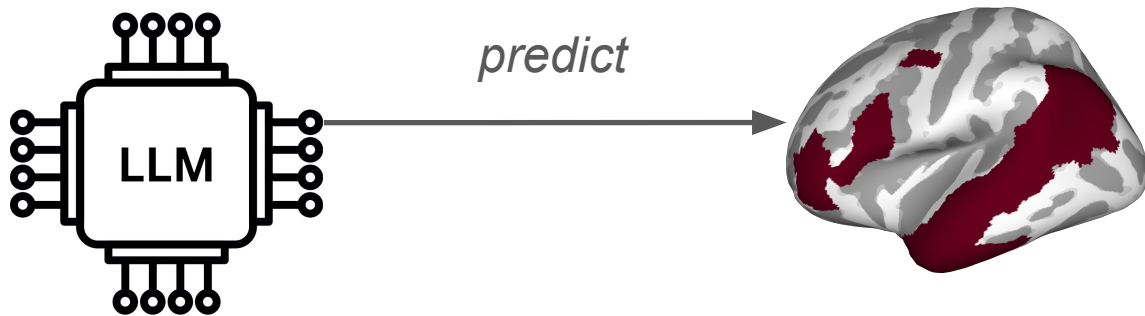# Brain-informed Fine-tuning
## for improved multilingual understanding in
# Language Models

Anuja Negi*, Subba Reddy Oota*, Anwar O Nunez-Elizalde,  Manish Gupta, Fatma Deniz
Technical University of Berlin / Bernstein Center for Computational Neuroscience Berlin/ Microsoft

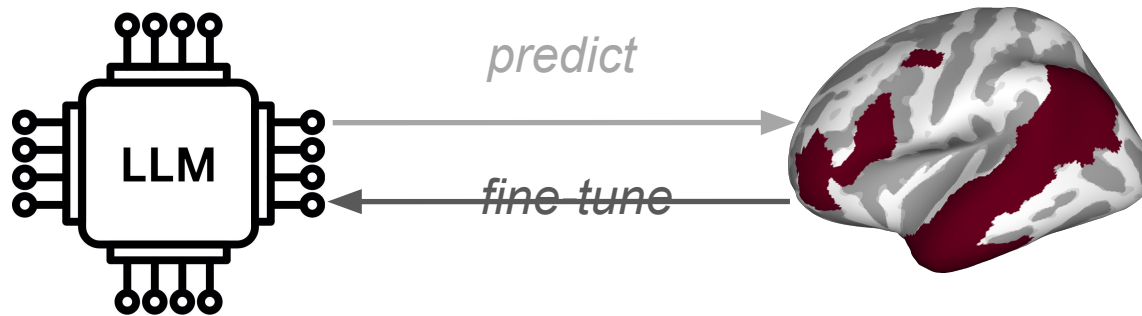# Language models accurately predict brain activity during language processing

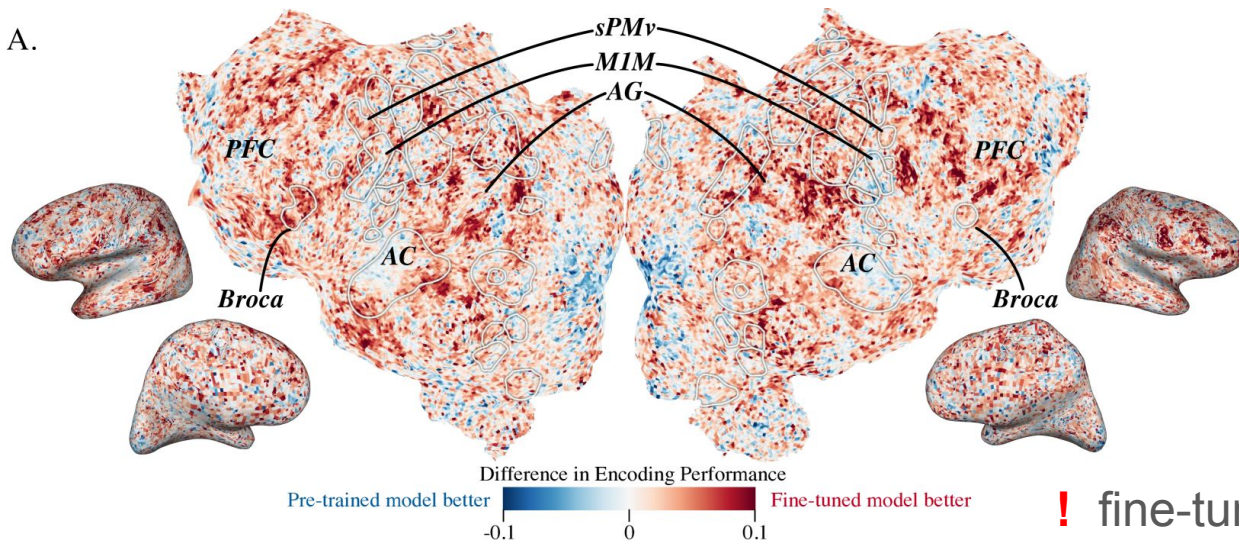Wehbe et al., 2014b; Jain & Huth, 2018; Toneva & Wehbe, 2019; Schrimpf et al., 2021; Caucheteux & King, 2022; Goldstein et al., 2022; Karamolegkou et al., 2023; Oota et al., 2025

