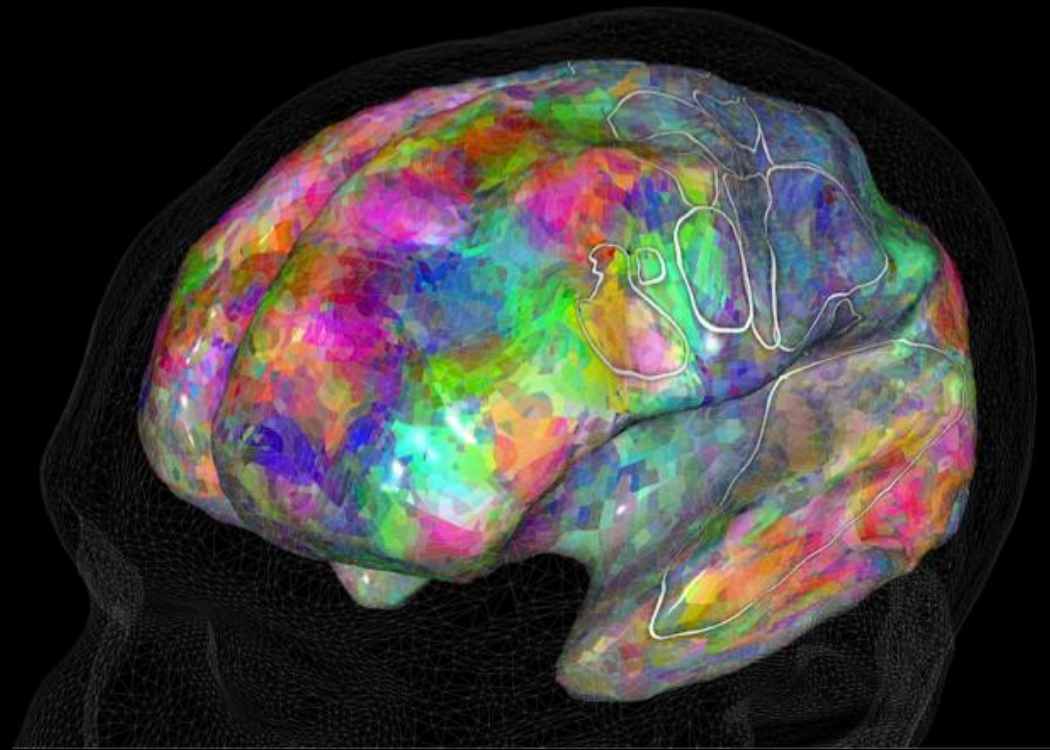


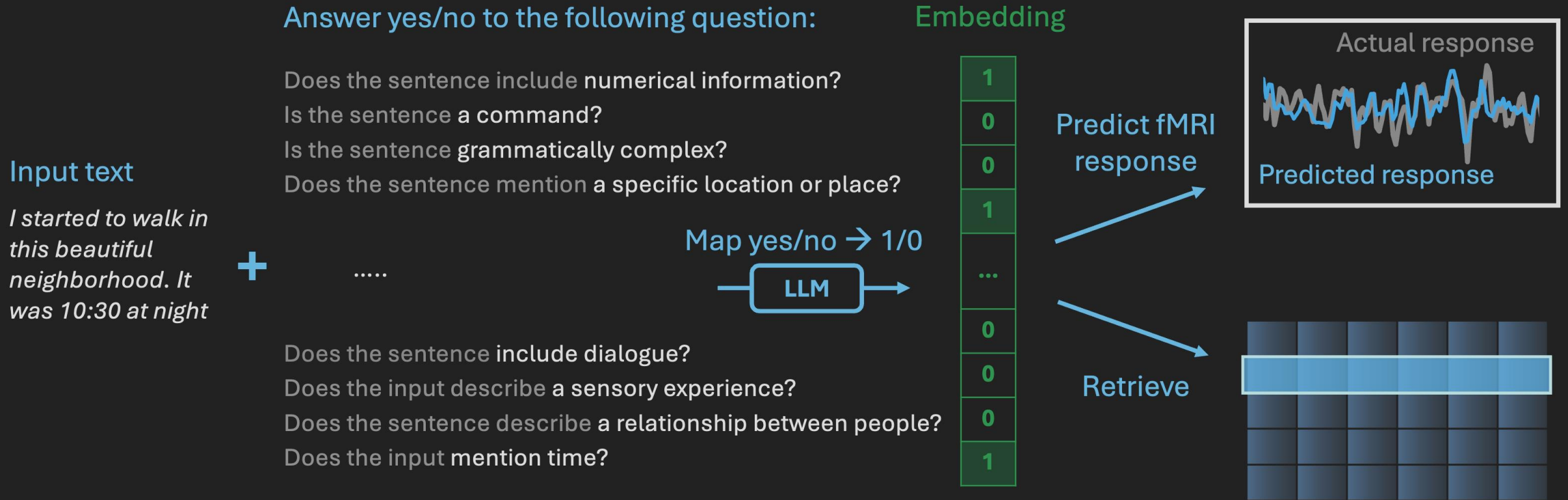
Crafting Interpretable Embeddings by Asking LLMs Questions

