

Self-MoE: Towards Compositional Large Language Models with Self-Specialized Experts

Junmo Kang¹
Jacob Hansen³

Leonid Karlinsky²
James Glass³
Rogerio Feris²

Hongyin Luo³
David Cox²
Alan Ritter¹

Zhen Wang⁴
Rameswar Panda²



Towards Compositional Large Language Models

Rapid advancement of LLMs

- as a generalist
- relying on substantial resources - data, compute, and parameters

Towards Compositional Large Language Models

Rapid advancement of LLMs

- as a generalist
- relying on substantial resources - data, compute, and parameters
- a monolithic structure

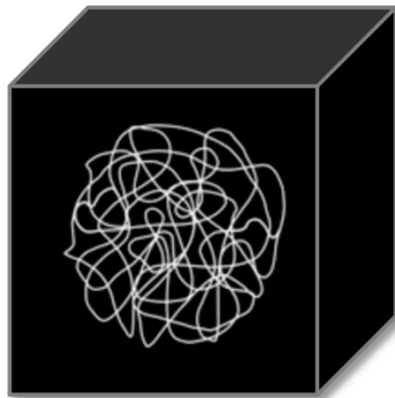
Monolithic

Generalist

Inefficient

Forgetting

Black-box



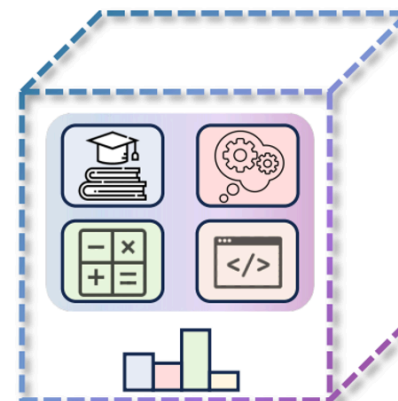
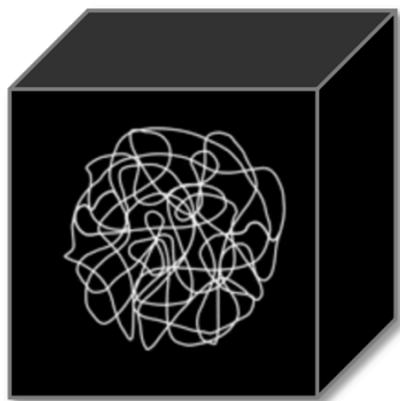
Towards Compositional Large Language Models

Rapid advancement of LLMs

- as a generalist
- relying on substantial resources - data, compute, and parameters
- a monolithic structure

