# MoE Lens - An Expert Is All You Need

Marmik Chaudhari[1], Idhant Gulati[1], Nishkal Hundia[2], Pranav Karra[1] & Shivam Raval[3]

[1]Penn State University  [2]University of Maryland  [3]Harvard University

**TL;DR:** We discover that specialization in Mixture of Experts models is localized among a few experts across various domains, with the top–weighted expert closely approximating full ensemble's next token predictions.

## Introduction

- The original aim of Mixture of Experts (MoEs) was to allow models to scale capacity without proportional increases in computational costs by sparsely activating a subset of parameters (experts) for each input token.

- However, scaling MoEs and expert count effectively has faced limitations due to its discrete expert routing such as training instability and forced generalization across multiple distinct semantic concepts for each expert, hindering specialization.

- Current MoEs face a fundamental tension between achieving fine–grained specialization through larger expert counts and the observed parameter inefficiency where only a small subset of experts handle most computation across domains.

- Significant progress in interpretability methods, such as Logit Lens, has shown success in uncovering the activations of language models – what insights can we gain when we apply these techniques to analyze expert specialization in MoEs?

## Methods

### 1. Tracing token routing distribution

Number of tokens from domain $D$ for which expert $E_i$ is selected as one of the Top-k experts

$$\text{Expert specialization}(E_i, D) = \frac{N_{E_i,D}^{(k)}}{N_D}$$

Total number of tokens from domain $D$

- We compute expert specialization by measuring the fraction of tokens from each domain routed to each expert.
- Domains were represented by distinct datasets: GSM8K & AIME, Github subset of Paloma, FrenchQA (FQuAD), subset of Chinese Fineweb Edu, Guttenberg English for Math, Code, French, Chinese and English domain respectively.
- We classify an expert as domain specialized if its routing frequency is significantly higher than a uniform routing baseline.
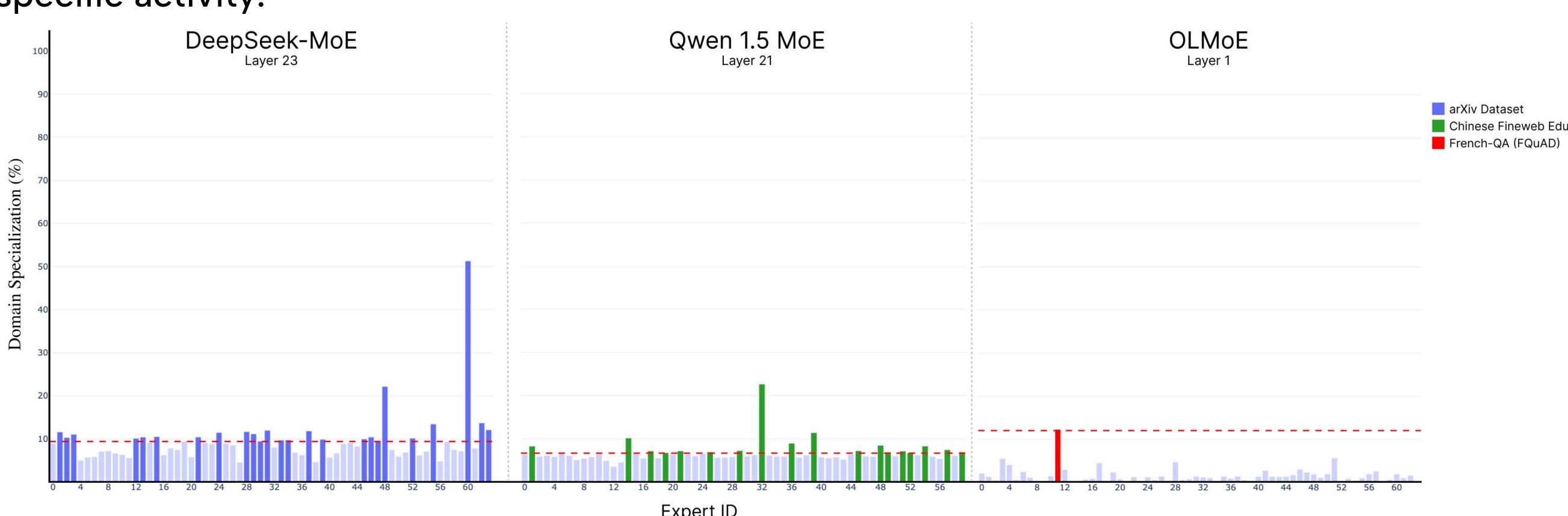
### 2. Logit Lens for analyzing expert contributions