Vinamra Benara*, Chandan Singh*,
Jack Morris, Richard Antonello,
Ion Stoica, Alexander G. Huth, Jianfeng Gao

Microsoft Research
UC Berkeley
UT Austin
Cornell University



LLMs let us build explanation-based embeddings



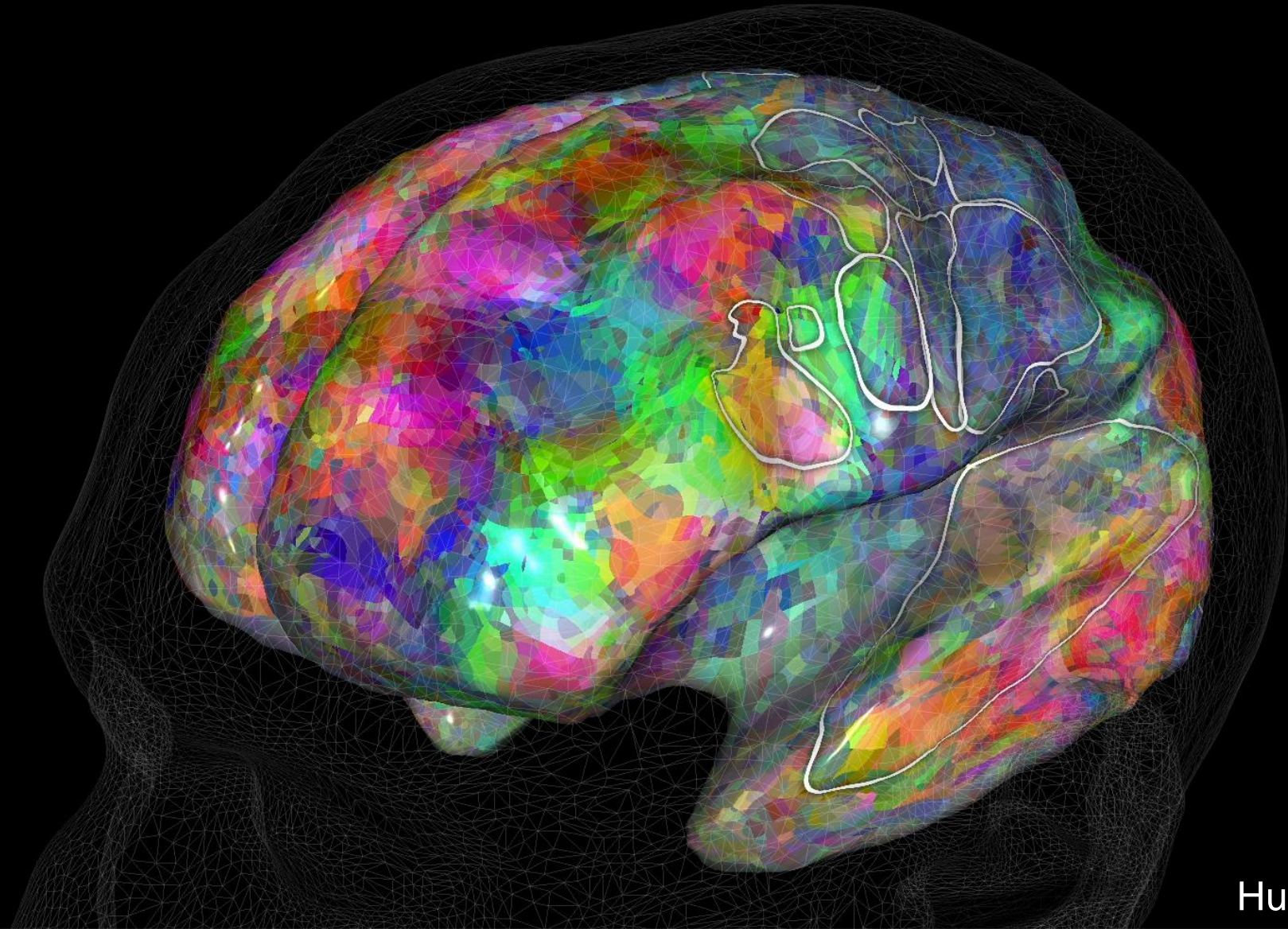
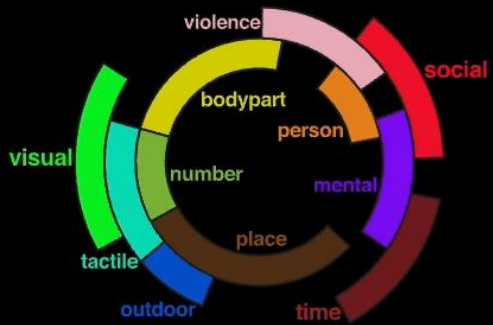
Limitations: computational cost, potentially inaccurate