Monolithic

Generalist
Inefficient
Forgetting
Black-box

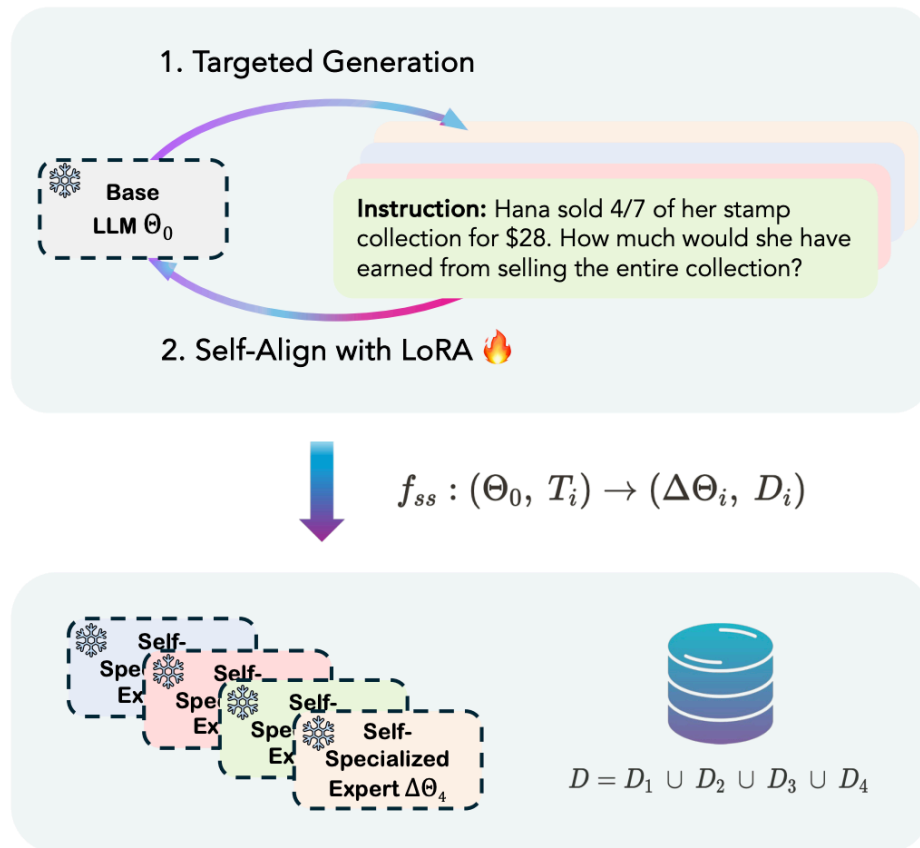


Compositional

Specialists
Efficient
Adaptable
Interpretable

Overview of Self-MoE

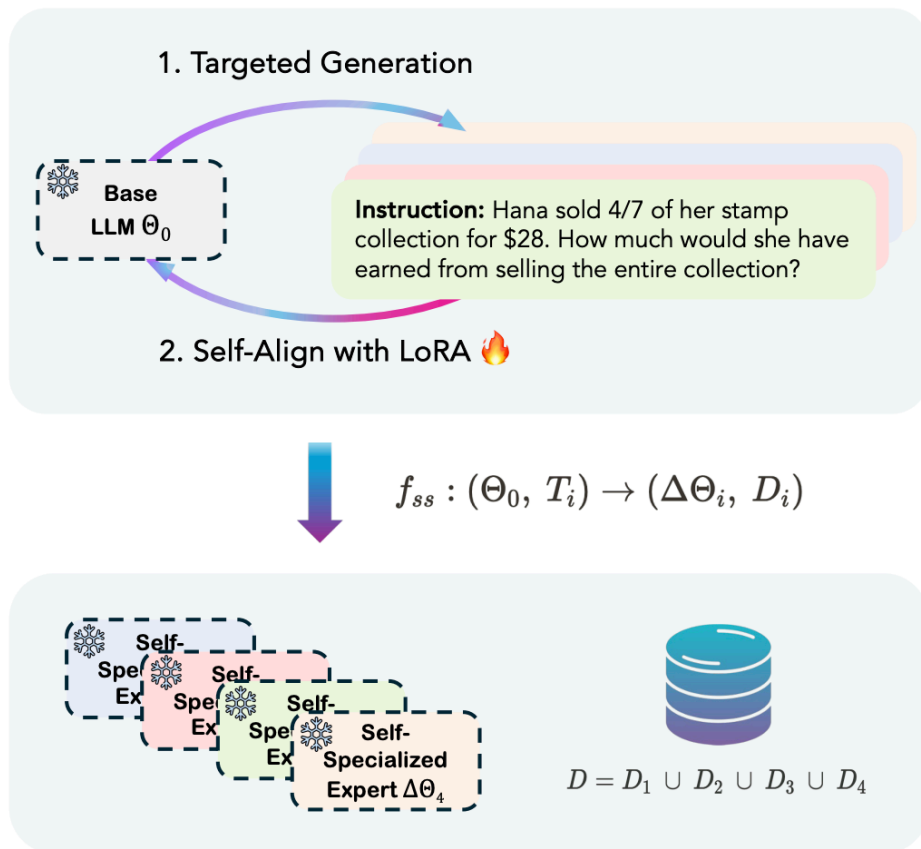
Self-Specialization



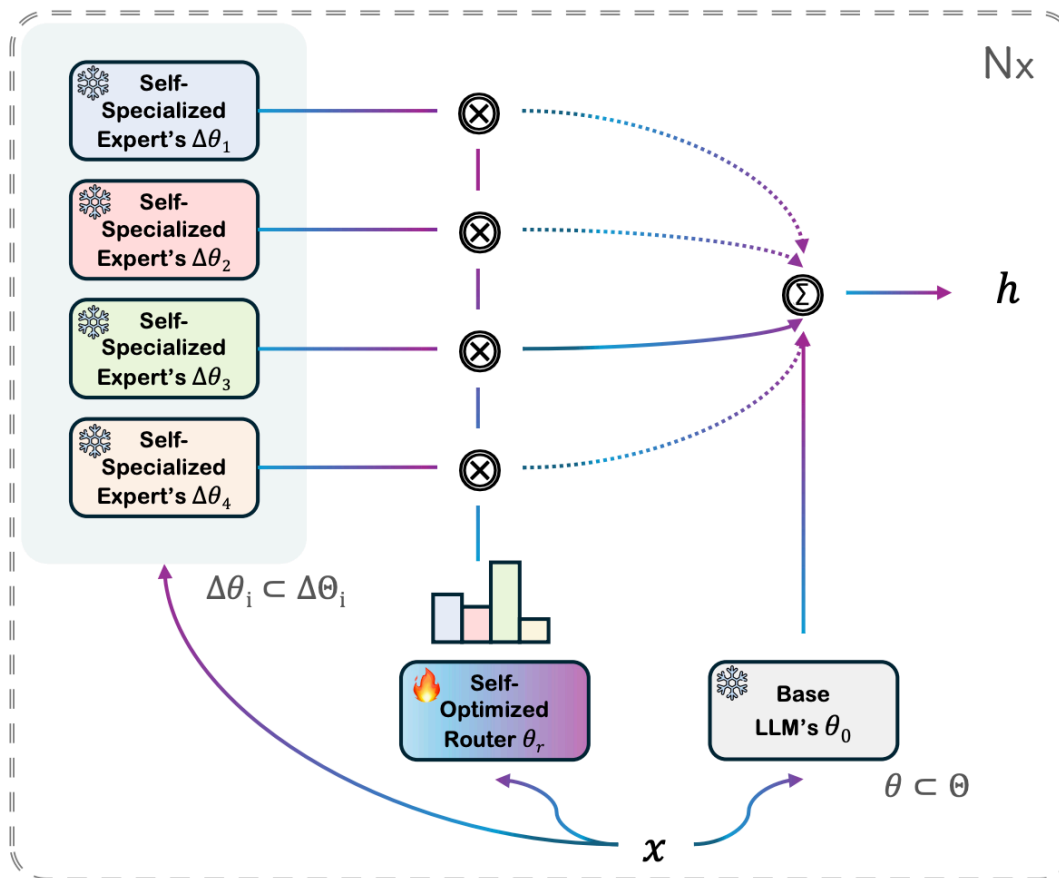
Approach

Overview of Self-MoE

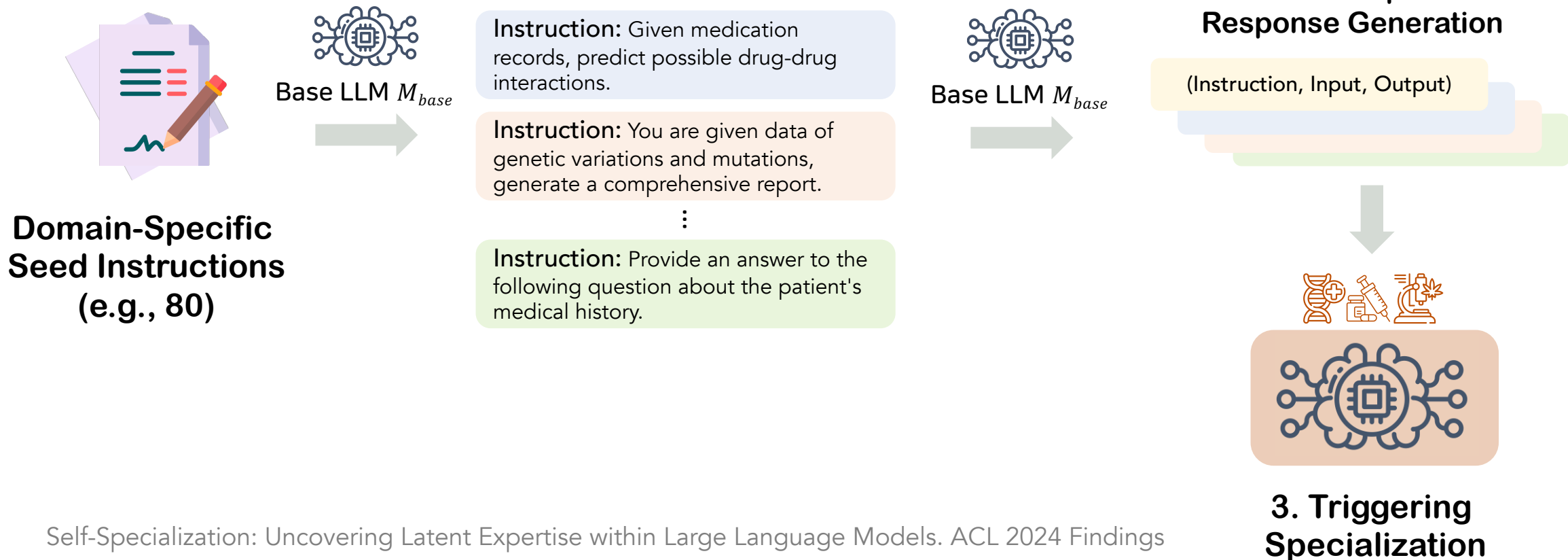
