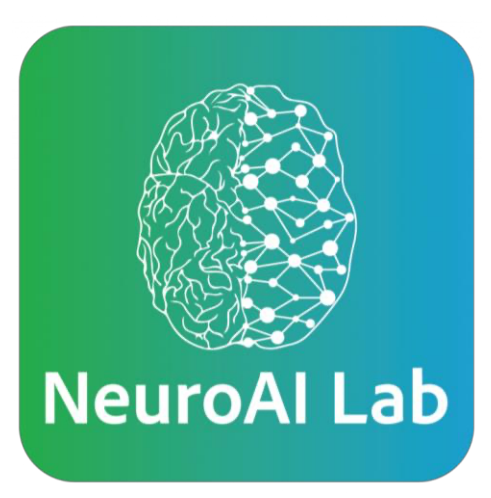# The LLM Language Network:
## How LLMs Outgrow the Human Language Network

**Badr AlKhamissi**[1]   Greta Tuckute[2]   Yingtian Tang[1]   Taha Binhuraib[3]

Antoine Bosselut*,[1]   Martin Schrimpf*,[1]

[1] EPFL   [2] MIT   [3] GT

*Equal Supervision

@bkhmsi

---

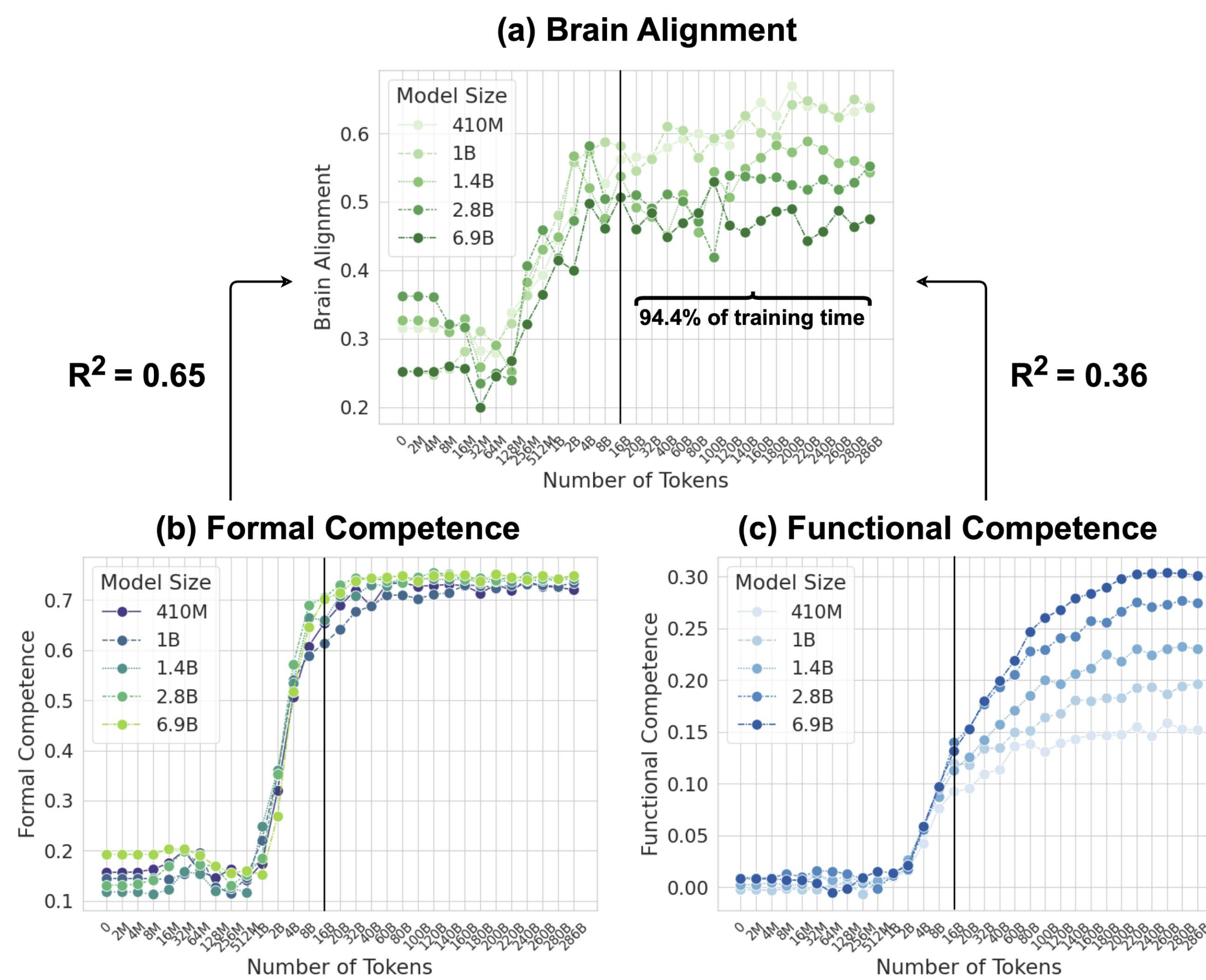## * Highlights

1. **Untrained models** align with brain via **context integration**

2. **Formal linguistic** competence drives alignment early, **saturates ~4B tokens**