Other works using yes-no questions:

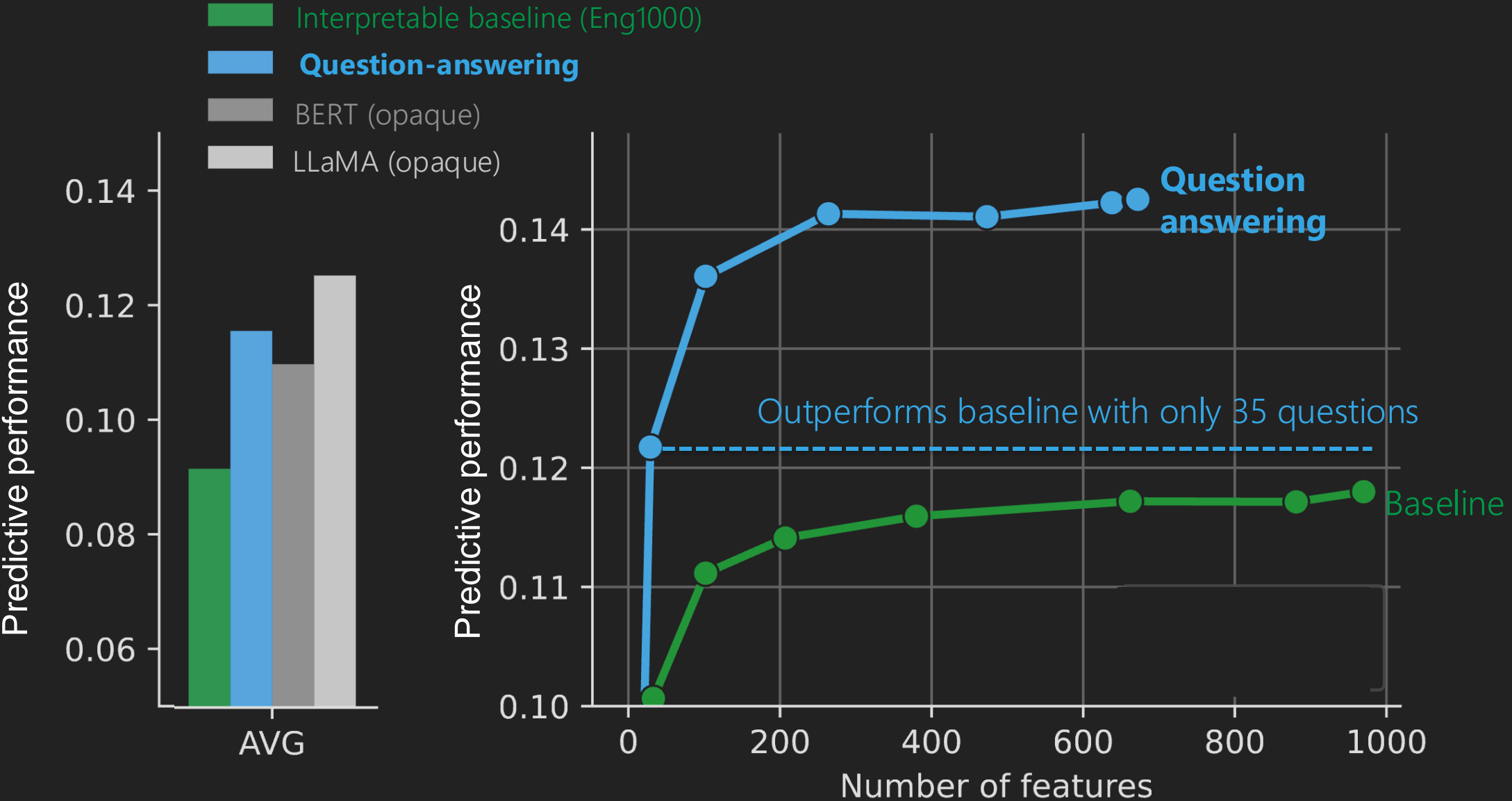
Style Embeddings (Patel et al. 2023); CHiLL (McInerney et al. 2023); Tree-Prompting (Morris et al. 2023); BC-LLM (Feng et al. 2024)