Self-Specialization



MiXSE (MiXture of Self-Specialized Experts)



Self-Specialization for Uncovering Domain Expertise

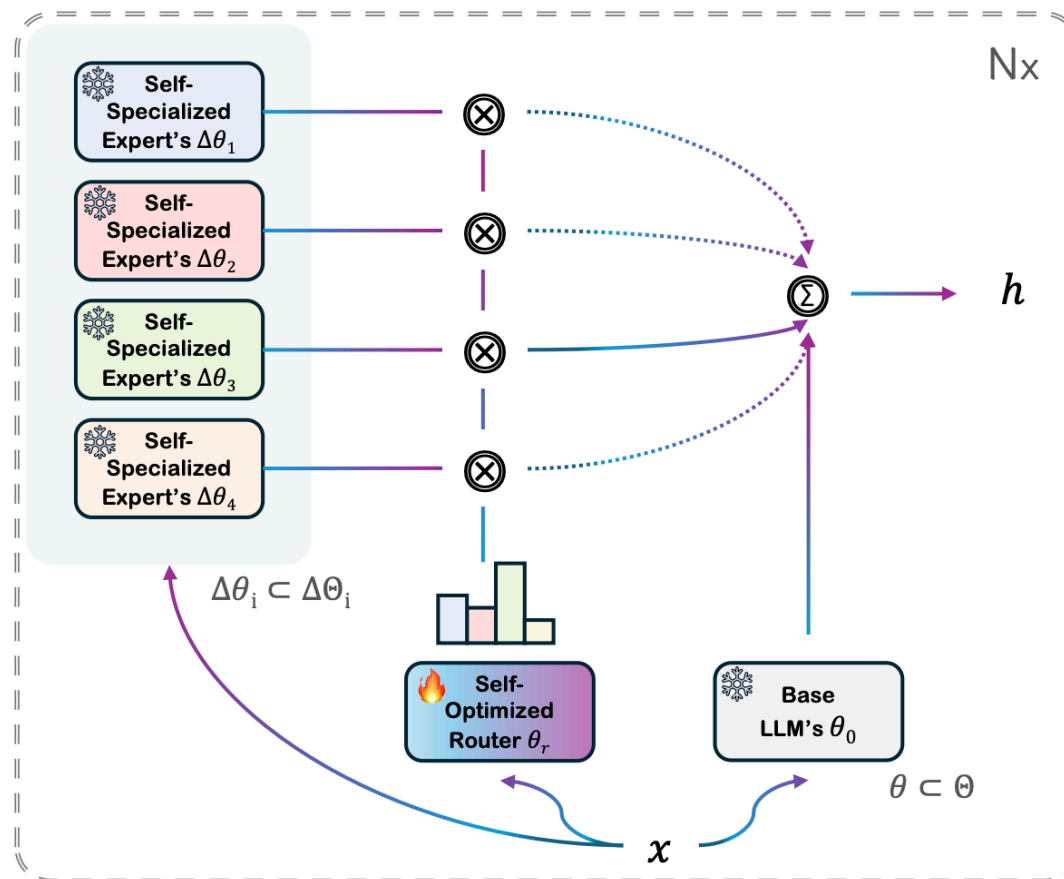


Overview of Self-MoE

$$\theta' = \theta + \Delta\theta = \theta + \theta_B \theta_A$$

$h = \theta x + \theta_B \theta_A x$ for a single expert LoRA

MiXSE (MiXture of Self-Specialized Experts)