*LLM: Large Language Model

# Fine-tuning language models with brain data



predict

fine tune

LLM

Schwartz et al., 2019

# Fine-tuning language models with brain data



A.

sPMv
M1M
AG
PFC
AC
Broca

PFC
AC
Broca

**Difference in Encoding Performance**
Pre-trained model better ___ Fine-tuned model better
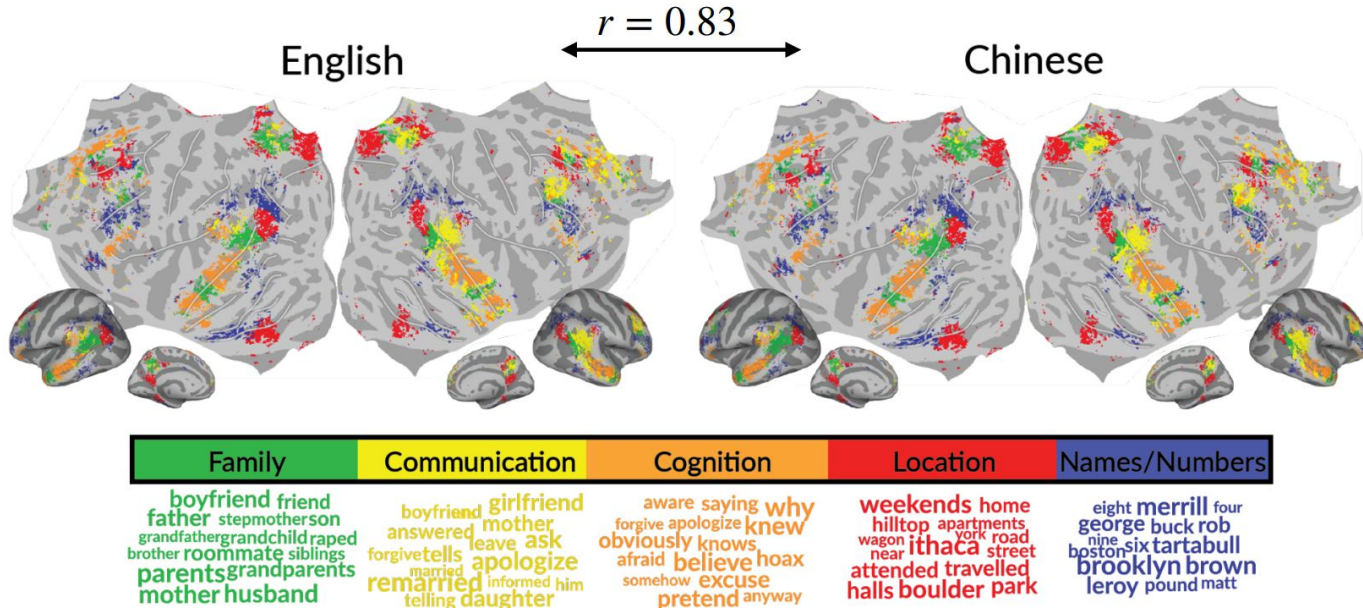-0.1     0     0.1

