

# Mixture of Experts (MoEs)

Exploring trade-offs in MoE routing designs

Florian Kowarsch





Florian Kowarsch

2021-2023

2023-2025

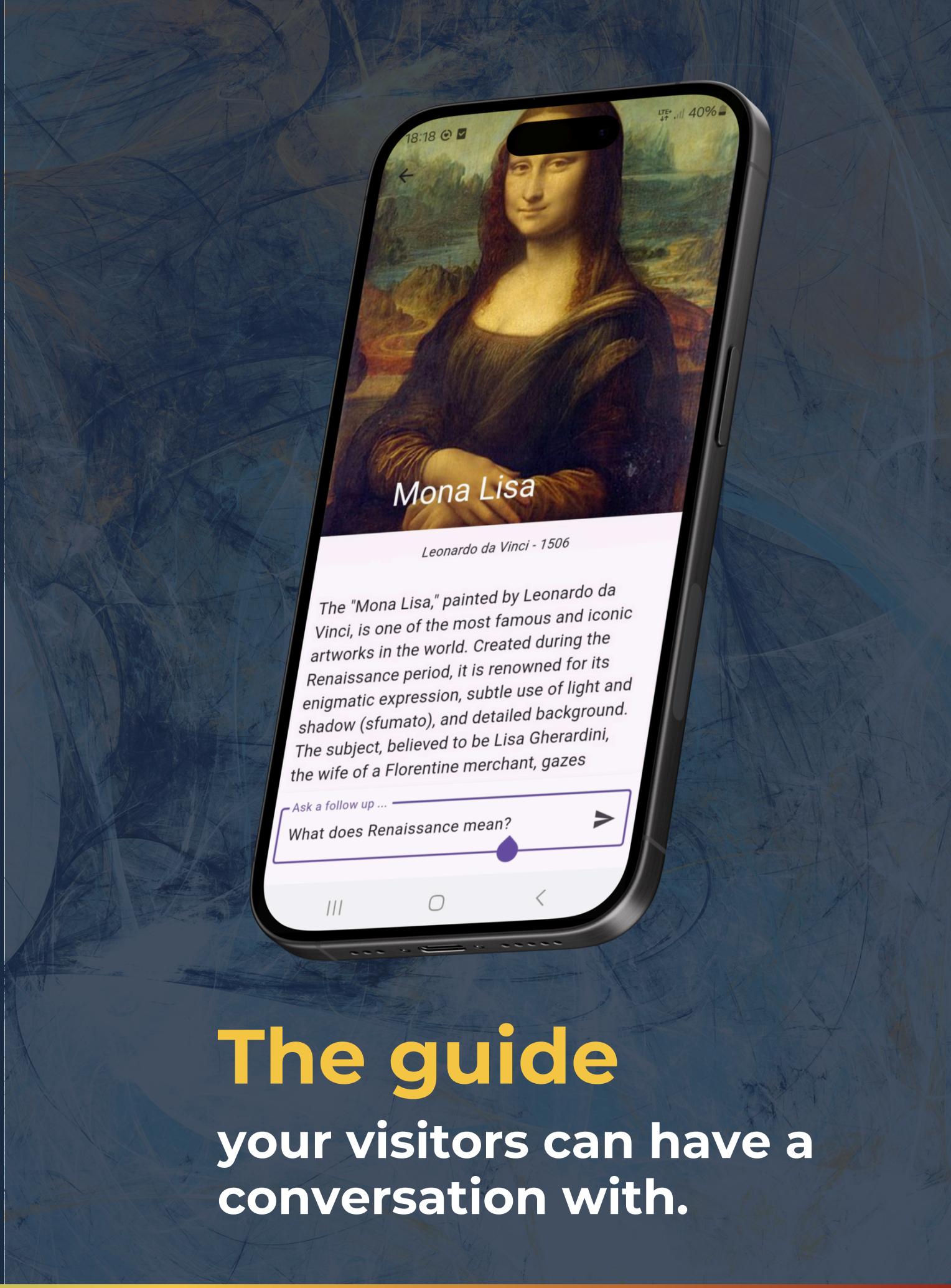


**SEVENTH  
SENSE**

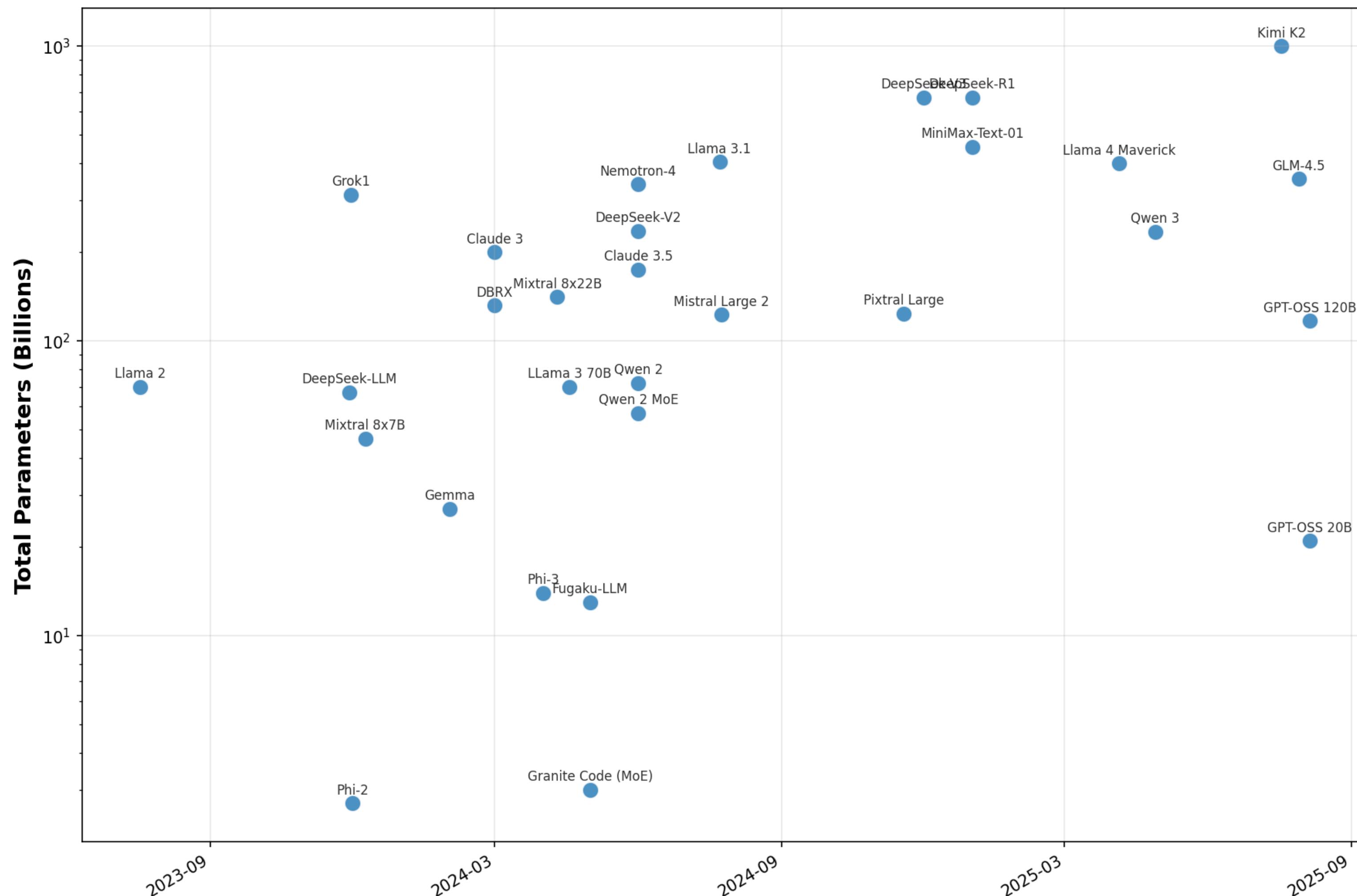


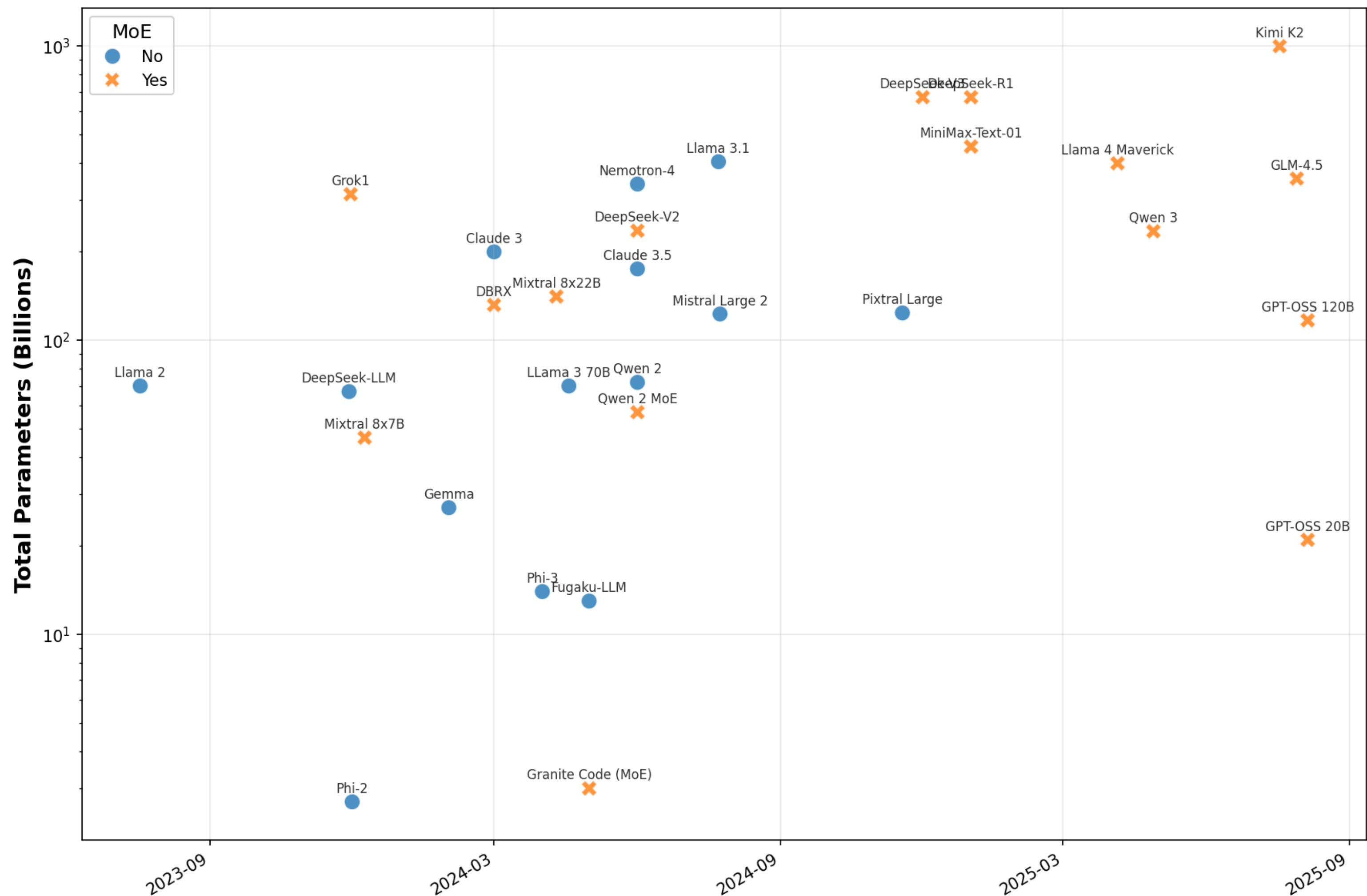
**Reducing Barriers,  
unleashing curiosity.**

[www.nuseum.ai](http://www.nuseum.ai)