Overview of Self-MoE

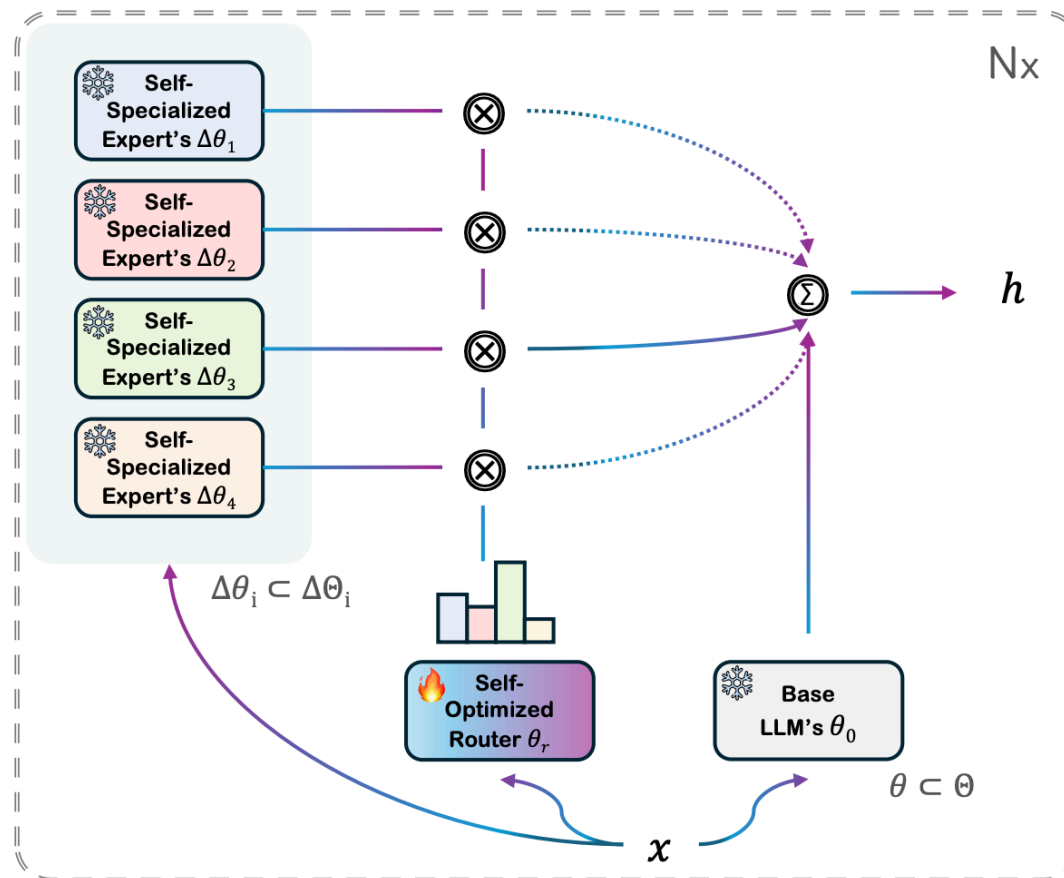
$$\theta' = \theta + \Delta\theta = \theta + \theta_B \theta_A$$

$$h = \theta x + \theta_B \theta_A x \text{ for a single expert LoRA}$$

Each Expert 1, 2, 3, 4

$$\begin{aligned} h &= \theta x + \sum \alpha_i \Delta\theta_i x \\ &= \theta x + \sum \alpha_i \theta_{Bi} \theta_{Ai} x \end{aligned}$$

MiXSE (MiXture of Self-Specialized Experts)