➔ improves <u>alignment</u> with the brain

➔ improves their semantic <u>downstream task performance</u>

❗ fine-tuning with monolingual brain data (English)
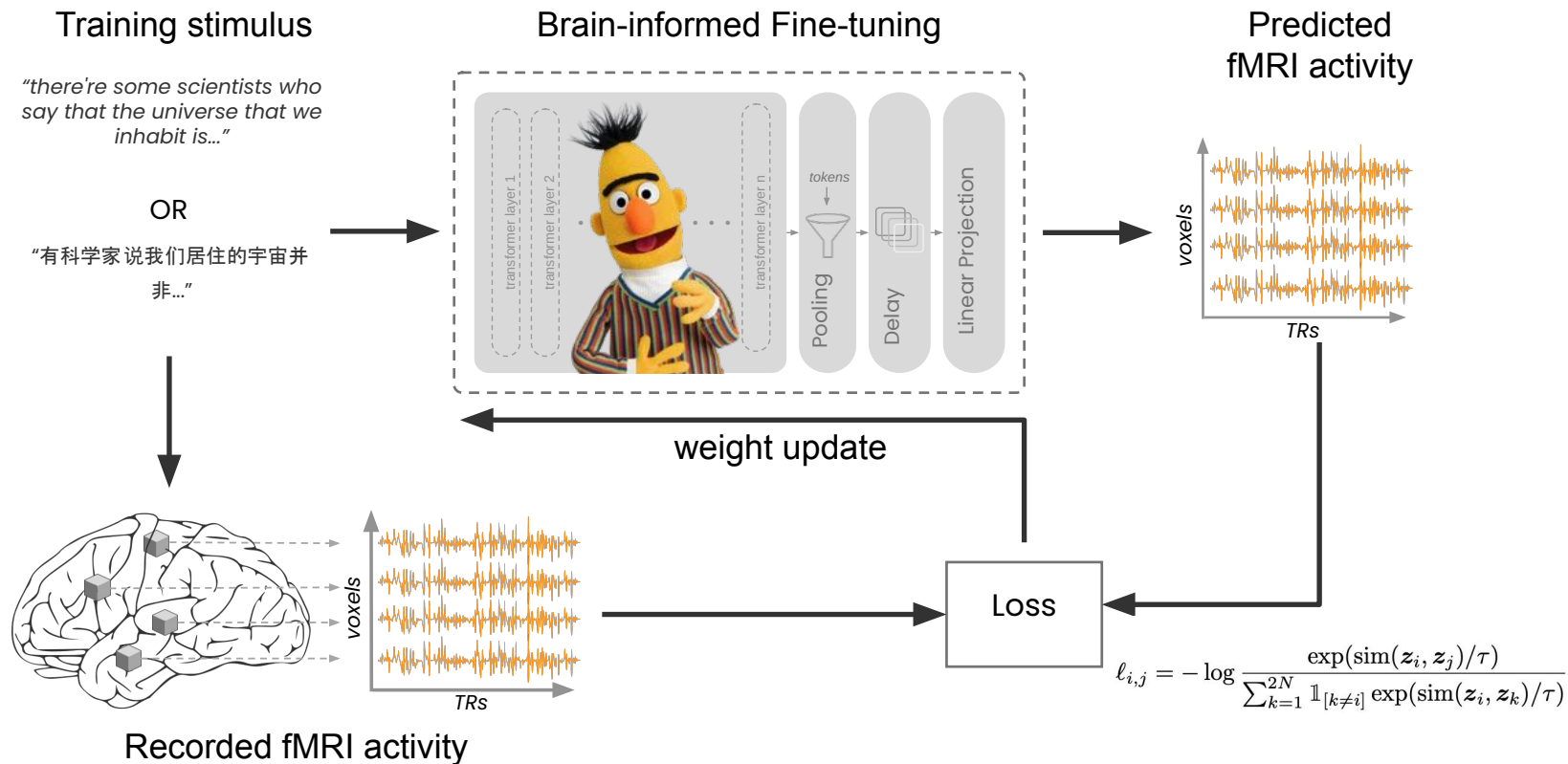
❗ only monolingual models were evaluated