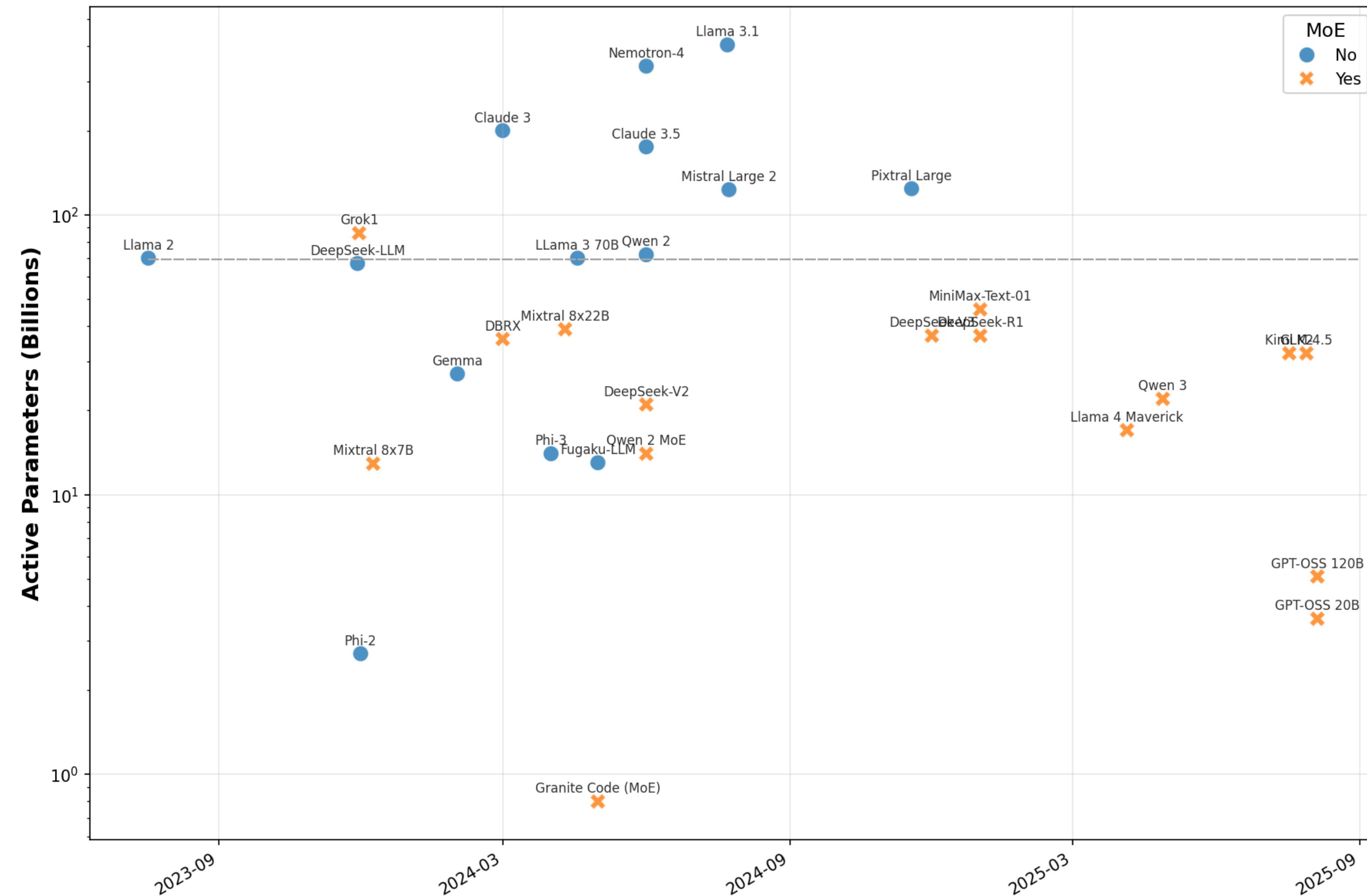


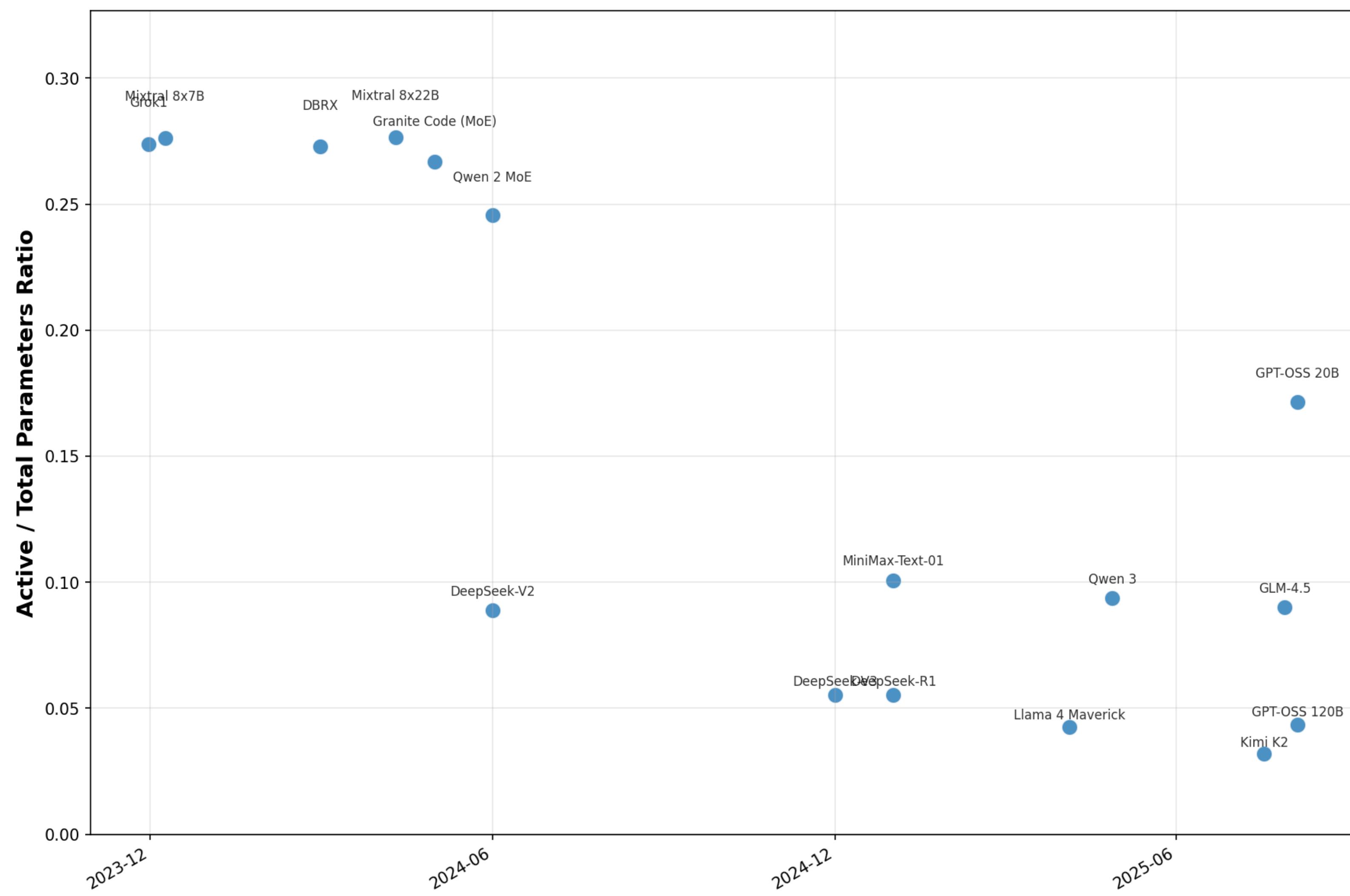
**The guide  
your visitors can have a  
conversation with.**





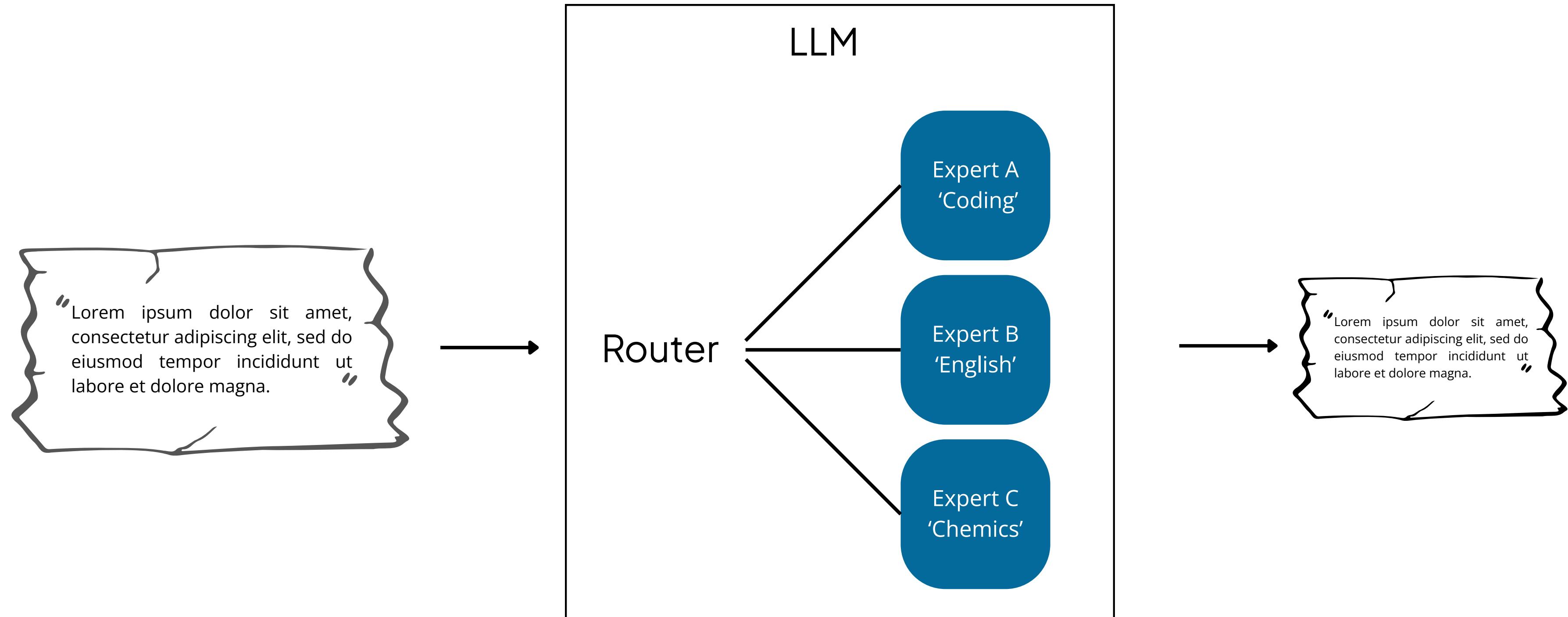
# Why MoEs?



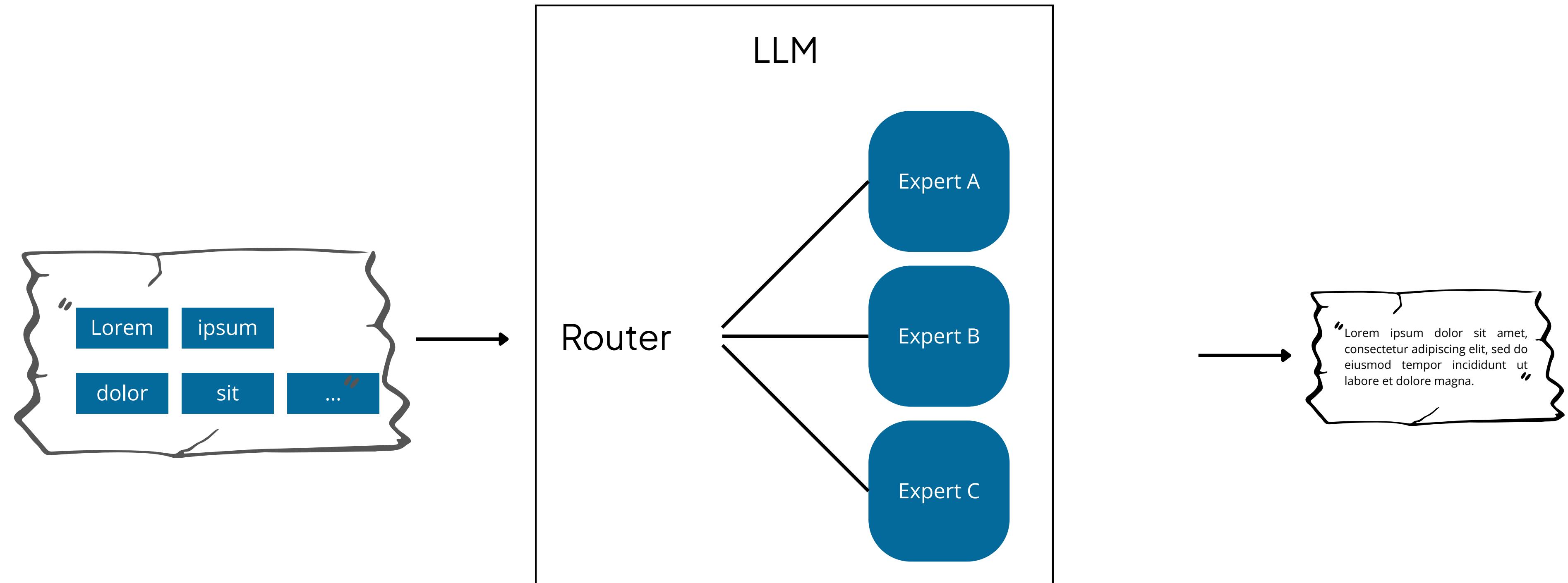


How does it work?