Approach

Overview of Self-MoE

$$\theta' = \theta + \Delta\theta = \theta + \theta_B \theta_A$$

$$h = \theta x + \theta_B \theta_A x \text{ for a single expert LoRA}$$

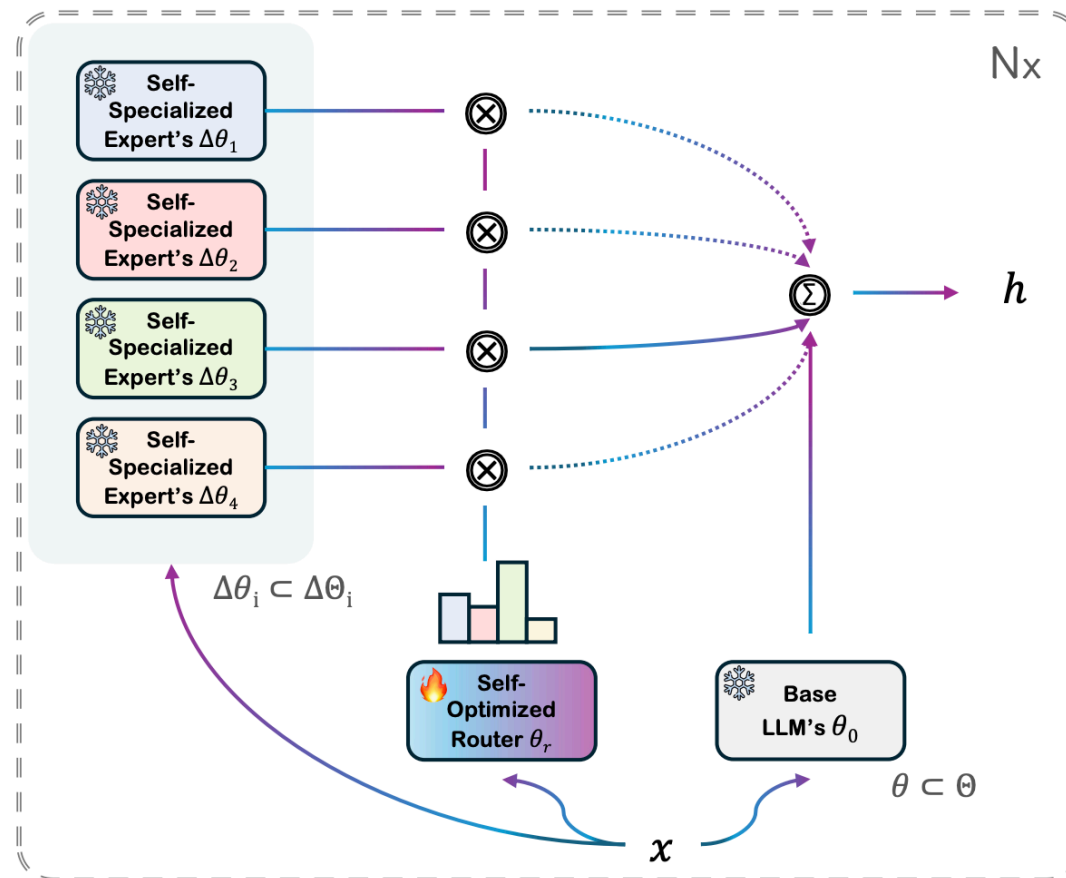
Each Expert 1, 2, 3, 4

$$\begin{aligned} h &= \theta x + \sum \alpha_i \Delta\theta_i x \\ &= \theta x + \sum \alpha_i \theta_{Bi} \theta_{Ai} x \end{aligned}$$

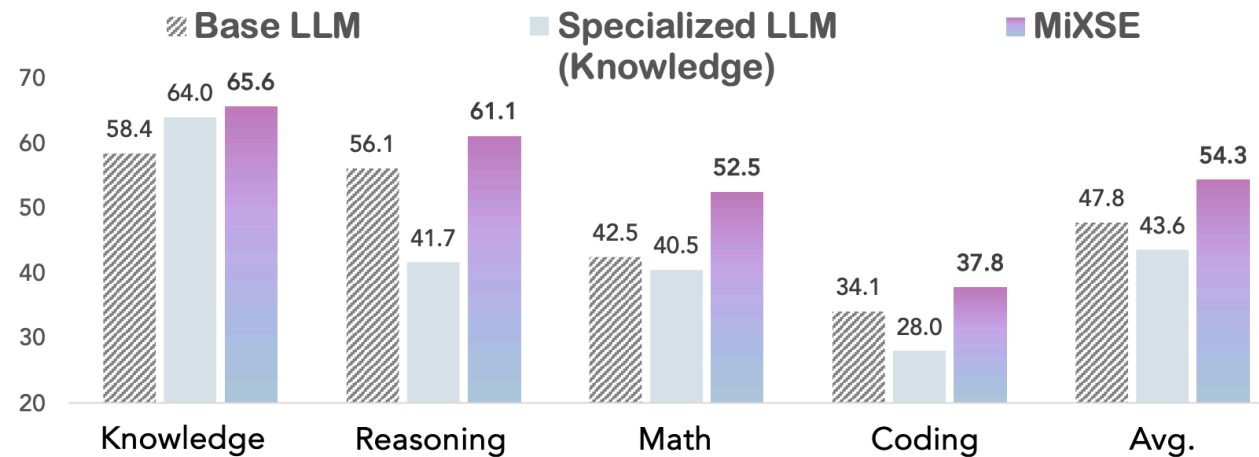
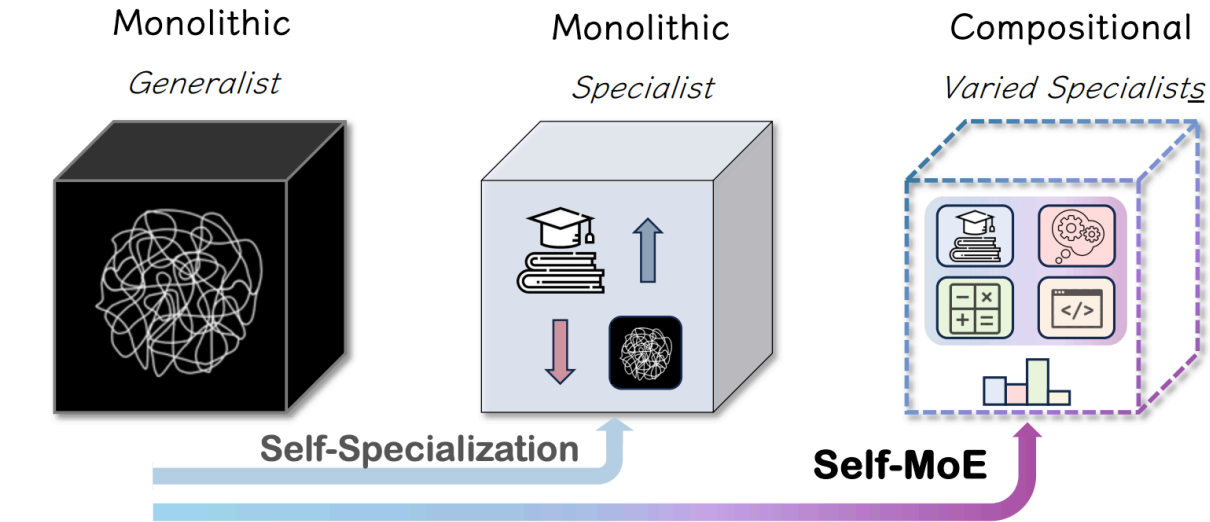
$$\alpha = \text{softmax}(\theta_r x) \text{ (topk)}$$