Explaining semantic representations is a fundamental goal of language neuroscience

voxel selectivity
colors show approximate semantic selectivity



A small set of questions predicts well



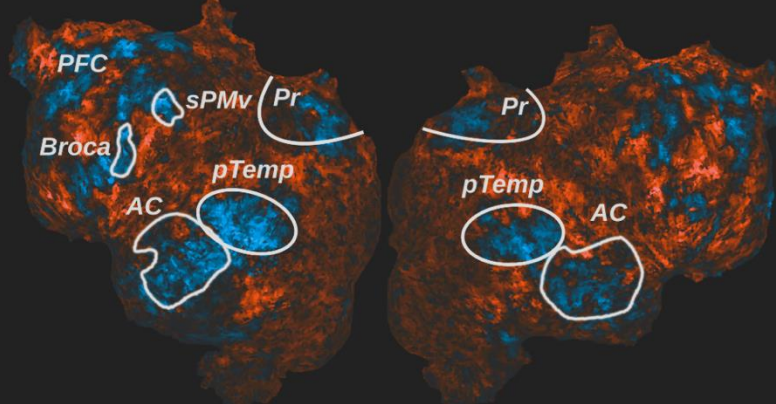
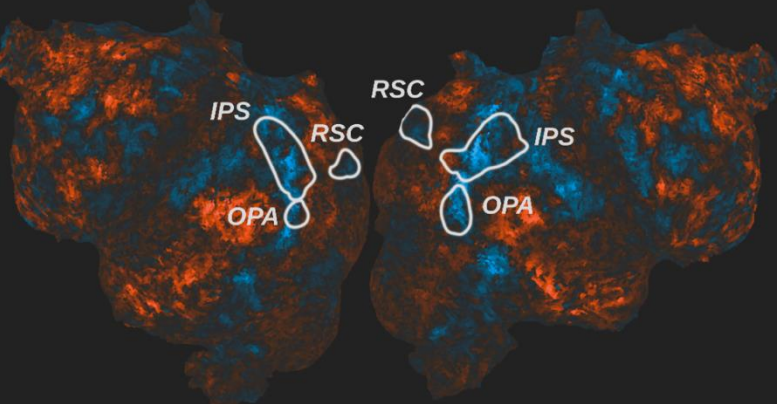
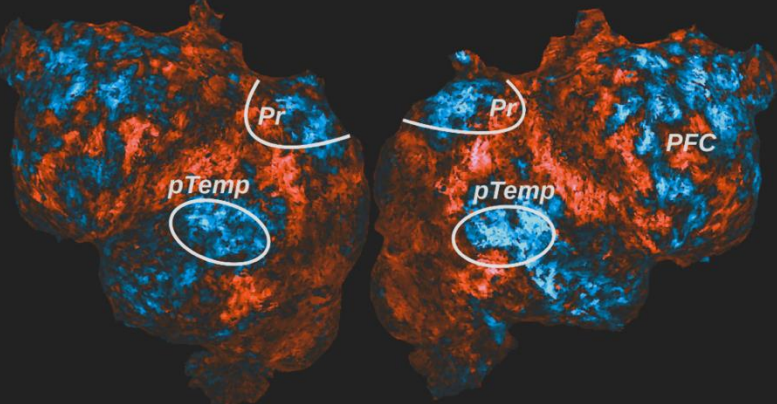
Each question yields a cortex-level selectivity map

Does the sentence describe a physical action?