Post-attention residual stream at layer $\ell$ for token $t$

$$\text{LogitLens}^{ext}(h_t^\ell) = \text{LayerNorm}(h_t^\ell + u_t^\ell)W_U$$

Hidden states at any intermediate layer $\ell$ for $t$-th token

Model's pretrained unembedding matrix

- For each layer $\ell$ and token $t$, we apply Logit Lens to project two different hidden states to the vocab space: individual expert output $E_i$, weighted sum of Top-k expert output ($H_t^\ell$).
- We also incorporate post-attention residual stream for projecting individual expert outputs.
- We perform these experiments on **DeepSeek–MoE, OLMoE, and Qwen 1.5 MoE** to validate our findings across different MoE implementations.

## Results

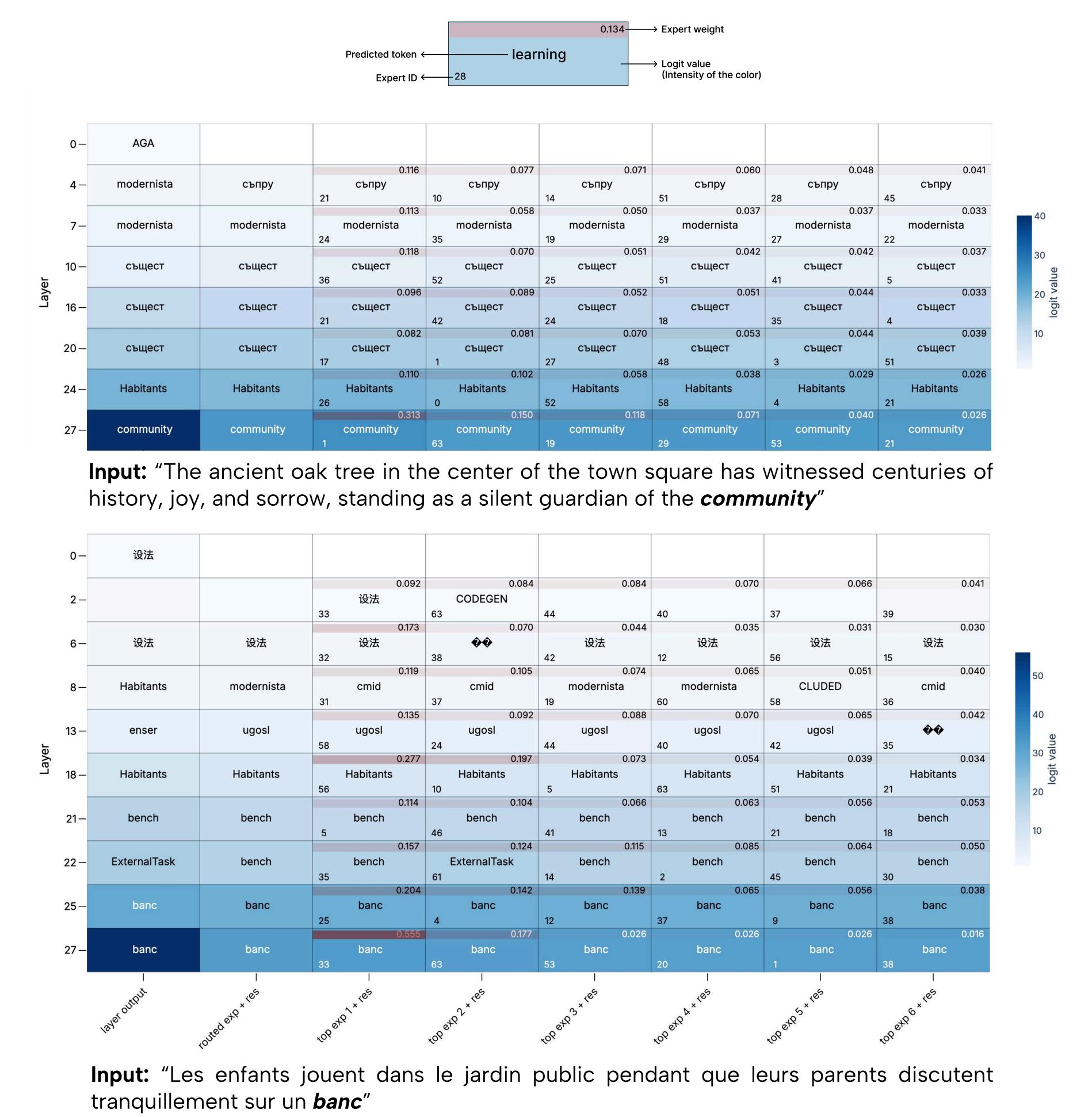### 1. Few experts dominate token routing across domains

The routing distribution reveals two key patterns in MoEs: **(1)** only a small number of experts show strong specialization for any domain. **(2)** most experts demonstrate minimal domain–specific activity.