Schwartz et al., 2019; Vattikonda et al. 2025; Moussa et al., 2025

# Shared semantic representations in bilinguals



r = 0.83

English                    Chinese

→ Bilingual language processing relies on shared semantic representations

| Family | Communication | Cognition | Location | Names/Numbers |

Family: boyfriend friend father stepmother son grandfather grandchild raped brother roommate siblings parents grandparents mother husband

Communication: boyfriend girlfriend mother answered leave ask forgive tells apologize married remarried informed him telling daughter

Cognition: aware saying why forgive apologize knew obviously knows afraid believe hoax somehow excuse pretend anyway

Location: weekends home hilltop apartments wagon york road near ithaca street attended travelled halls boulder park

Names/Numbers: eight merrill four george buck rob nine six tartabull boston brooklyn brown leroy pound matt

Chen, Gong, Tseng, Klein, Gallant, Deniz 2024

Can **fine-tuning** language models with **bilingual brain data** elicit multilingual capabilities in them?
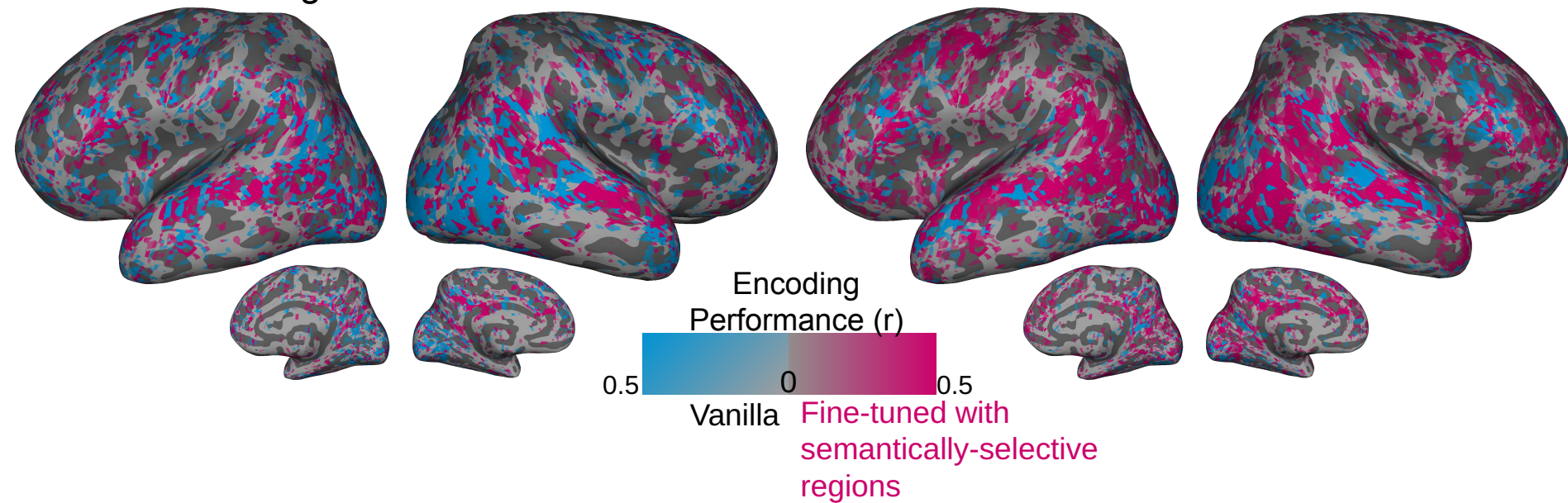
# Brain-informed fine-tuning with bilingual brain data



Training stimulus

*"there're some scientists who say that the universe that we inhabit is..."*

OR

*"有科学家说我们居住的宇宙并非..."*

Brain-informed Fine-tuning

tokens

transformer layer 1
transformer layer 2
transformer layer n
Pooling
Delay
Linear Projection

Predicted fMRI activity

voxels

TRs

weight update

voxels

TRs

Recorded fMRI activity

Loss

$$\ell_{i,j} = -\log \frac{\exp(\mathrm{sim}(\boldsymbol{z}_i, \boldsymbol{z}_j)/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(\mathrm{sim}(\boldsymbol{z}_i, \boldsymbol{z}_k)/\tau)}$$