Routing Weights

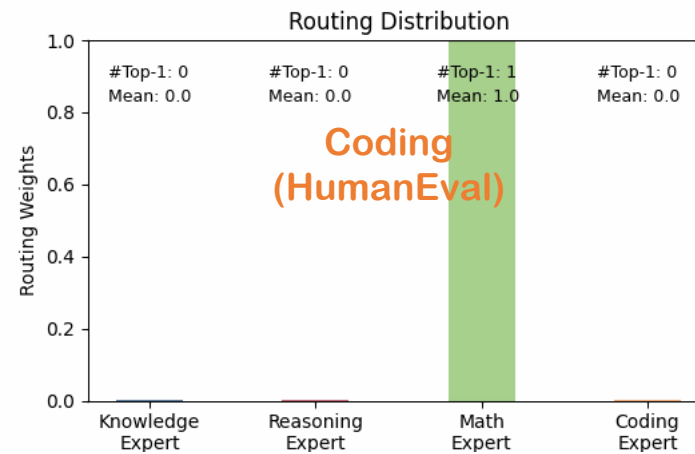
MiXSE (MiXture of Self-Specialized Experts)



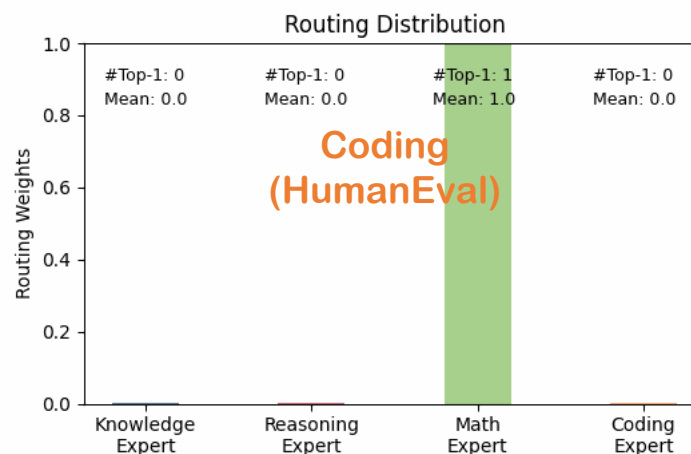
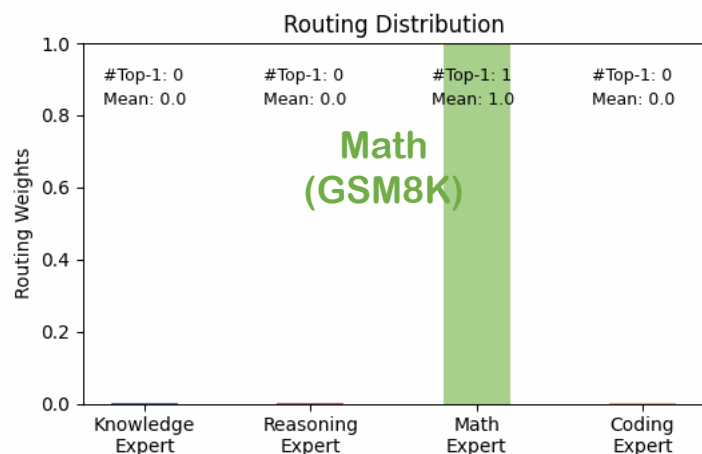
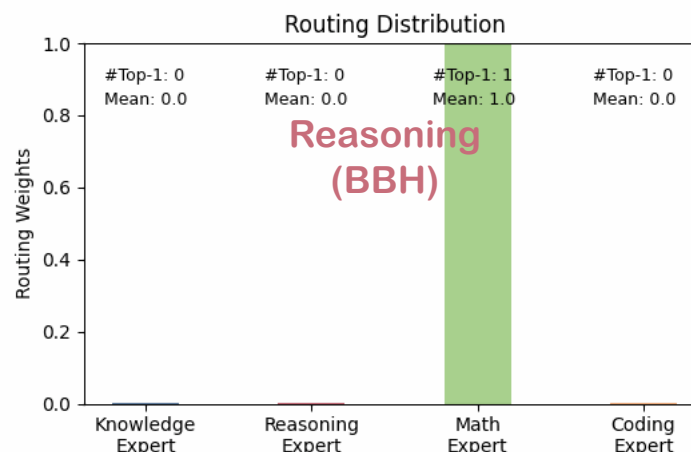
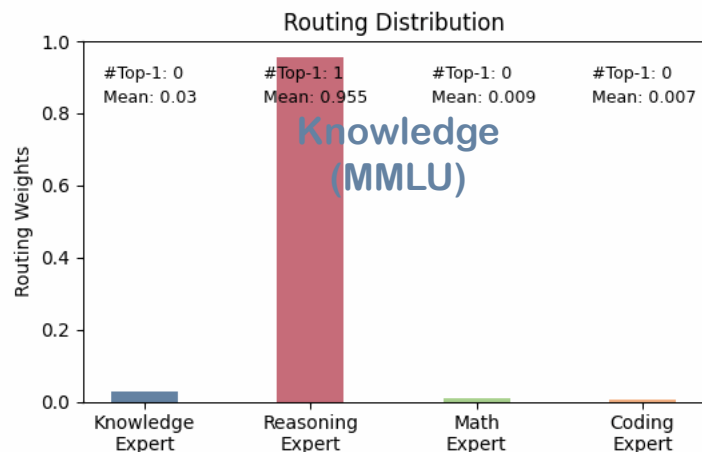
Improved Capabilities across All Domains