# The naive Concept



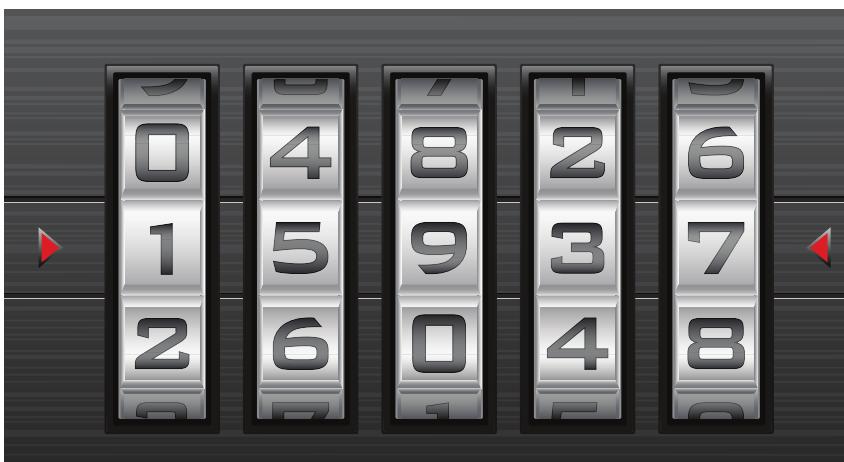
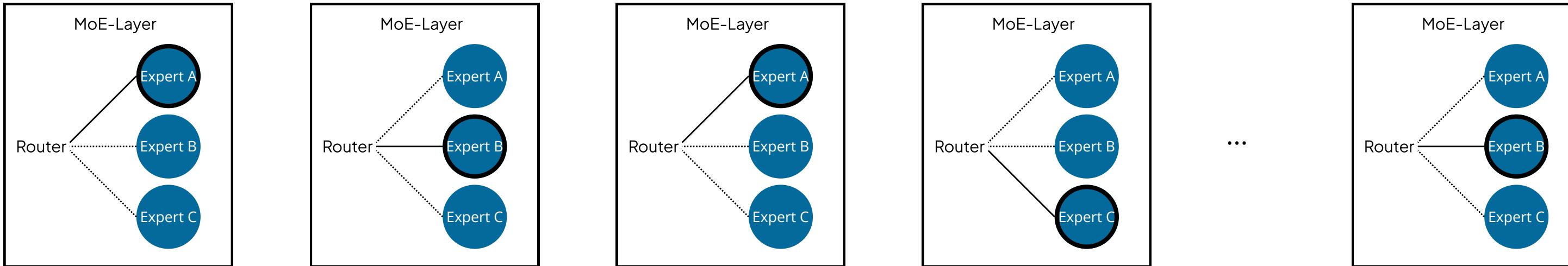
# #1 Token-Level Routing



## #2 There are many MoE-Layers



# #2 There are many MoE-Layers



**gpt-oss-120b**

128 experts per layer  
36 MoE layers

→  $128^{36}$  paths

# #3 There are no domain expert

Fuzhao Xue<sup>1†</sup> Zian Zheng<sup>1</sup> Yao Fu<sup>2</sup> Jinjie Ni<sup>1</sup> Zangwei Zheng<sup>1</sup>

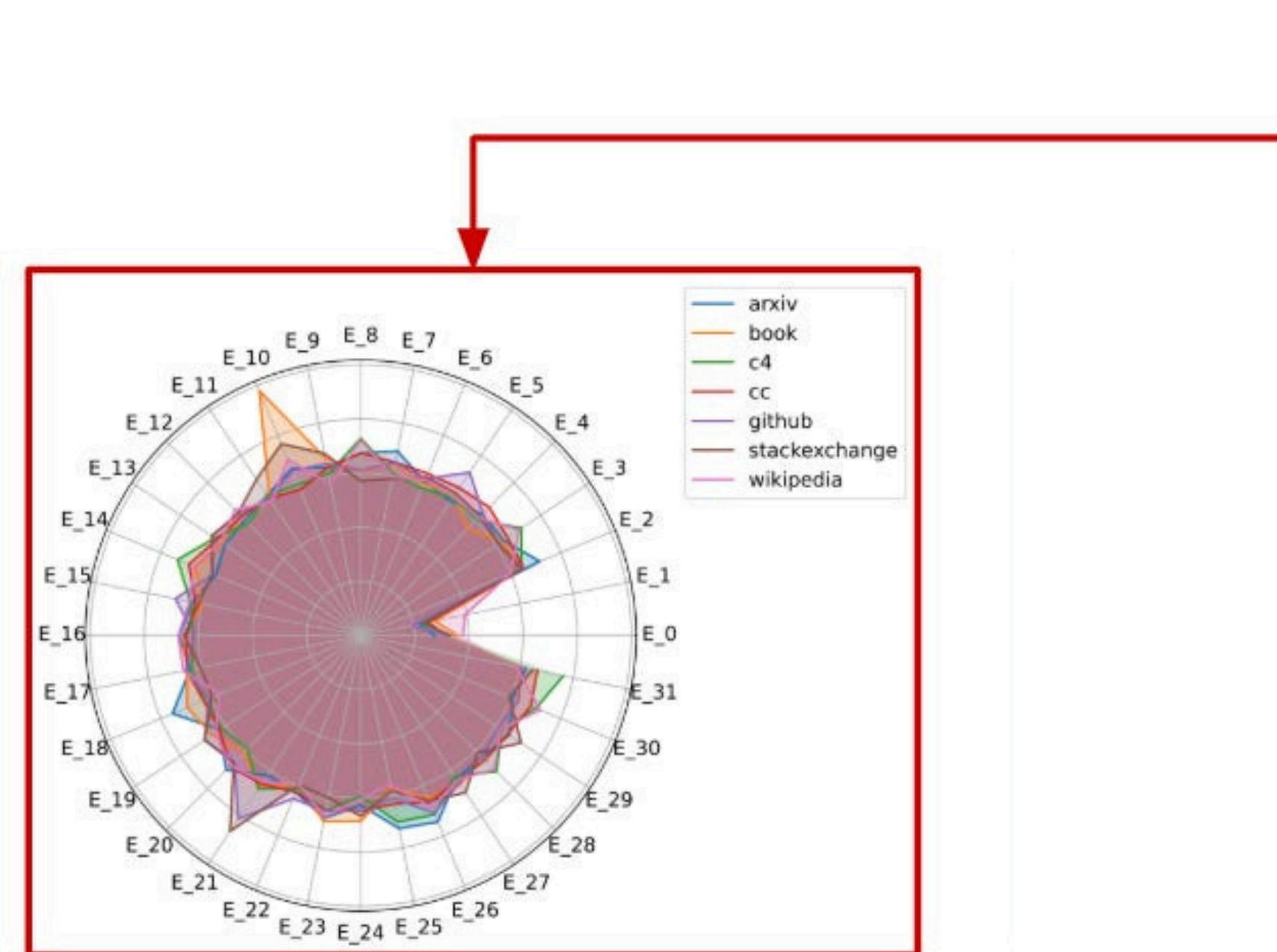
Wangchunshu Zhou<sup>3</sup> Yang You<sup>1</sup>

<sup>1</sup>National University of Singapore

<sup>2</sup>University of Edinburgh

<sup>3</sup>ETH Zurich

No indications of expert specialization across domains of data



# #3 There are no domain expert

Fuzhao Xue<sup>1†</sup> Zian Zheng<sup>1</sup> Yao Fu<sup>2</sup> Jinjie Ni<sup>1</sup> Zangwei Zheng<sup>1</sup>

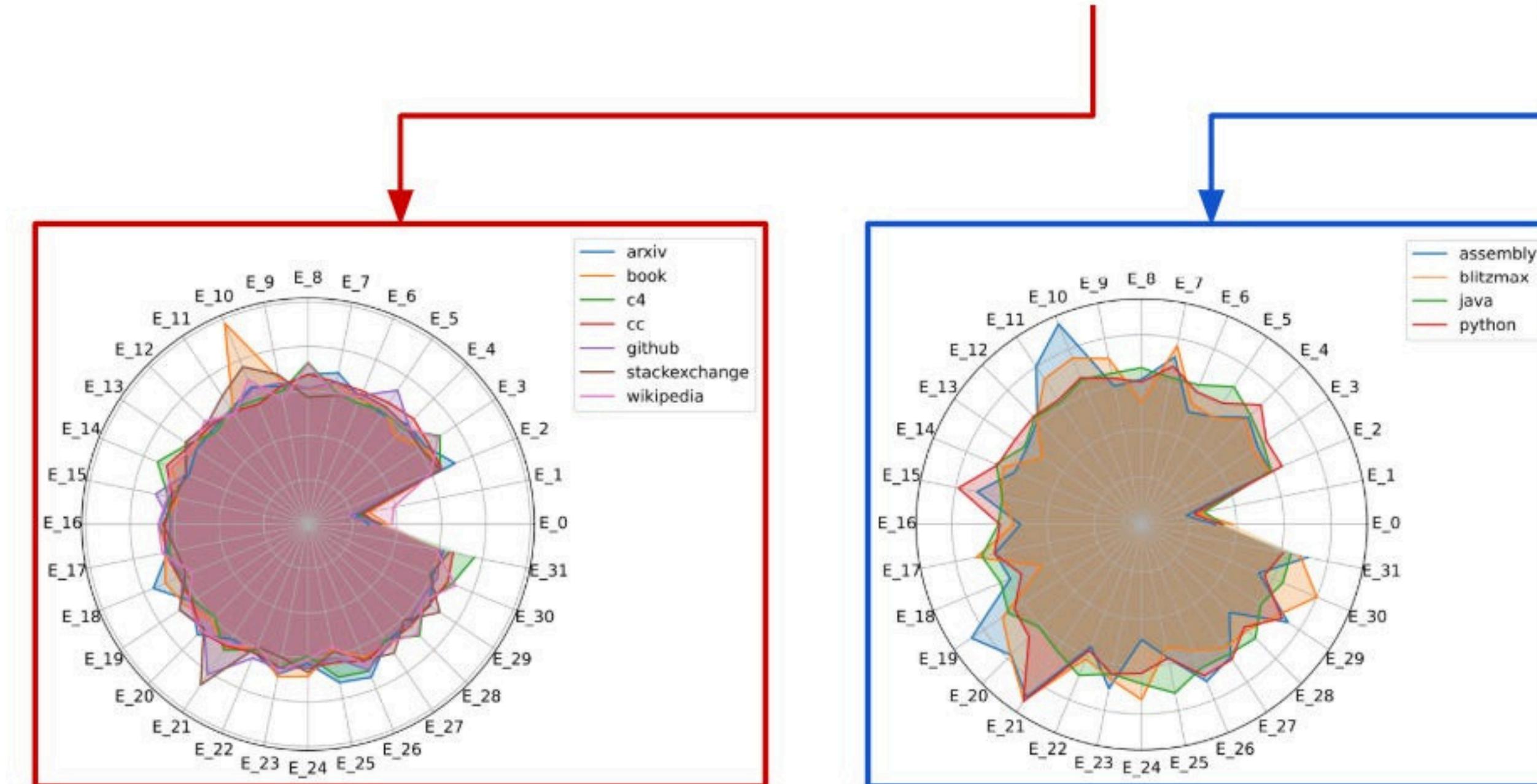
Wangchunshu Zhou<sup>3</sup> Yang You<sup>1</sup>