Negi*, Oota*, Nunez-Elizalde, Gupta, Deniz 2025

# **Evaluating** brain-informed fine-tuned models

**Evaluation**

🧠 Brain Alignment
(Voxelwise Encoding)

📄 NLP Evaluation
(Downstream Tasks)

# Brain-informed fine-tuning improves brain alignment



BERT-en fine-tuned with English brain data
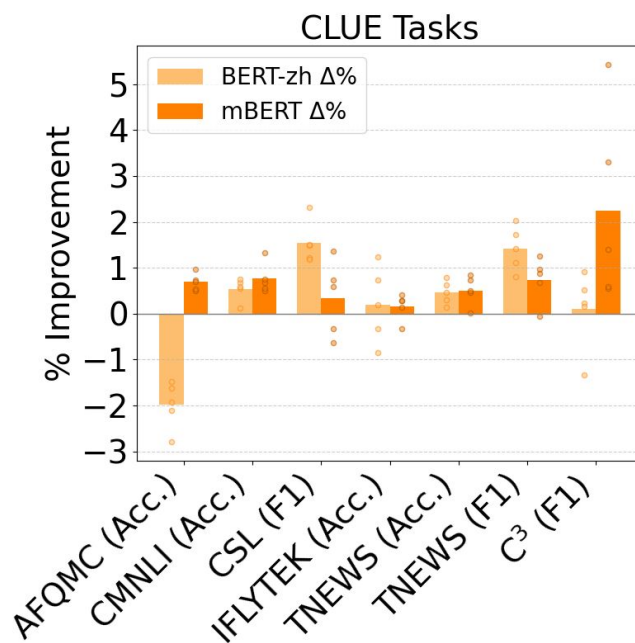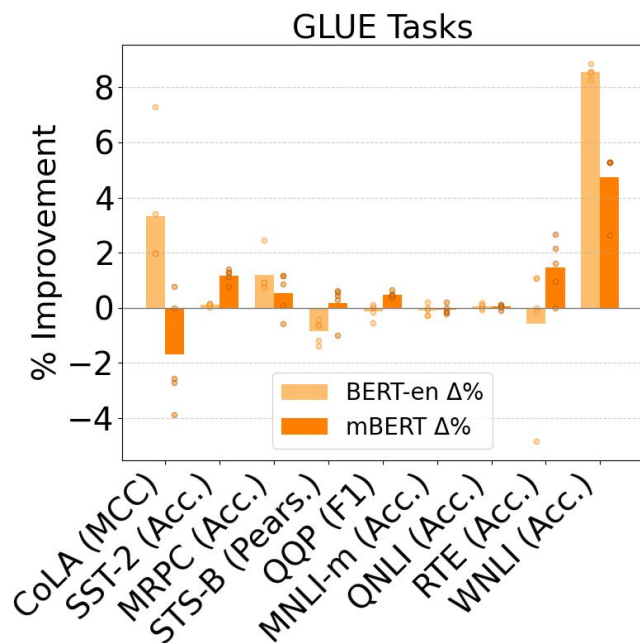
BERT-zh fine-tuned with Chinese brain data

Encoding Performance (r)

0.5 — 0 — 0.5

Vanilla

Fine-tuned with semantically-selective regions

~ *70% semantic voxels prefer a fine-tuned model over vanilla model*

Negi*, Oota*, Nunez-Elizalde, Gupta, Deniz 2025

# Evaluating brain-informed fine-tuned models



Evaluation
- Brain Alignment (Voxelwise Encoding)
- NLP Evaluation (Downstream Tasks)

GLUE Benchmark (en)

CoLA, SST, MRPC, STS, QQP, MNLI, QNLI, RTE, WNLI

CLUE Benchmark (zh)

AQFMC, CMNLI, IFLYTEK, TNEWS, CHID, C3

● Paraphrase/Semantic Similarity  ● Natural Language Inference
● Classification/Sentiment  ● Coreference/Structure