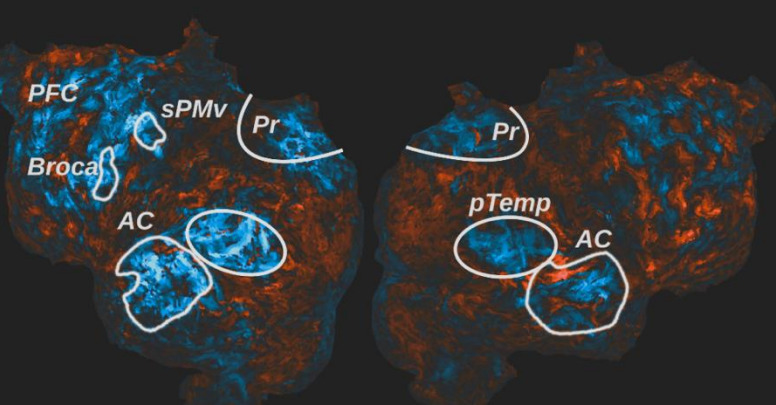
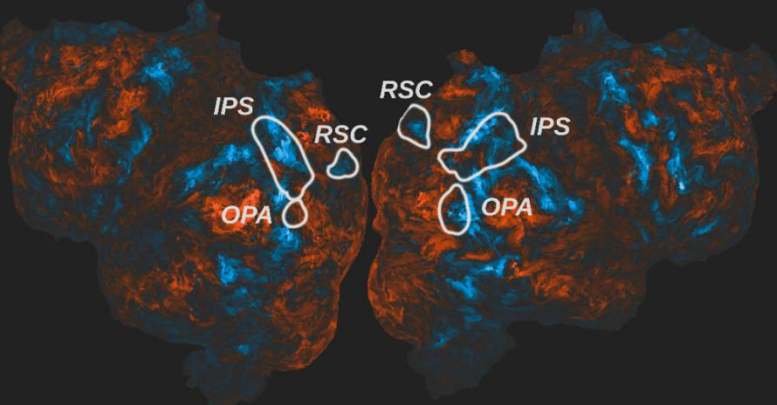
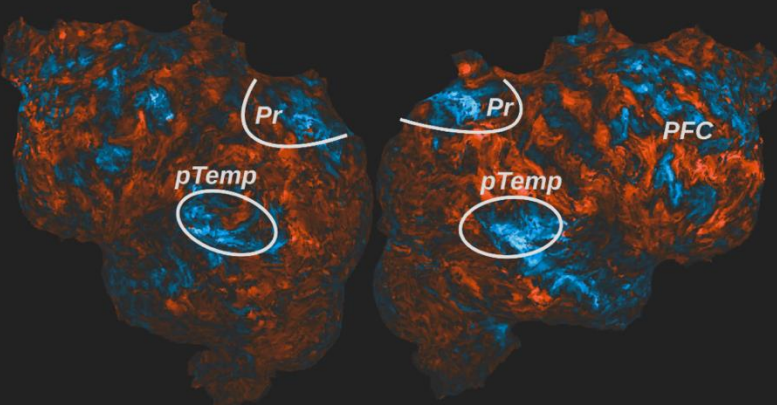
Does the sentence involve a description of a physical environment or setting?

Is the sentence grammatically complex?

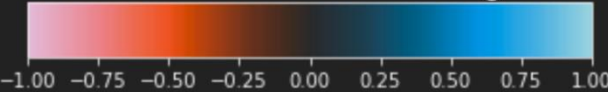
S02



S03



Rescaled Feature Weight



Ongoing work: Causally test these explanations

Generative explanation-mediated validation (antonello*, singh* et al. 2024, arXiv)

Build story from voxel explanations

Begin a story about **food preparation**... Use key phrases such as “slicing cucumber” or “zesting lemons”.

As I meticulously sliced the cucumber, the rhythmic sound...

Continue the story, but now make it about **locations**...

I packed the food and headed downtown on fifth street...

...