<sup>1</sup>National University of Singapore

<sup>2</sup>University of Edinburgh

<sup>3</sup>ETH Zurich

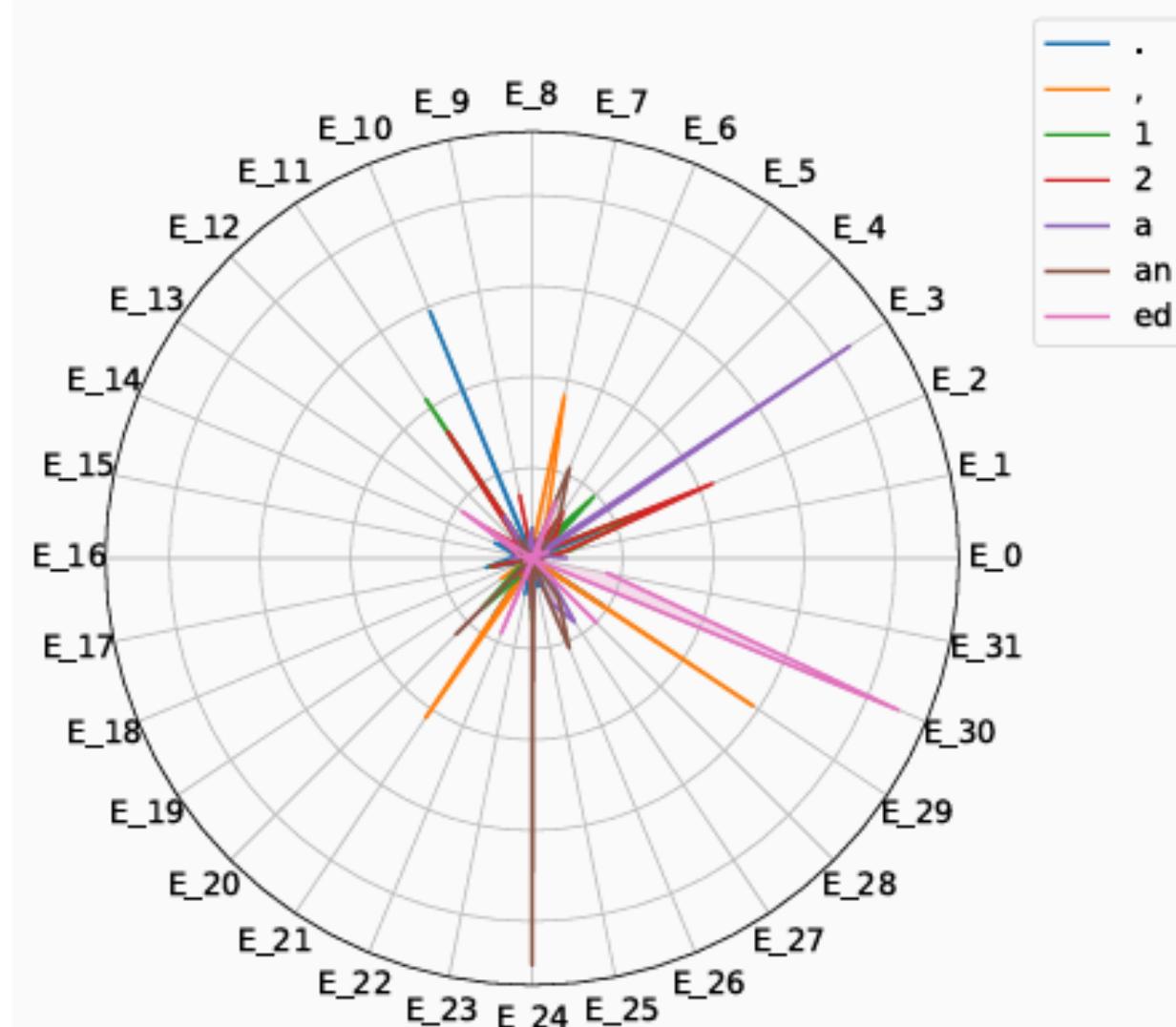
No indications of expert specialization across domains of data or coding languages.



<sup>1</sup>National University of Singapore

<sup>2</sup>University of Edinburgh

<sup>3</sup>ETH Zurich

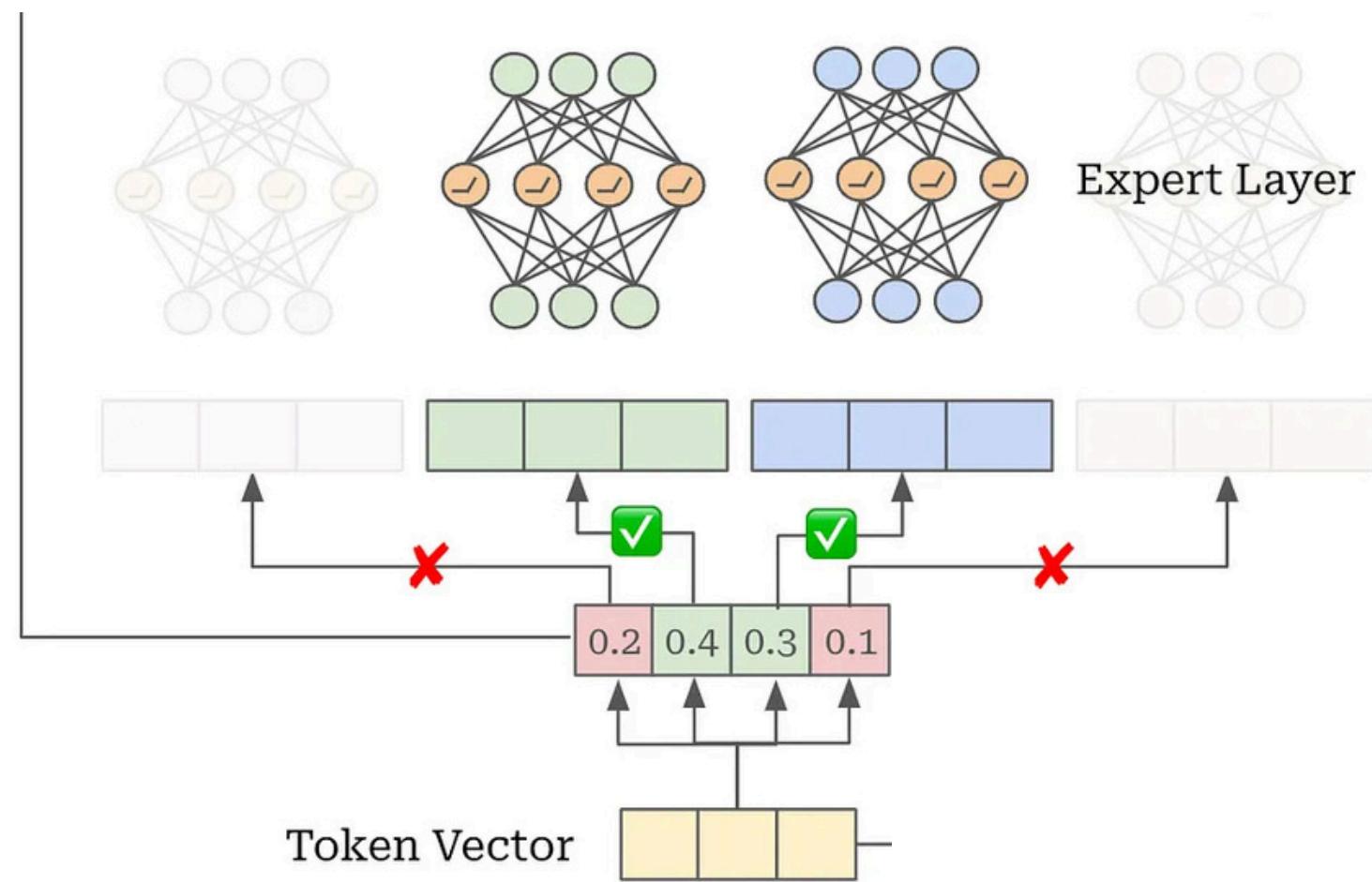


**Table 9: Top Tokens selected by each expert.**

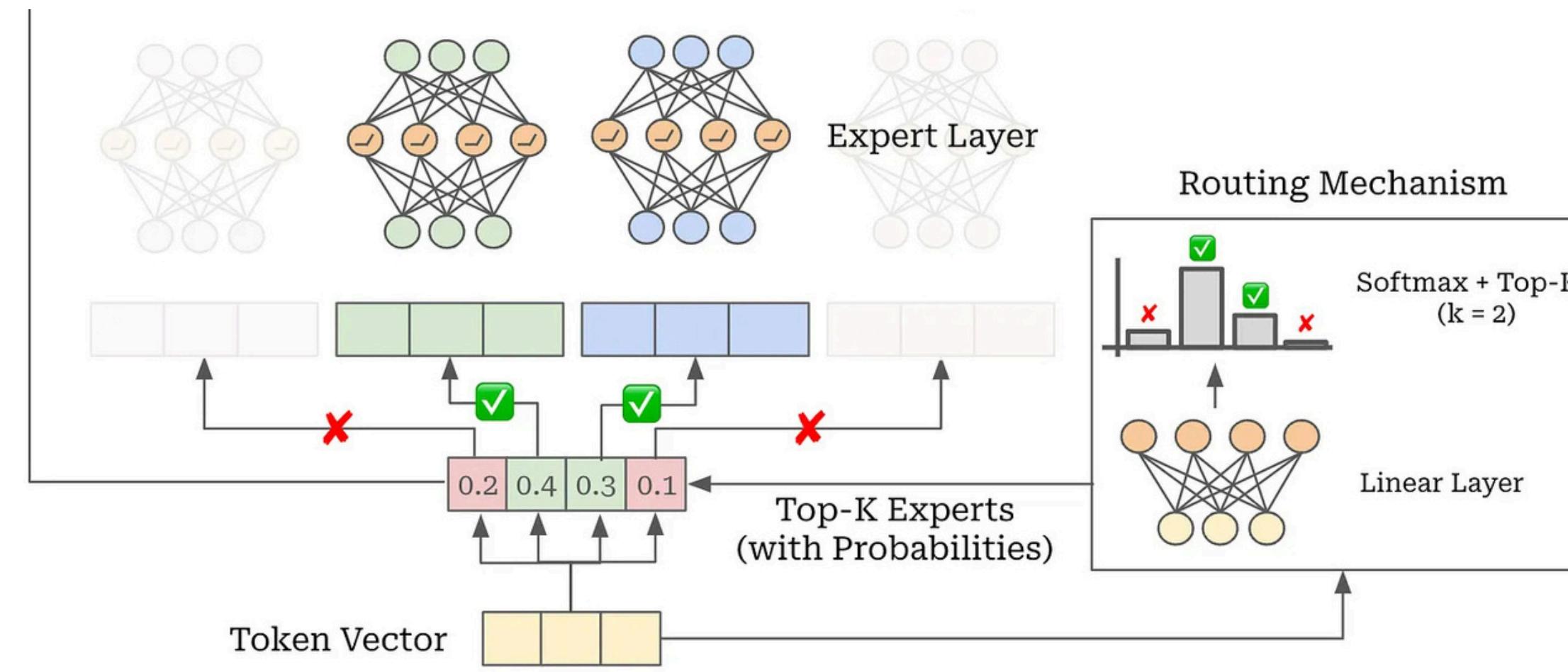
Expert ID	Top Tokens
0	\n, ‘, ’, s, -, \$, y, _, 2
1	\n, 1, 2, W, S, ., -, C, {
21	, , and, ., ., \n, =, \t, the, , n
30	}, ed, d, have, ing, , has, s, " had
31	to, can, s, of, ing, will, not, e, ed, would