3. **Functional competence** emerges later, with **weaker brain correlation**

4. **Correlation** between models' **brain alignment** and their **next-word-prediction** performance, as well as their **behavioral alignment fades over time.**

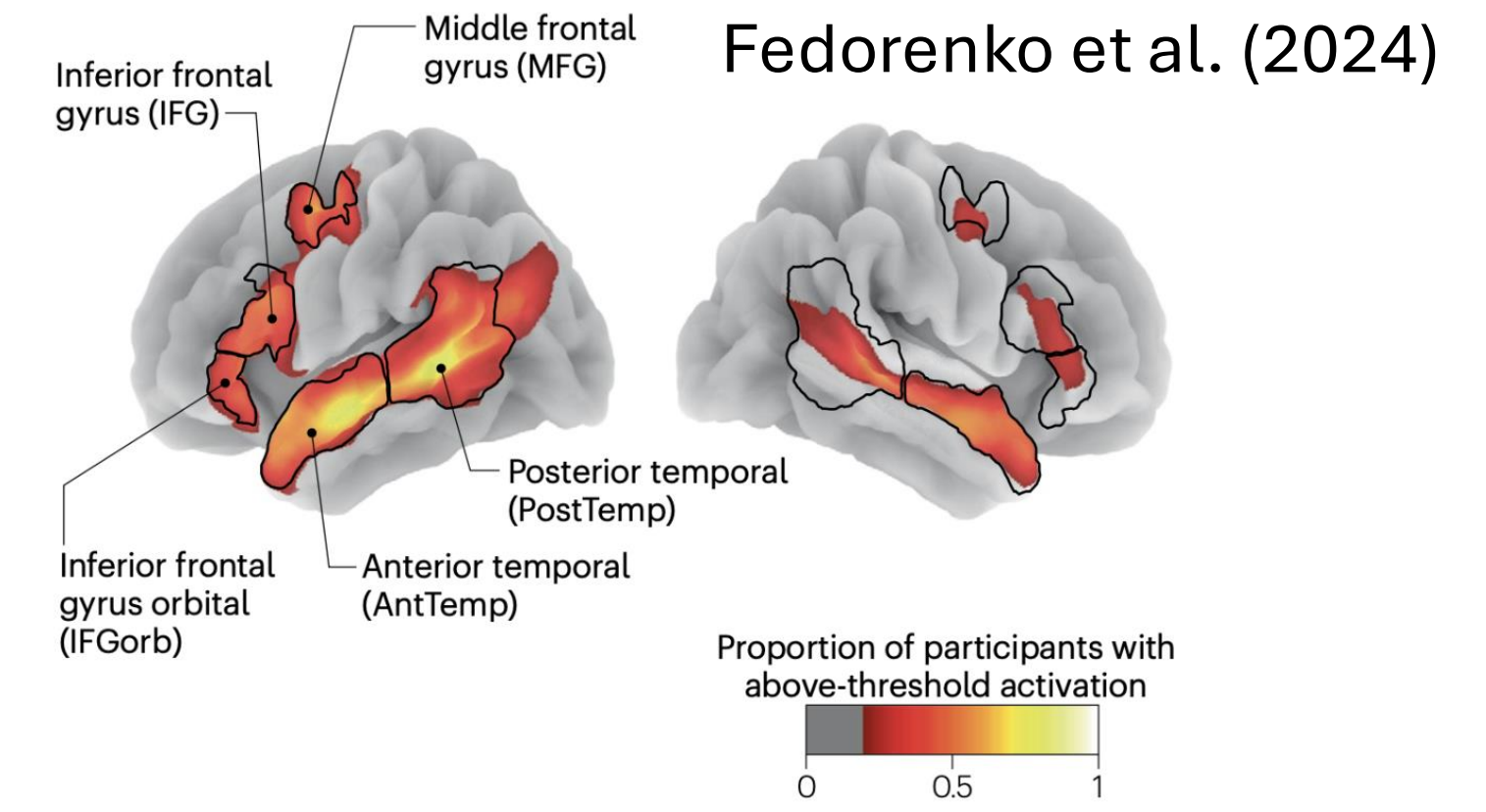5. **Model size ≠ better brain alignment** (when controlling features).