Negi*, Oota*, Nunez-Elizalde, Gupta, Deniz 2025

# Fine-tuning improves linguistic task performance



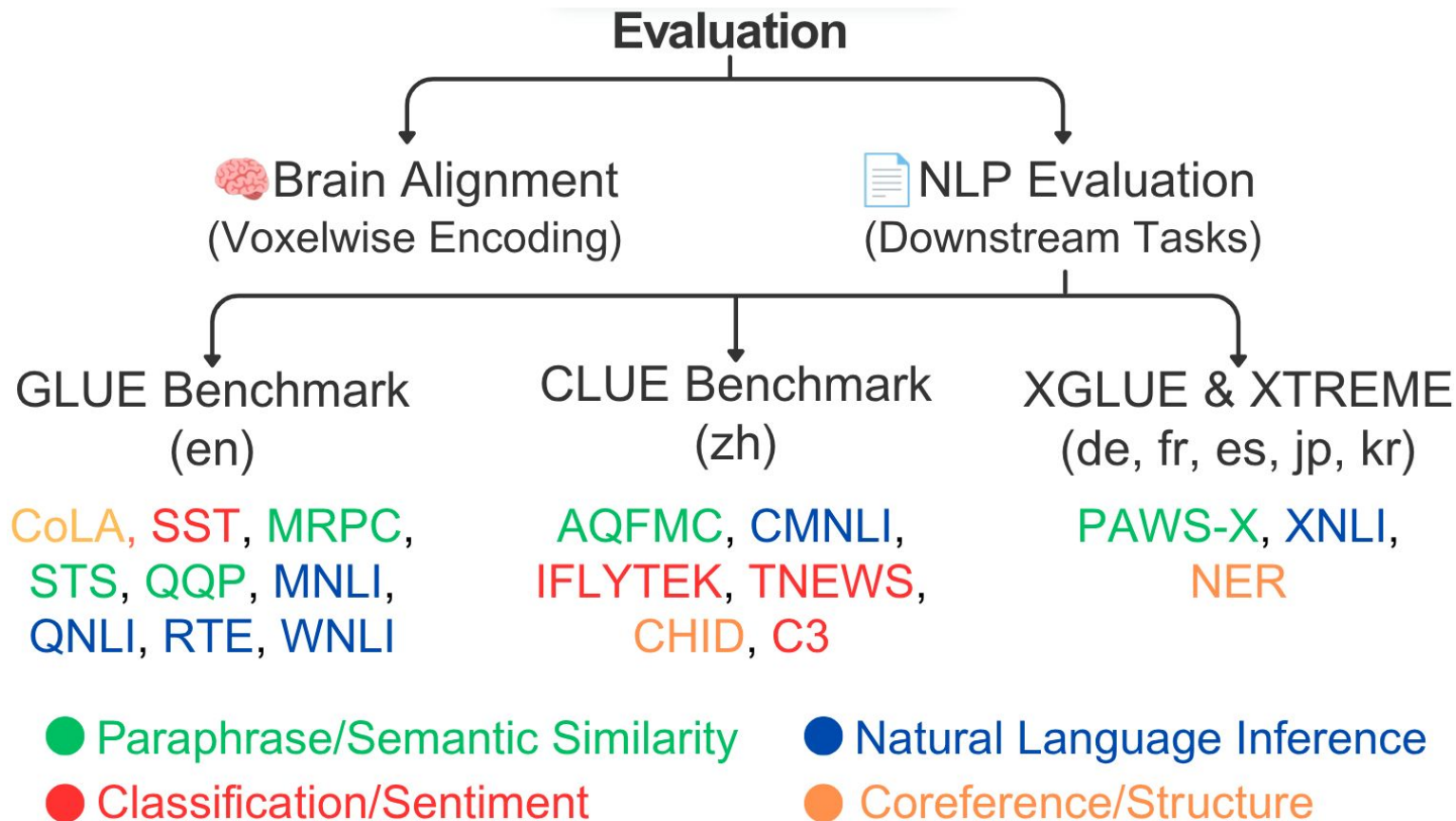GLUE Tasks

CLUE Tasks

## Monolingual model
⬆️ 7/9 on English benchmark (GLUE)
⬆️ 6/7 on Chinese benchmark (CLUE)
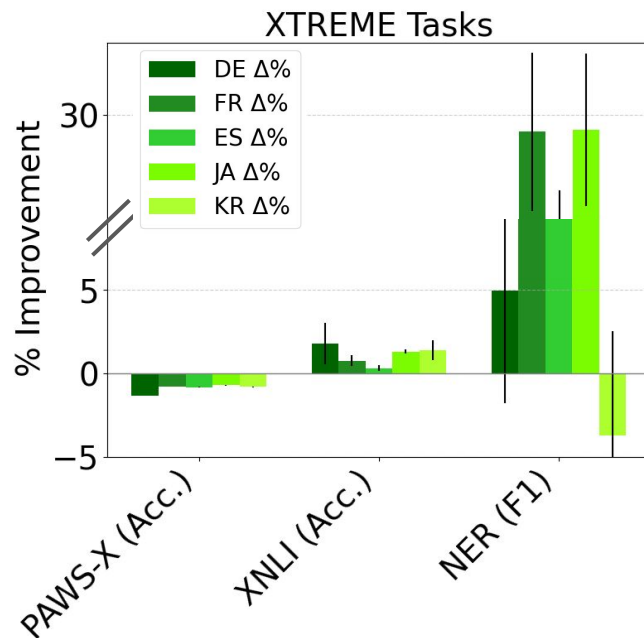
## Multilingual model
⬆️ 8/9 on English benchmark (GLUE)
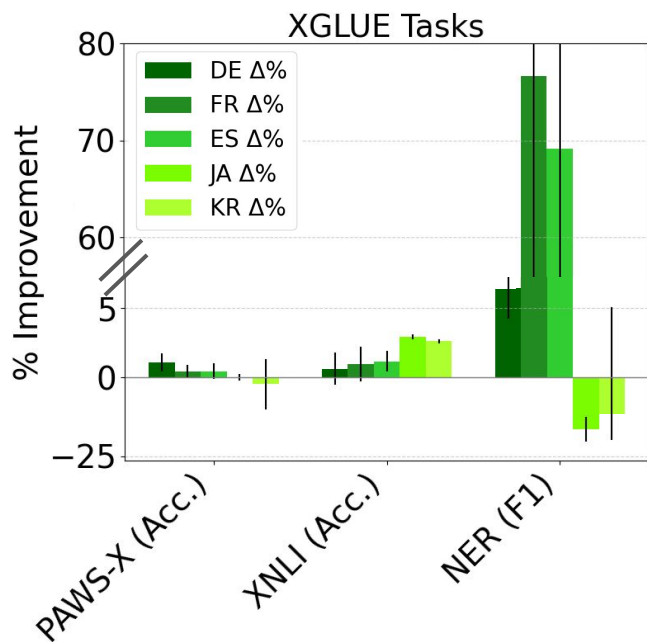⬆️ 7/7 on Chinese benchmark (CLUE)

# Fine-tuning enables cross-language transfer



GLUE Tasks

CLUE Tasks

**Monolingual model**
⬆️ 8/9 on English benchmark (GLUE)
⬆️ 6/7 on Chinese benchmark (CLUE)

**Multilingual model**
⬆️ 9/9 on English benchmark (GLUE)
⬆️ 7/7 on Chinese benchmark (CLUE)

Negi*, Oota*, Nunez-Elizalde, Gupta, Deniz 2025

# Evaluating brain-informed fine-tuned models

# Fine-tuning improves language-agnostic representations



⬆️ 2/3 on DE, FR, ES

# Cross-linguistic transfer is because of bilingual brain



⬆️ 6/7 bilingual > monolingual

Potentially driven by shared semantics:

- ○ in bilingual brains (Chen et al., 2024)
- ○ across different languages in the brain (de Varda et al., 2025)

# Conclusions

➔ First study to perform brain-informed fine-tuning using bilingual brain data.

➔ Brain-informed fine-tuning improves
  ◆ brain alignment
  ◆ downstream task performance across within-, cross-, and unseen language settings.

➔ Improvements are driven specifically by fine-tuning with bilingual brain data, not brain data in general.

➔ Potential of leveraging bilingual brain representations for developing language-agnostic models.

**Future Work:** Explore which linguistic properties the model captures (e.g., syntax, morphology, discourse) to improve model training and evaluation.