Let's understand how its trained

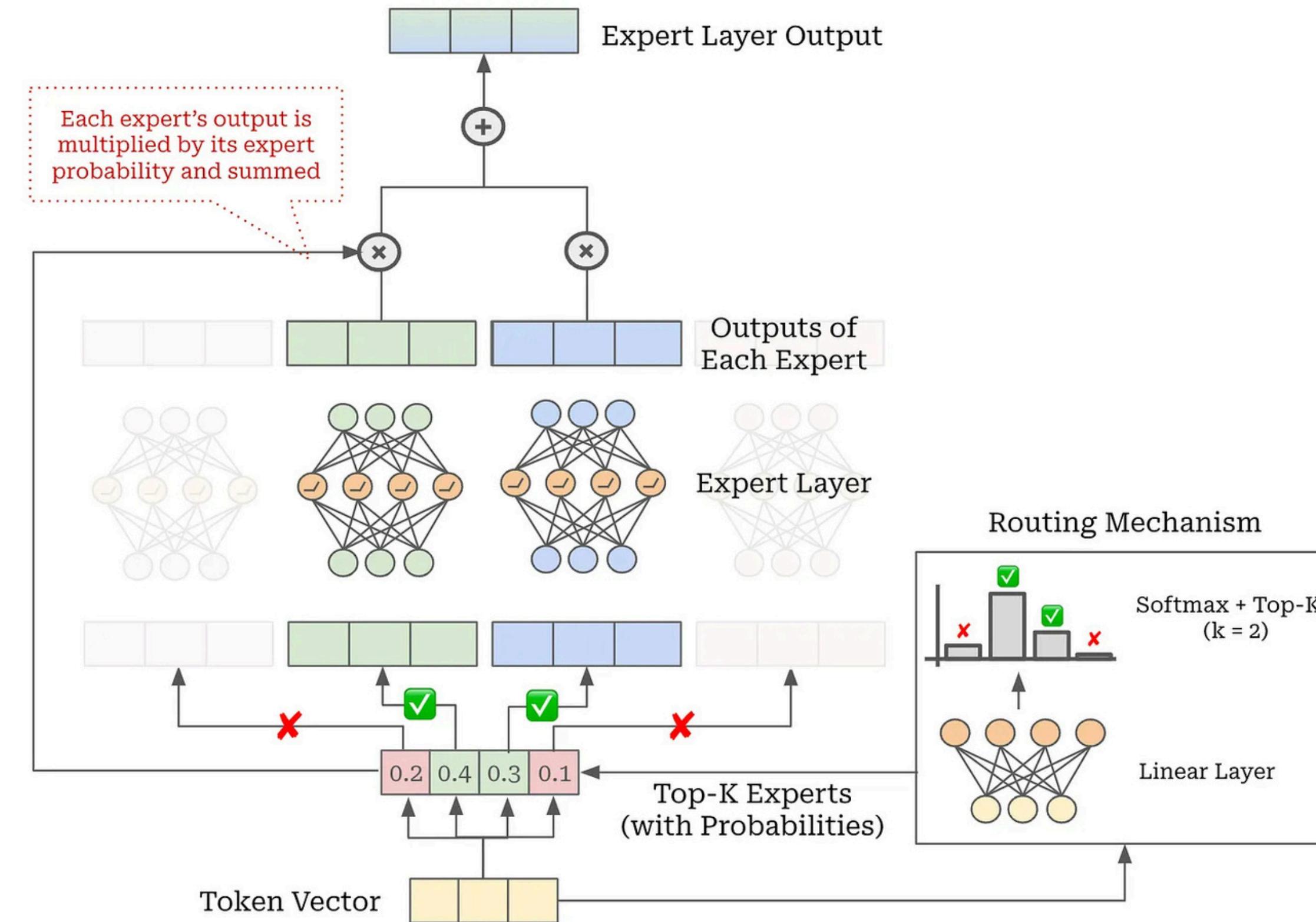
# MoE Layer



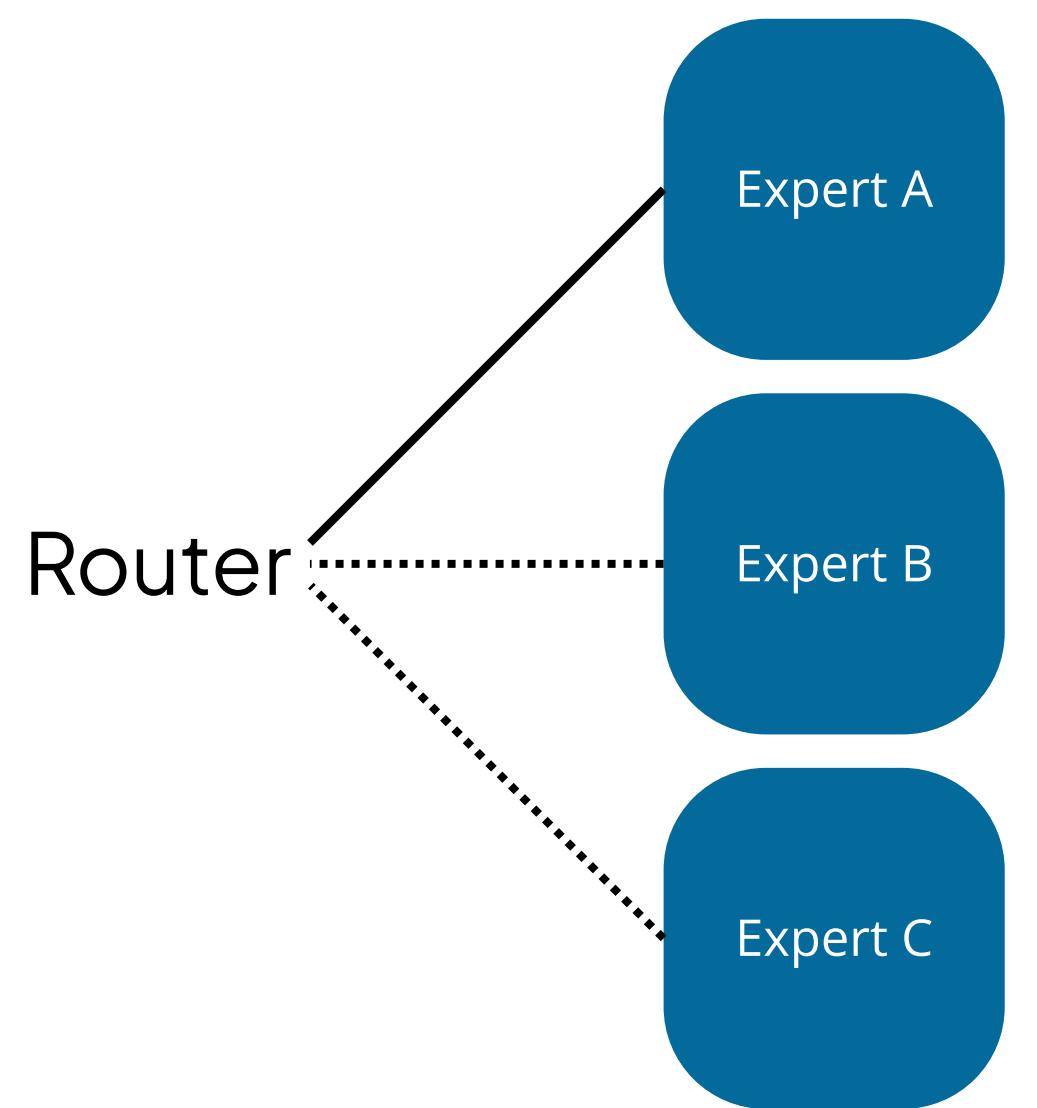
# MoE Layer



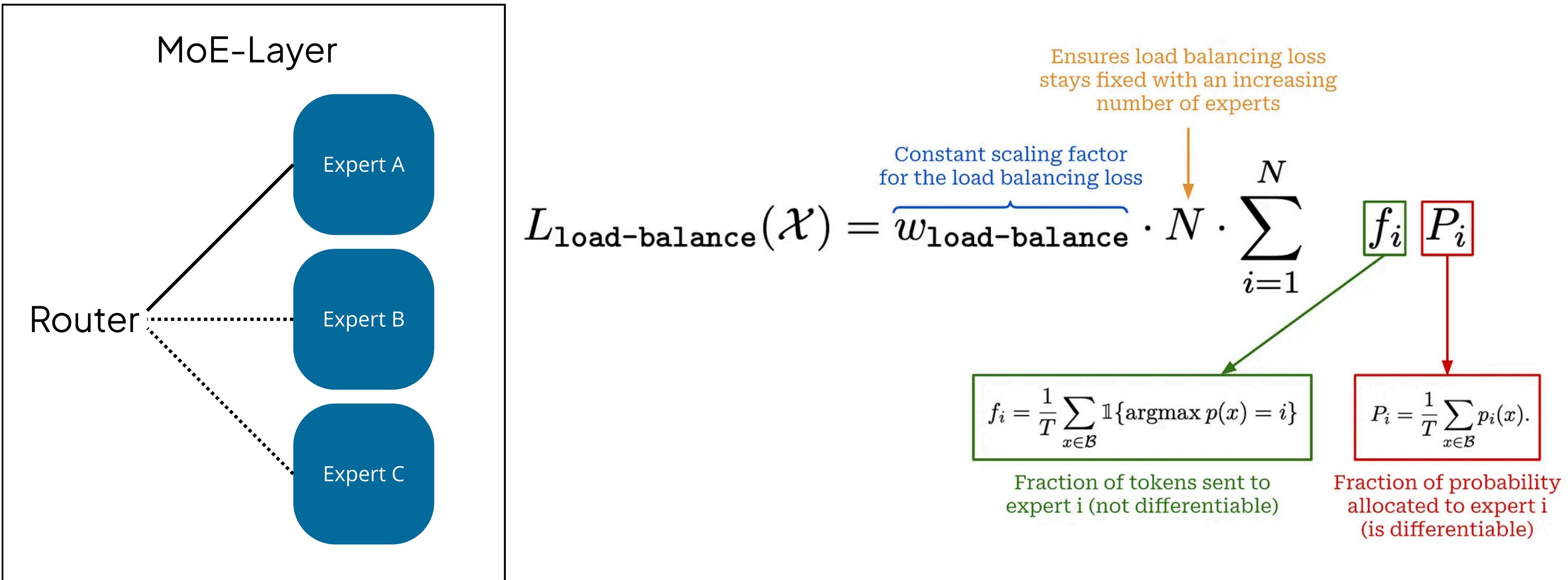
# MoE Layer



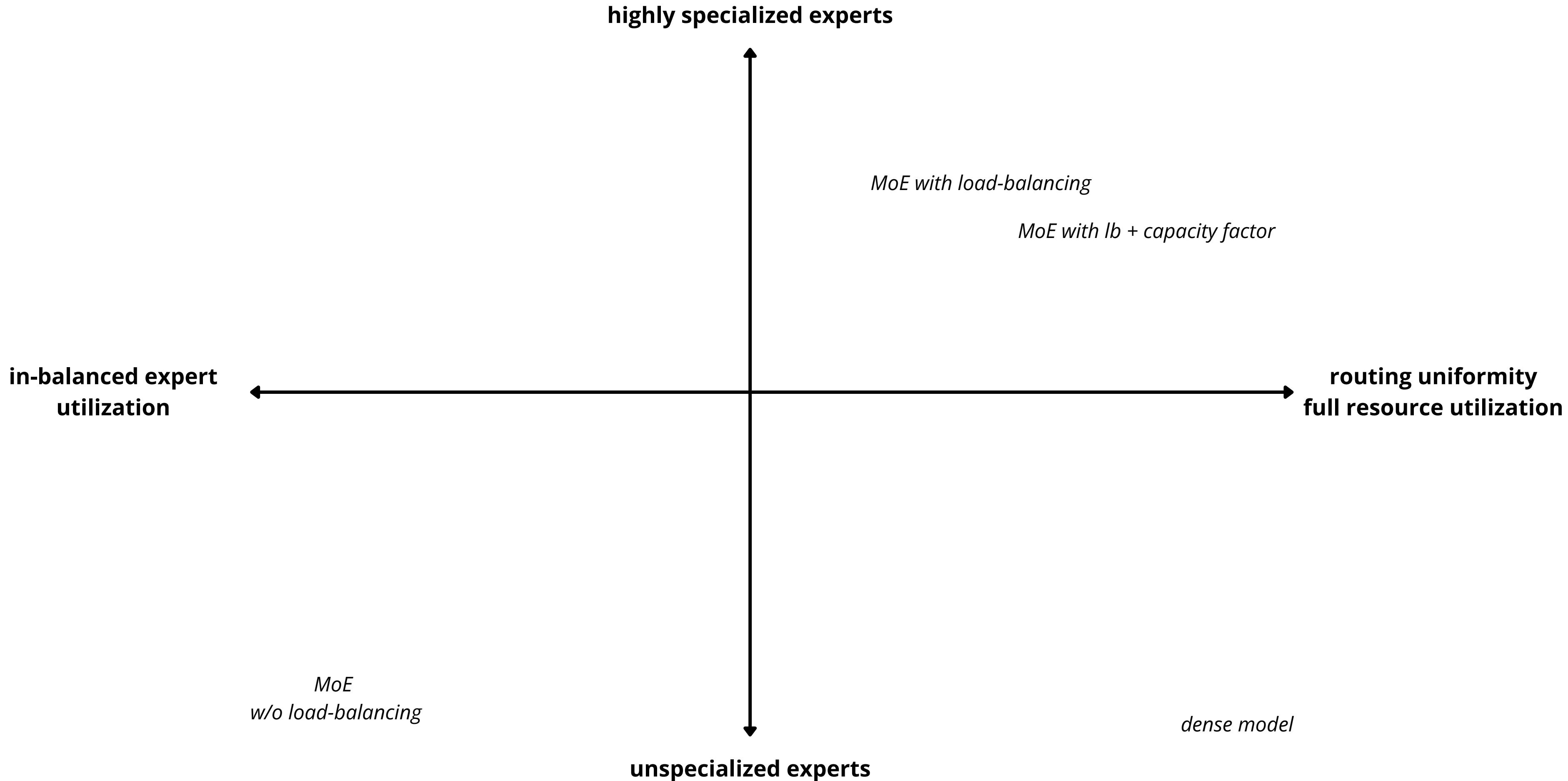
## MoE-Layer



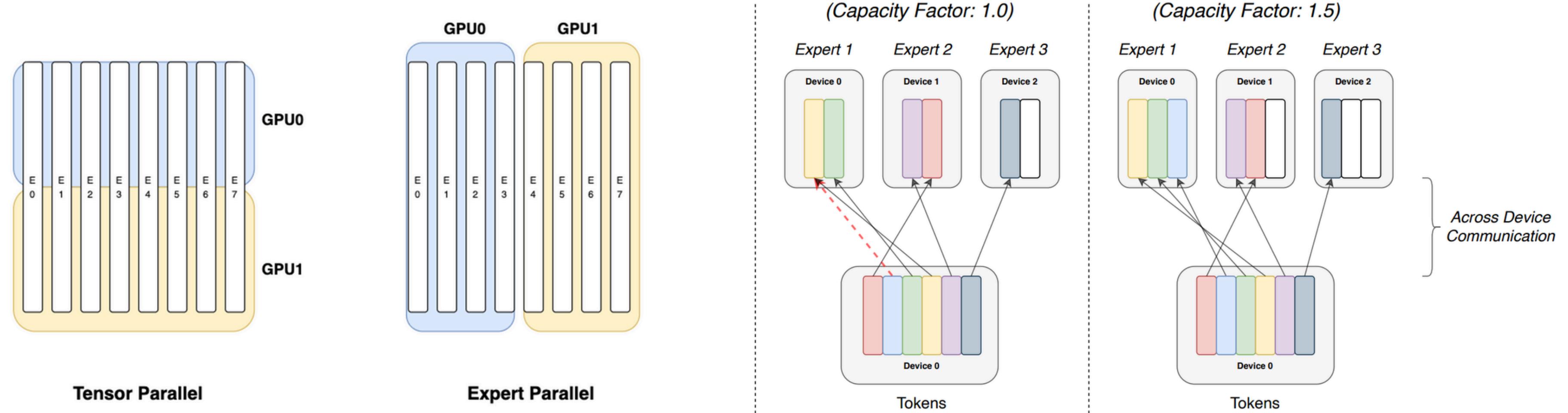
# Routing Collapse



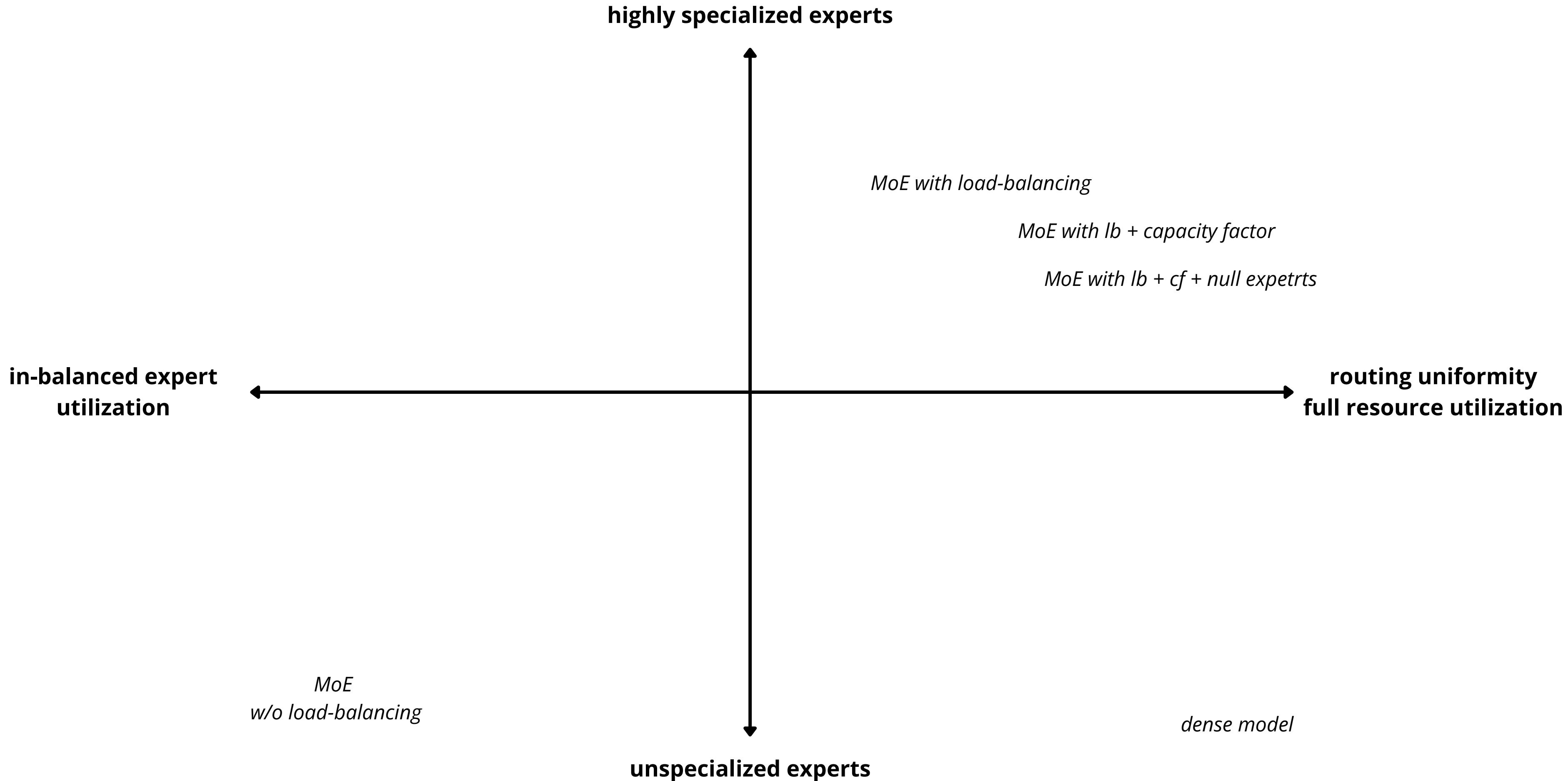
# Mixture of Experts trade-off



# Expert Capacity



# Mixture of Experts trade-off



# *AdaMoE*: Token-Adaptive Routing with Null Experts for Mixture-of-Experts Language Models

Zihao Zeng<sup>1\*</sup>, Yibo Miao<sup>1\*</sup>, Hongcheng Gao<sup>2</sup>, Hao Zhang<sup>3</sup>, Zhijie Deng<sup>1†</sup>

<sup>1</sup>Qing Yuan Research Institute, SEIEE, Shanghai Jiao Tong University

<sup>2</sup>University of Chinese Academy of Sciences <sup>3</sup>University of California, San Diego

{zengzihao, miaoyibo, zhijied}@sjtu.edu.cn

gaohongcheng23@mails.ucas.ac.cn, haozhang@ucsd.edu

Published as a conference paper at ICLR 2025

## MOE++: ACCELERATING MIXTURE-OF-EXPERTS METHODS WITH ZERO-COMPUTATION EXPERTS

Peng Jin<sup>1,2</sup>, Bo Zhu<sup>3</sup>, Li Yuan<sup>1,2,4</sup>✉, Shuicheng Yan<sup>3,5</sup>✉

<sup>1</sup>Pengcheng Laboratory

<sup>2</sup>School of Electronic and Computer Engineering, Shenzhen Graduate School, Peking University