---

## 1 Brain Alignment Tracks Formal More Than Functional Competence



(a) Brain Alignment

94.4% of training time

$R^2 = 0.65$   $R^2 = 0.36$

(b) Formal Competence

(c) Functional Competence

---

## * Human Language Network



Fedorenko et al. (2024)

Specialized area within the brain responsible for understanding and producing language.

---

## 2 Methods

1. Benchmarked **34 checkpoints**
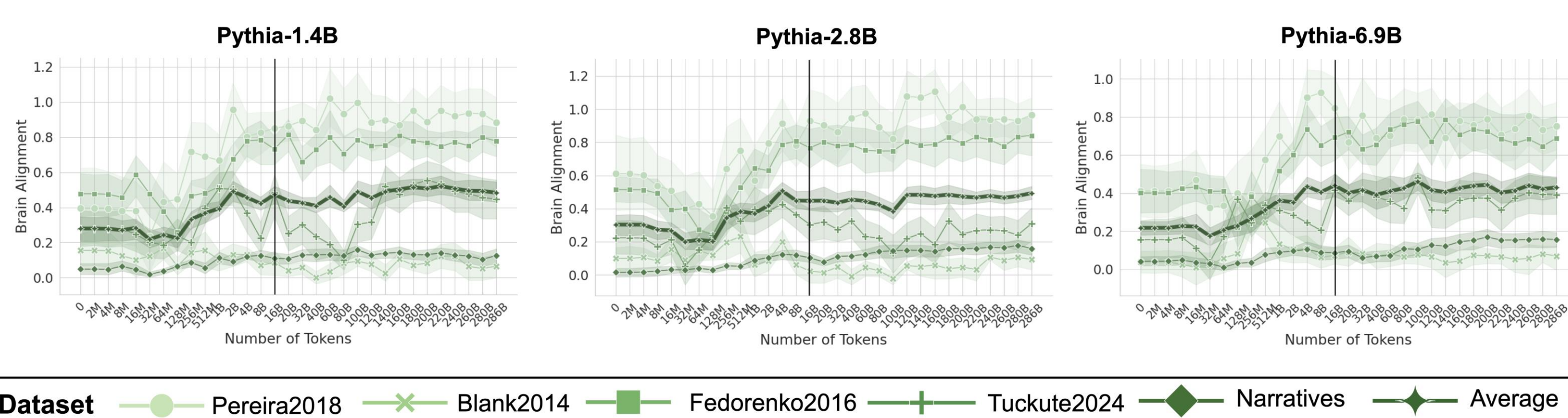2. Spanning **~300B tokens**
3. Across **8 different model sizes**
4. On **5 brain-recording datasets**, and **1 behavioral dataset**
5. And on **2 formal linguistic benchmarks and 6 functional**