**Figure 1:** Expert Specialization across Multiple MoE Architectures. We visualize domain–specific routing patterns for DeepSeek–MoE, Qwen 1.5 MoE, and OLMoE models across three different domains. The y–axis shows domain specialization percentage for each expert, with the red dashed line indicating the uniform routing baseline (~9.375% for DeepSeek–MoE, ~6.67% for Qwen 1.5 MoE, and ~12.5% for OLMoE)

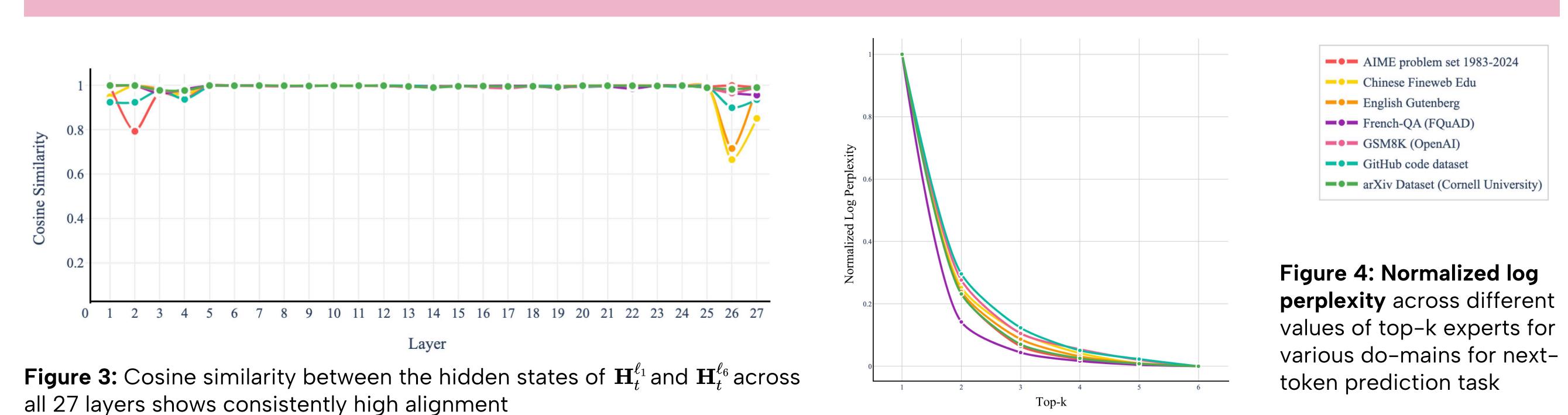### 2. Top-weighted expert approximates full ensemble's output

The extended Logit Lens gives us evidence that solely projecting the hidden state of the top weighted expert's output, $H_t^{\ell_1}$ across layers decode to roughly the same next token prediction as the output at the end of that layer, i.e, the hidden state of the Top-k weighted expert outputs, $H_t^{\ell_6}$.



**Input:** "The ancient oak tree in the center of the town square has witnessed centuries of history, joy, and sorrow, standing as a silent guardian of the **community**"



**Input:** "Les enfants jouent dans le jardin public pendant que leurs parents discutent tranquillement sur un **banc**"

**Figure 2.1–2:** Logit Lens visualization for DeepSeek–MoE on the three different input sequences. Each cell shows the top–1 token prediction across layers (rows) for layer output, routed experts with residual stream for various top-k values. Color intensity indicates prediction confidence. The lower-left subscript indicates expert indices and the top-right superscript indicates expert weight.

### 3. Cosine similarity and perplexity



**Figure 3:** Cosine similarity between the hidden states of $\mathbf{H}_t^{\ell_1}$ and $\mathbf{H}_t^{\ell_6}$ across all 27 layers shows consistently high alignment

**Figure 4: Normalized log perplexity** across different values of top-k experts for various do–mains for next–token prediction task

We also observe very high cosine similarity between $H_t^{\ell_1}$ and $H_t^{\ell_6}$ across all layers and each domain indicating that the top–weighted expert is contributing the most in shaping final output representation whereas the contributions of other experts are minimal in the hidden space.

The perplexity moderately increases when reducing Top-k $= 6 \rightarrow 1$ which validates the claim that the top–weighted expert, when combined with the residual stream, produces representations closely aligned with the layer output.