<sup>3</sup>Kunlun 2050 Research & Skywork AI

<sup>4</sup>Rabbitpre Intelligence <sup>5</sup>National University of Singapore

jp21@stu.pku.edu.cn, yuanli-ece@pku.edu.cn

Code: <https://github.com/SkyworkAI/MoE-plus-plus>

## Token-Adaptive Router in *AdaMoE*

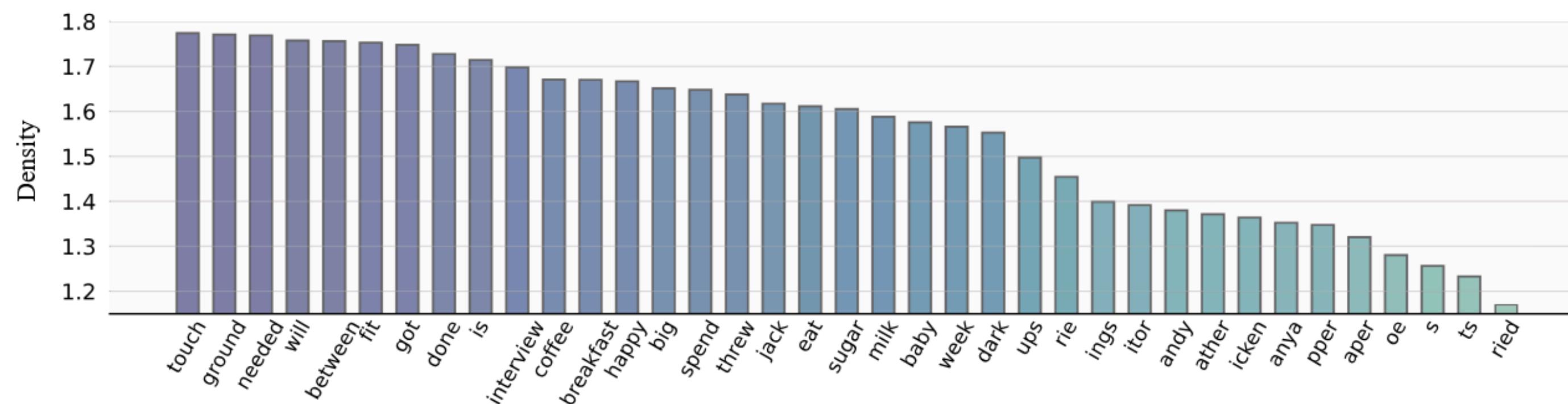
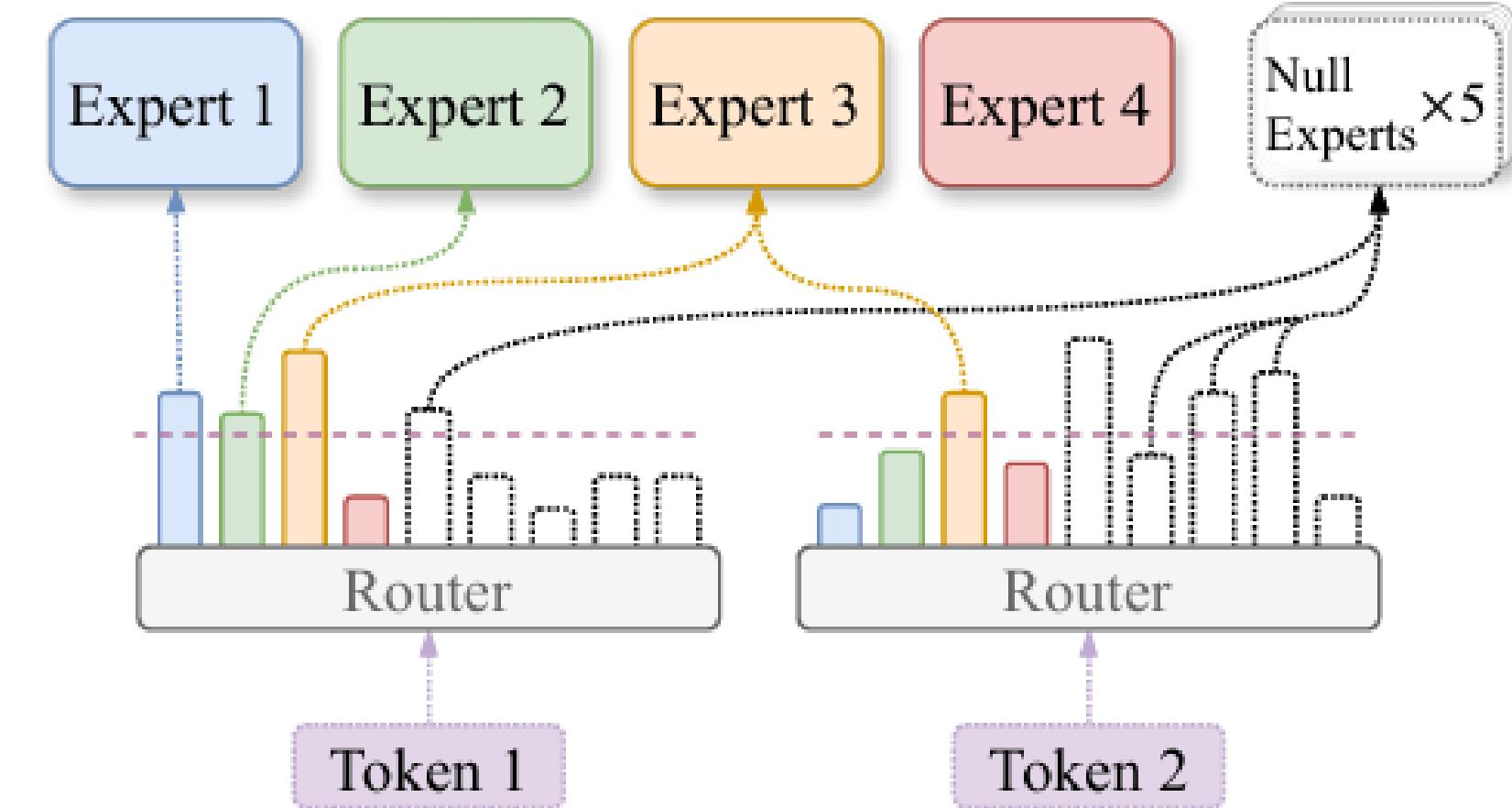
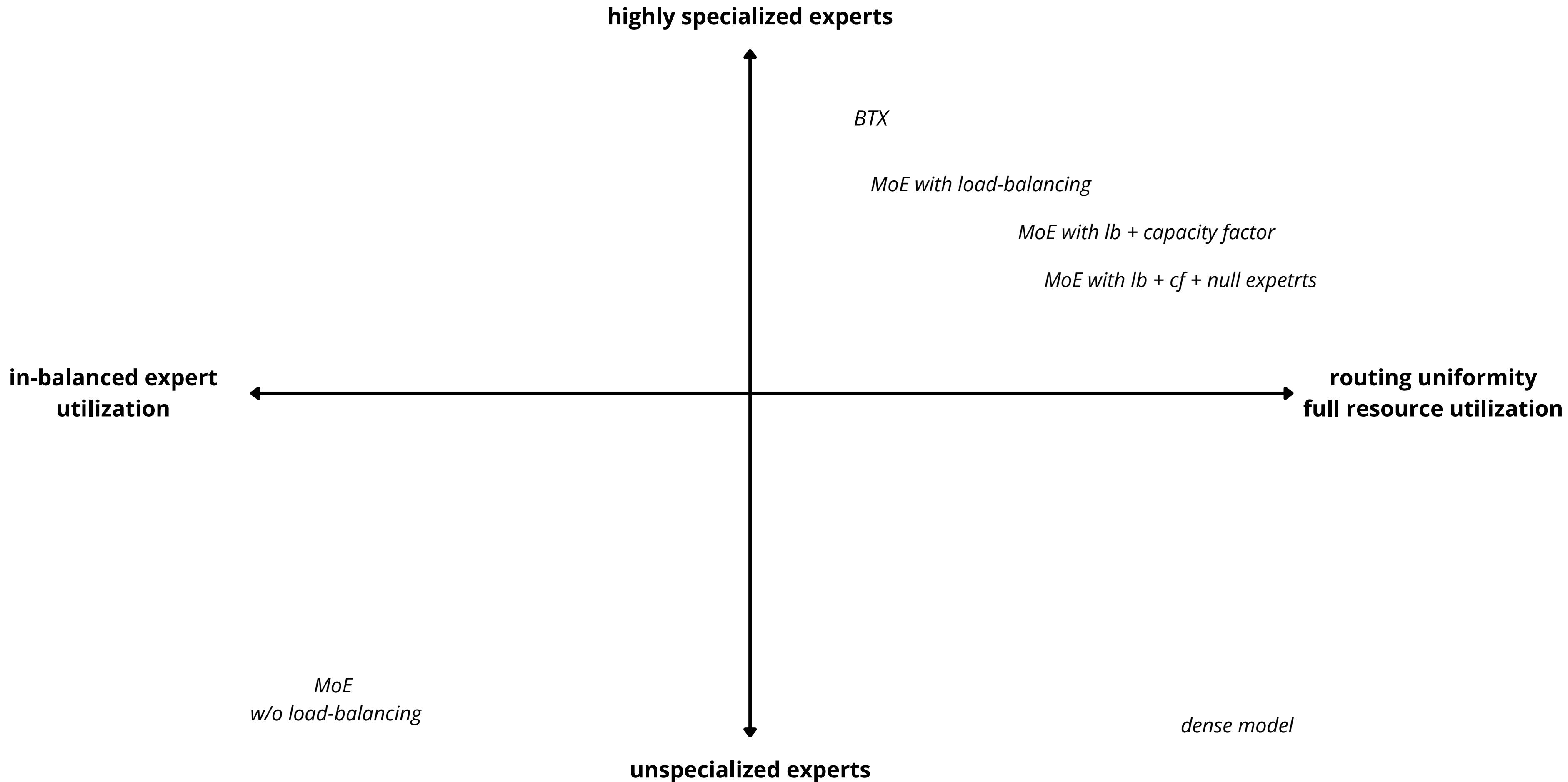


Figure 5: The visualization of the number of FFN experts activated per token at the token level.

# Mixture of Experts trade-off



# Branch-Train-MiX: Mixing Expert LLMs into a Mixture-of-Experts LLM

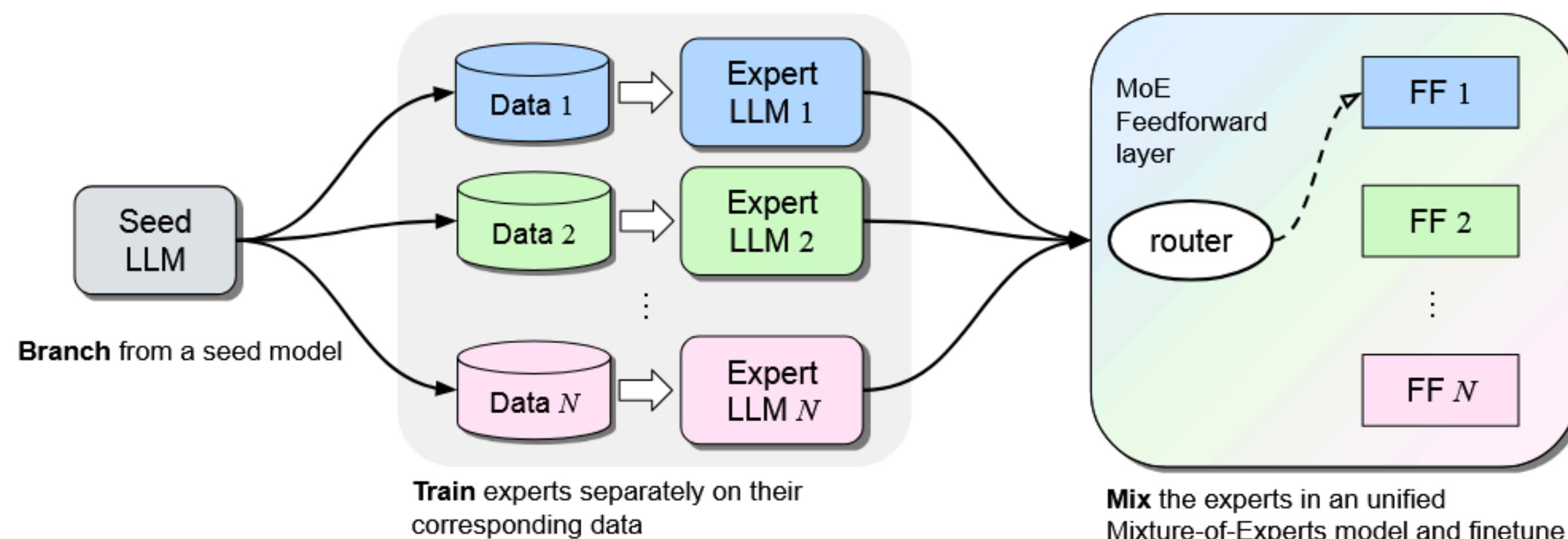
Sainbayar Sukhbaatar, Olga Golovneva, Vasu Sharma, Hu Xu, Xi Victoria Lin, Baptiste Rozière, Jacob Kahn, Daniel Li, Wen-tau Yih, Jason Weston, Xian Li

FAIR at Meta

We investigate efficient methods for training Large Language Models (LLMs) to possess capabilities in multiple specialized domains, such as coding, math reasoning and world knowledge. Our method, named Branch-Train-MiX (BTX), starts from a seed model, which is branched to train experts in embarrassingly parallel fashion with high throughput and reduced communication cost. After individual experts are asynchronously trained, BTX brings together their feedforward parameters as experts in Mixture-of-Expert (MoE) layers and averages the remaining parameters, followed by an MoE-finetuning stage to learn token-level routing. BTX generalizes two special cases: the Branch-Train-Merge method and sparse upcycling, which omit the MoE layer. By combining the strengths of both approaches, BTX achieves the best performance.

Date: March 13, 2024

Correspondence: {sainbar,xian}



**Figure 1 The Branch-Train-MiX (BTX) method** has three steps: **1) branch** from a pretrained seed LLM by making multiple copies of it; **2) train** those copies separately on different subsets of data to obtain expert LLMs; **3) mix** those expert LLMs by combining them into a single LLM using mixture-of-experts feedforward (FF) layers, and finetuning the overall unified model.

# Self-Selecting Experts

# Self-Selecting Experts

## Autonomy-of-Experts Models

Ang Lv<sup>1</sup> Ruobing Xie<sup>2</sup> Yining Qian<sup>3</sup> Songhao Wu<sup>1</sup> Xingwu Sun<sup>2,4</sup> Zhanhui Kang<sup>2</sup> Di Wang<sup>2</sup> Rui Yan<sup>1,5,6</sup>

## Mixture-of-Experts with Expert Choice Routing

Yanqi Zhou, Tao Lei, Hanxiao Liu, Nan Du, Yanping Huang, Vincent Zhao, Andrew Dai, Zhifeng Chen, Quoc Le, and James Laudon

Google, Mountain View, CA, USA  
{yanqiz, taole, hanxiao1, dunan, huangyp, vzhao, adai, zhifengc, qvl, jlaudon}@google.com