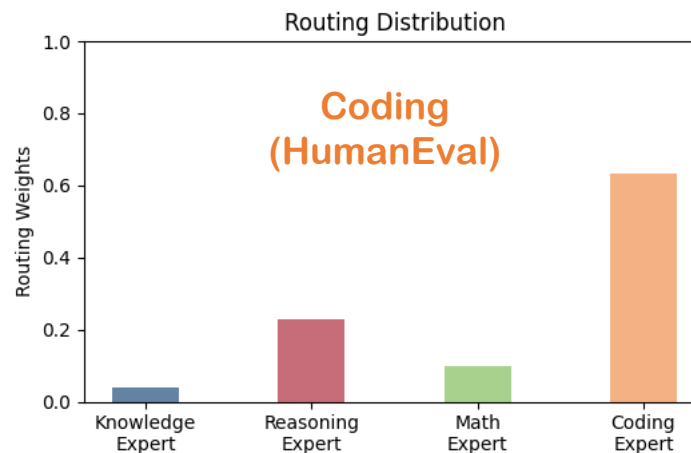
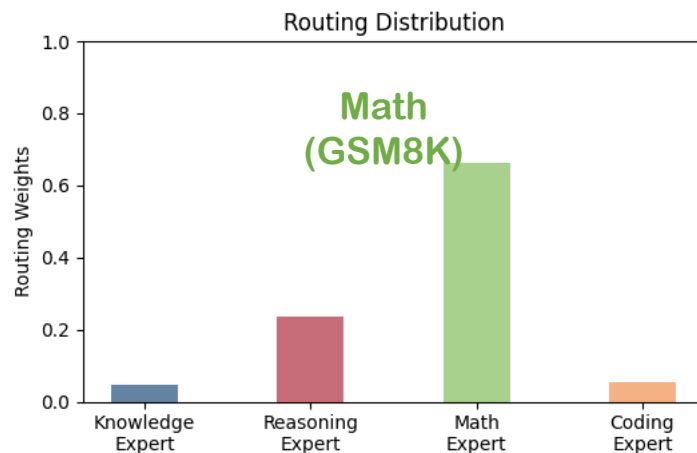
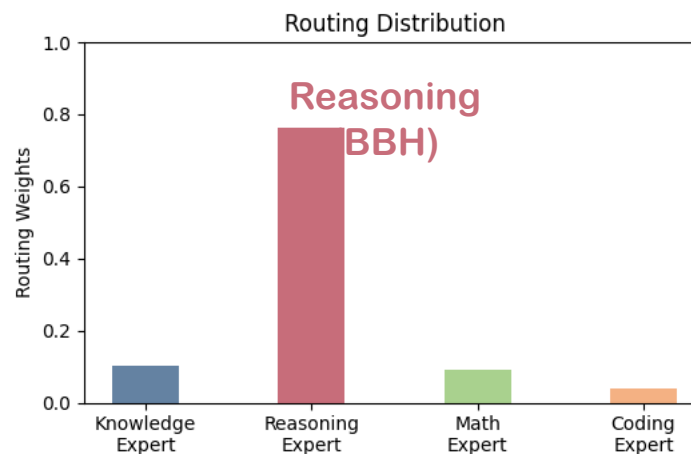
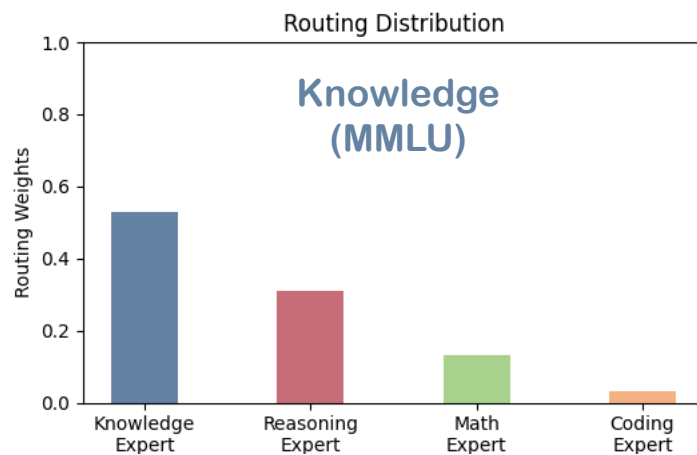
Routing Analysis (Token-by-Token)



Routing Analysis (Token-by-Token)



Routing Analysis (Mean over Tokens)