---

## Research Questions

**What drives brain alignment of LLMs?**
**Is it primarily linked to formal or functional linguistic competence?**
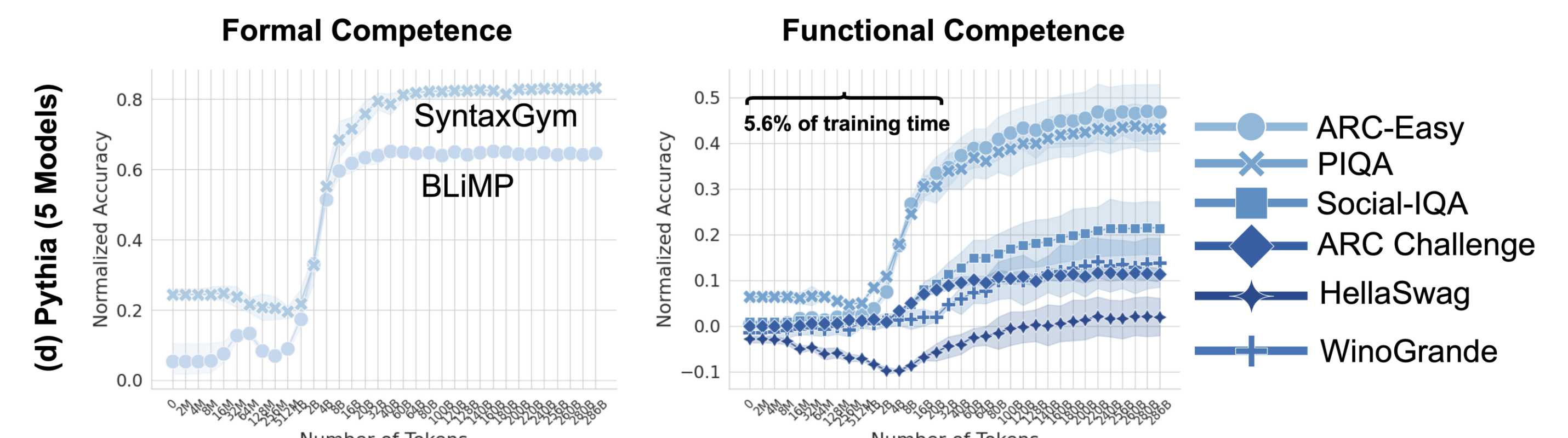**Do LLMs diverge from humans as they surpass human-level prediction?**
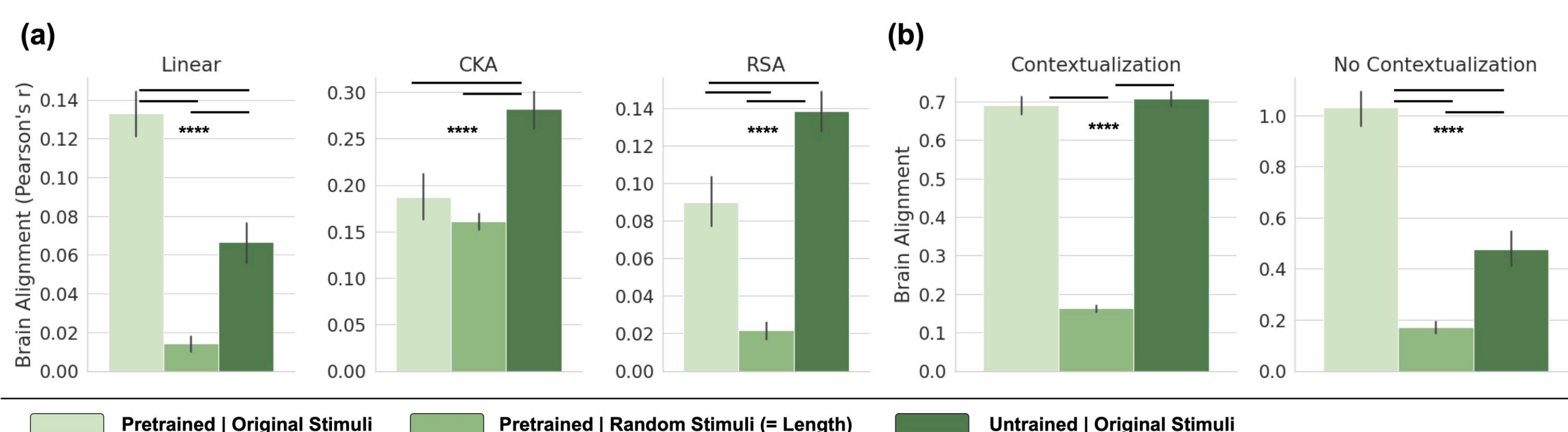
---

## 3 Brain Alignment Over Training Per Dataset



Pythia-1.4B   Pythia-2.8B   Pythia-6.9B

Dataset: Pereira2018   Blank2014   Fedorenko2016   Tuckute2024   Narratives   Average

---

## 4 Formal & Functional Competence Over Training Per Dataset



(d) Pythia (5 Models)

Formal Competence: SyntaxGym, BLiMP

Functional Competence: 5.6% of training time

ARC-Easy   PIQA   Social-IQA   ARC Challenge   HellaSwag   WinoGrande

---

## 5 Rigorous Brain-Scoring



(a) Linear   CKA   RSA   (b) Contextualization   No Contextualization

Pretrained | Original Stimuli   Pretrained | Random Stimuli (= Length)   Untrained | Original Stimuli

---

## 6 Context Integration drives Brain Alignment of Untrained Models



(a) Architecture: MLP, GRU, LSTM, MLP+Mean, Transformer-v1, Transformer-v2

(b) Pos+Attn+MLP, Attn+MLP, Attn, Pos+Attn, MLP, Pos+MLP, Tokens

(c) MLP, LayerNorm, Multihead Attention, LayerNorm, Tokens, Pos Embeddings

(d) Formal, Functional

---

## 7 Brain Tracks Formal > Functional Competence



(a) Pythia (5 Models)   $R^2 = 0.65$   $R^2 = 0.36$
(b) Pythia-2.8B   $R^2 = 0.51$   $R^2 = 0.40$

Legend: Brain Alignment, Formal Competence, Functional Competence

---

## 8 Brain Aligns with NWP & Behavior Early



(a) Pythia (8 Models)   r = 0.26**   r = 0.81****   r = n.s.   r = 0.84****
(b) Pythia-2.8B   r = n.s.   r = 0.63*   r = -0.54*   r = 0.89****

Training Stage: Early, Late

---

## 9 Model Size ≠ Better Brain Alignment



Pythia Model Size: 14M   70M   160M   410M   1B   1.4B   2.8B   6.9B