## SELF-MOE: TOWARDS COMPOSITIONAL LARGE LANGUAGE MODELS WITH SELF-SPECIALIZED EXPERTS

Junmo Kang\*  
Georgia Tech

Leonid Karlinsky  
MIT-IBM Watson AI Lab

Hongyin Luo  
MIT

Zhen Wang  
UCSD

Jacob Hansen  
MIT

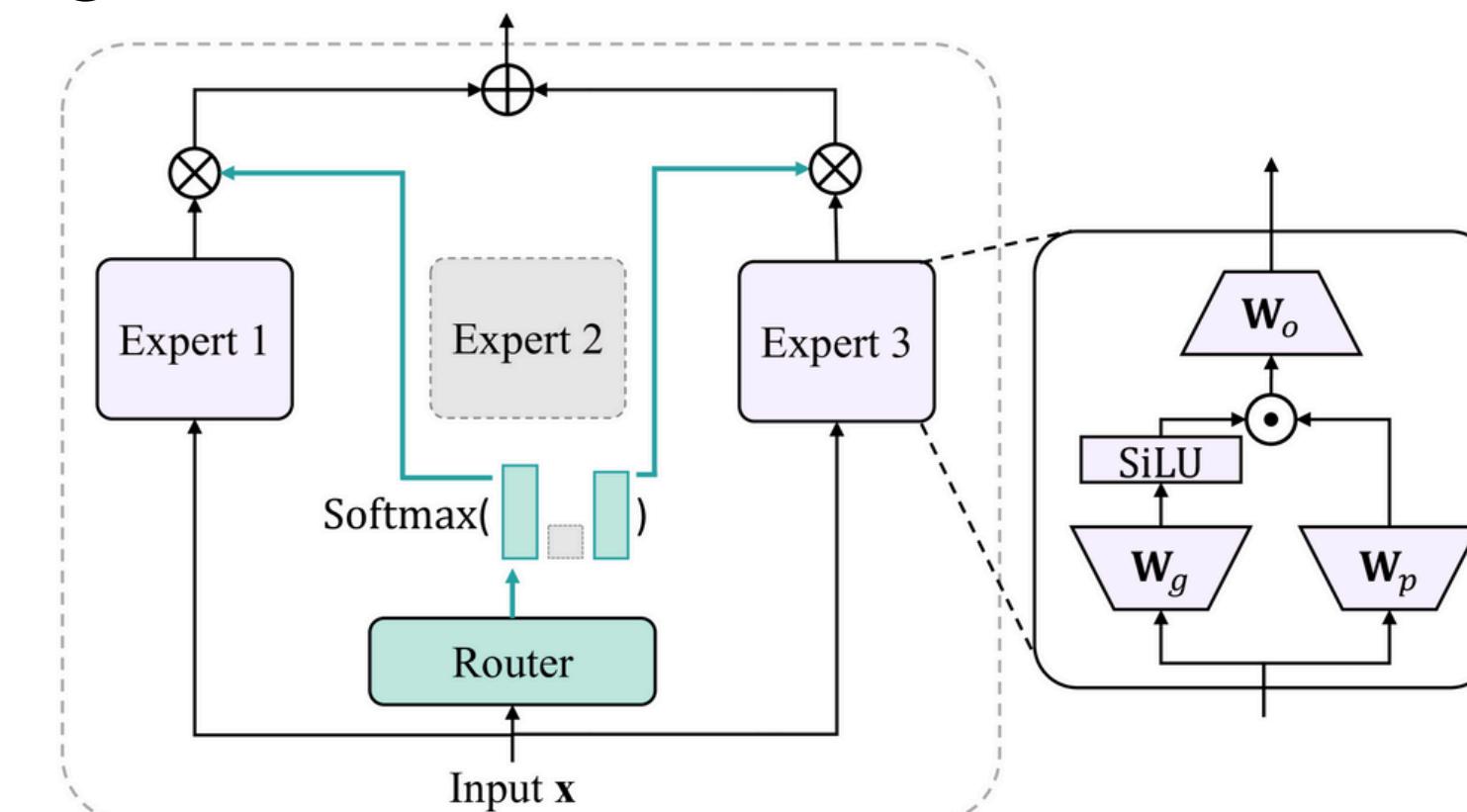
James Glass  
MIT

David Cox  
MIT-IBM Watson AI Lab

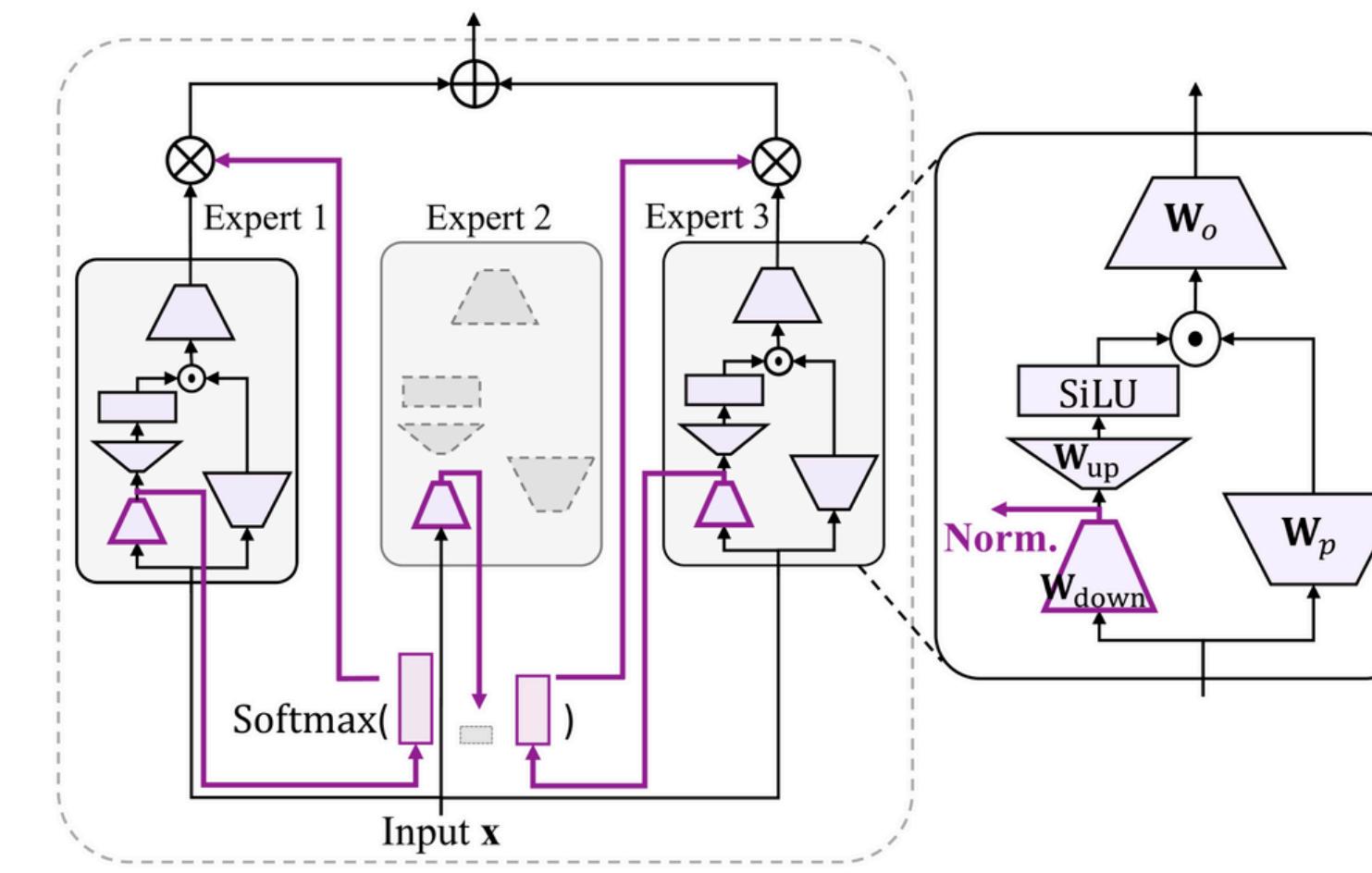
Rameswar Panda  
MIT-IBM Watson AI Lab

Rogerio Feris  
MIT-IBM Watson AI Lab

Alan Ritter  
Georgia Tech



(a) Mixture-of-Experts

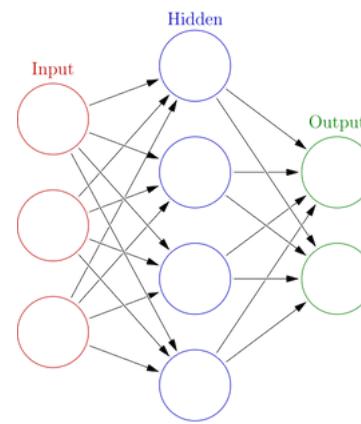


(b) Autonomy-of-Experts

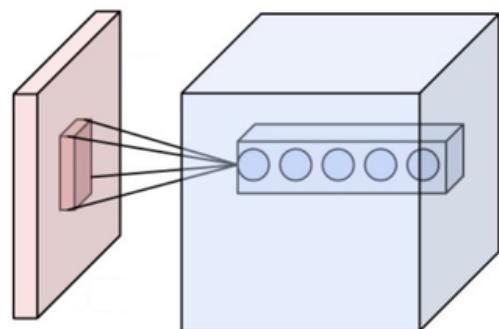
The Bigger Picture

# Dense Net → System of Modular Sparse Nets

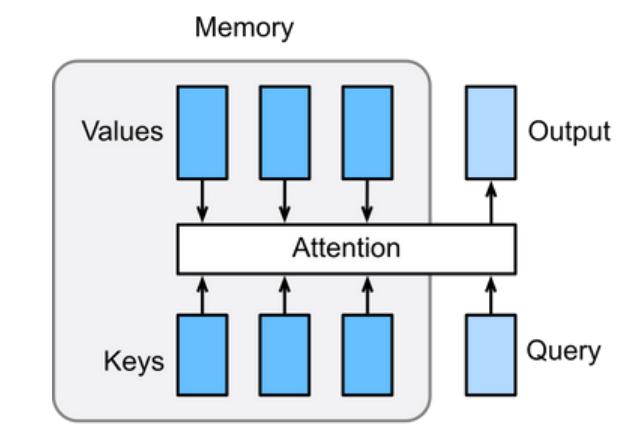
Fully Connected Networks



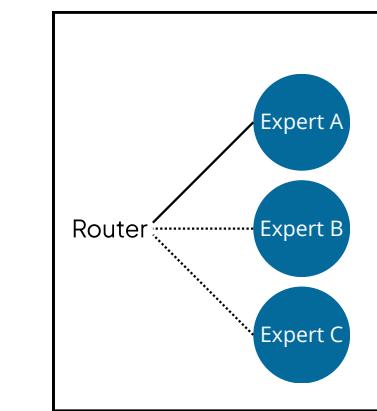
CNNs



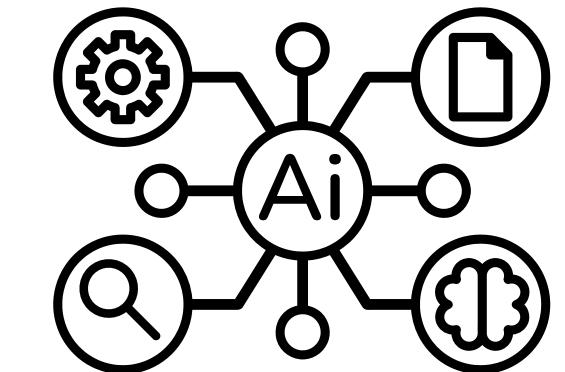
Transformer



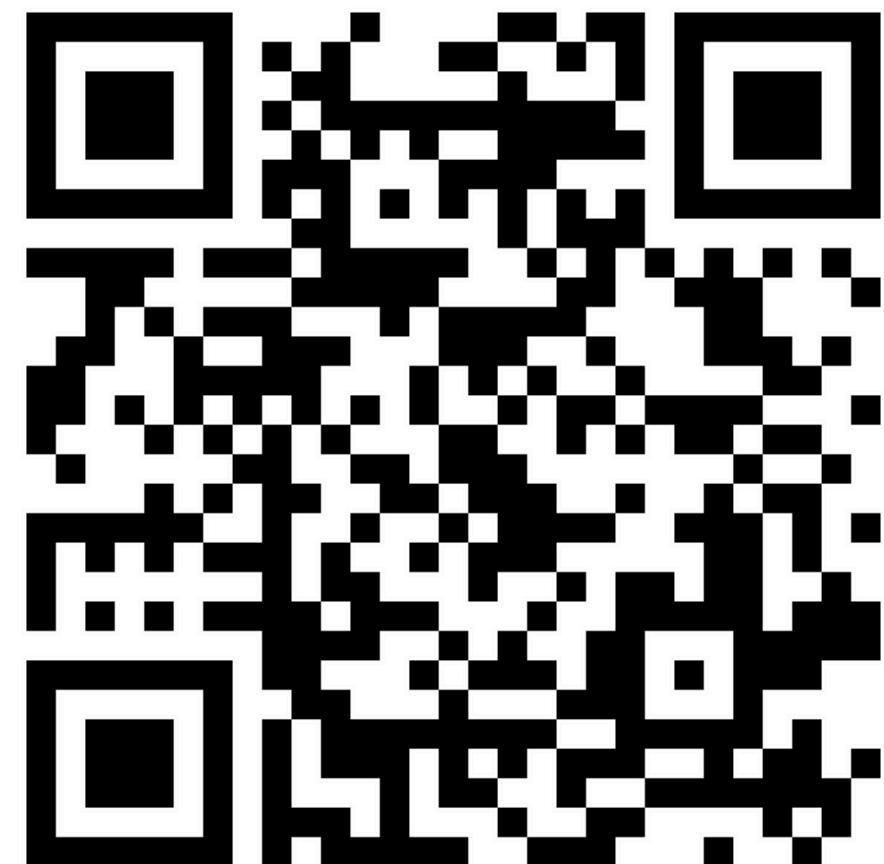
MoE



System of Modular Experts



# Thank you



Link To The Slides