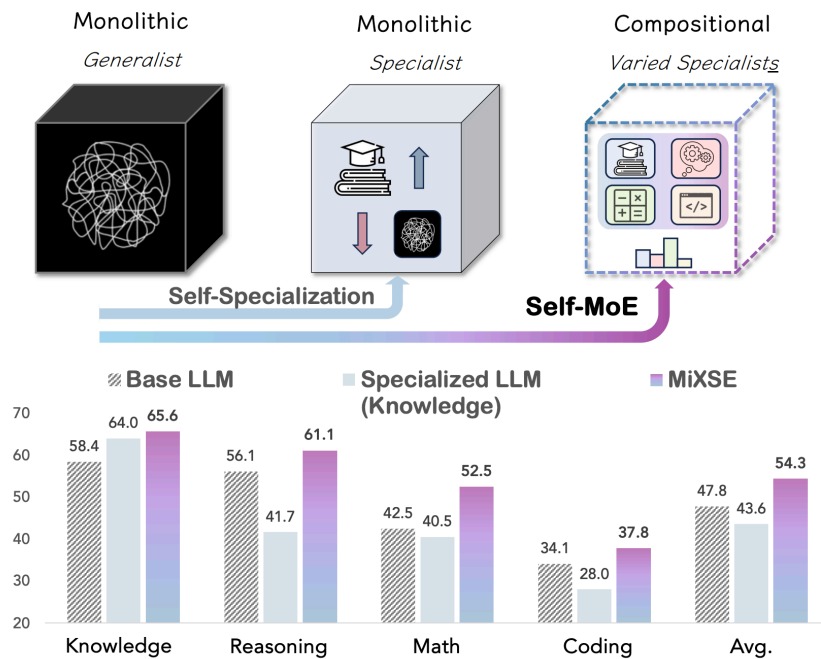


Key Takeaways

Q. Can we build compositional LLMs that achieve uncompromising multiple expertise with minimal resources?

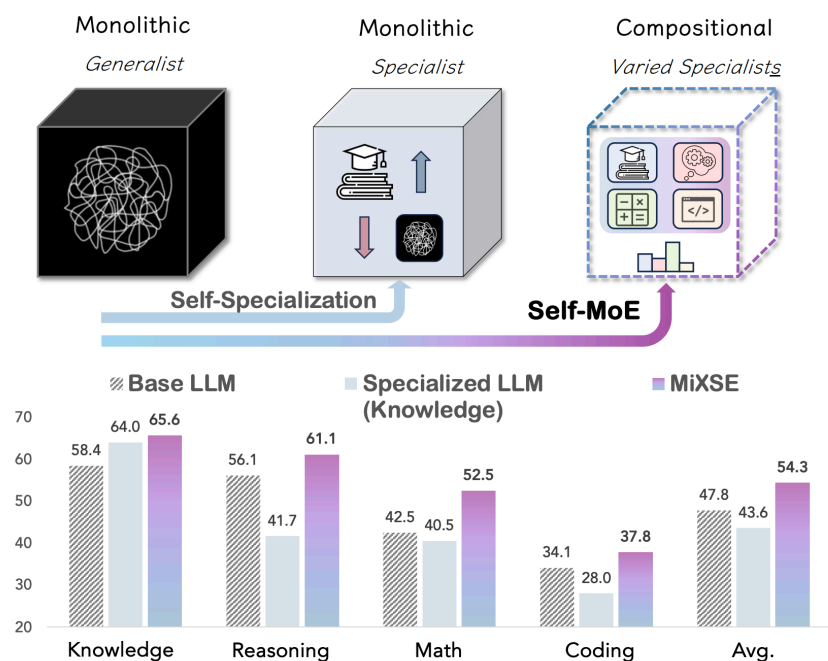
1. Highlighting Limitations of Monolithic Models

Focusing on a specific capability often comes at the cost of degrading performance in other domains



Key Takeaways

Q. Can we build compositional LLMs that achieve uncompromising multiple expertise with minimal resources?



1. Highlighting Limitations of Monolithic Models

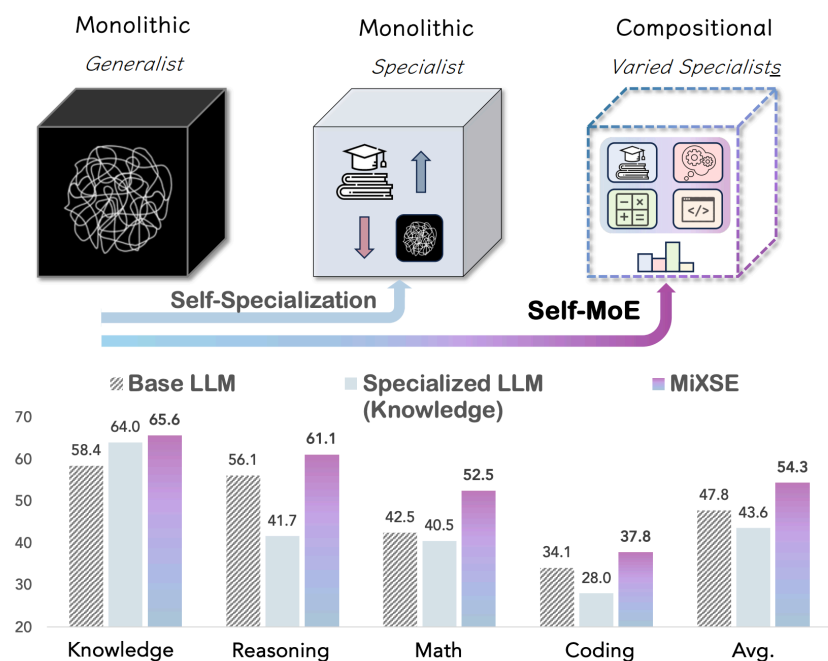
Focusing on a specific capability often comes at the cost of degrading performance in other domains

2. Introducing Self-MoE for Modular Specialization

Self-MoE transforms a monolithic LLM into a lightweight, modular system of self-specialized experts, without requiring extensive human supervision, compute resources, or added overhead in active parameters.

Key Takeaways

Q. Can we build compositional LLMs that achieve uncompromising multiple expertise with minimal resources?



1. Highlighting Limitations of Monolithic Models

Focusing on a specific capability often comes at the cost of degrading performance in other domains

2. Introducing Self-MoE for Modular Specialization

Self-MoE transforms a monolithic LLM into a lightweight, modular system of self-specialized experts, without requiring extensive human supervision, compute resources, or added overhead in active parameters.

3. Findings

- Consistent improvement over a base LLM, outperforming various baselines
- Ablation studies validate the impact of modularity, routing strategies, and self-generated synthetic data
- Analyses explore routing distributions, forgetting issues, and the applicability to various base LLMs

Thank you