# Self-Specialization: Uncovering Latent Expertise within Large Language Models

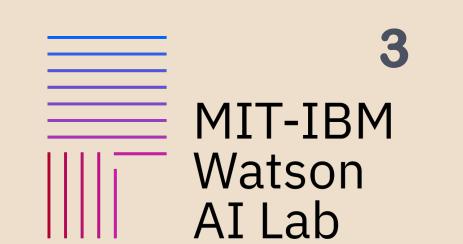
Junmo Kang<sup>1</sup> Hongyin Luo<sup>2</sup> Yada Zhu<sup>3</sup> Jacob Hansen<sup>2</sup> James Glass<sup>2</sup> David Cox<sup>3</sup> Alan Ritter<sup>1</sup> Rogerio Feris<sup>3</sup> Leonid Karlinsky<sup>3</sup>

Junmo.kang@gatech.edu



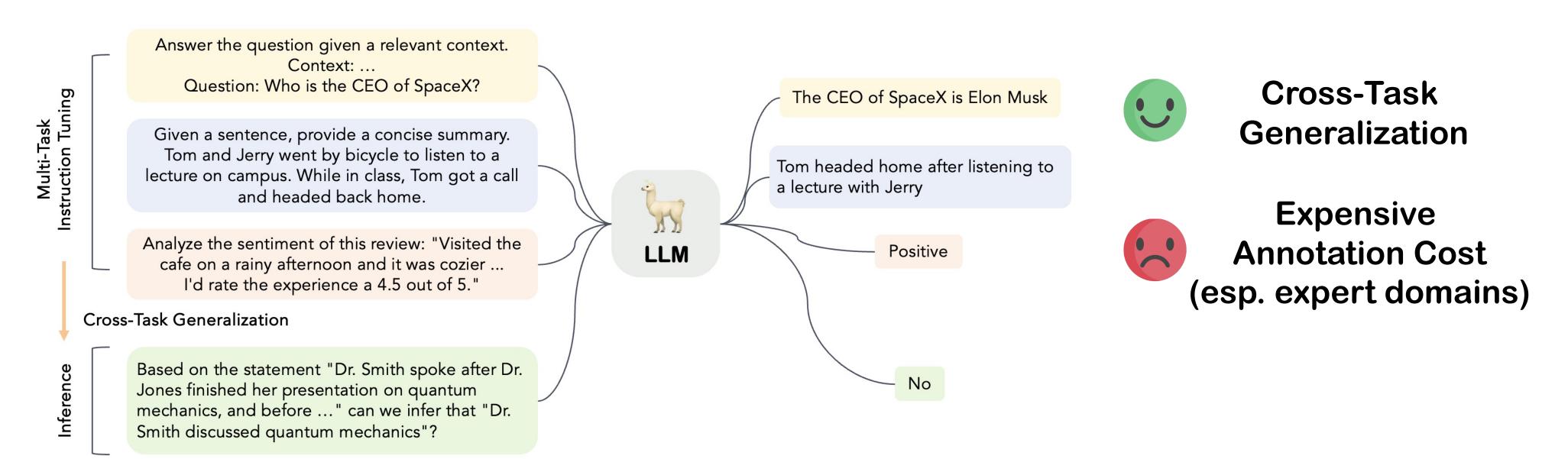
45.10 (+1.23)

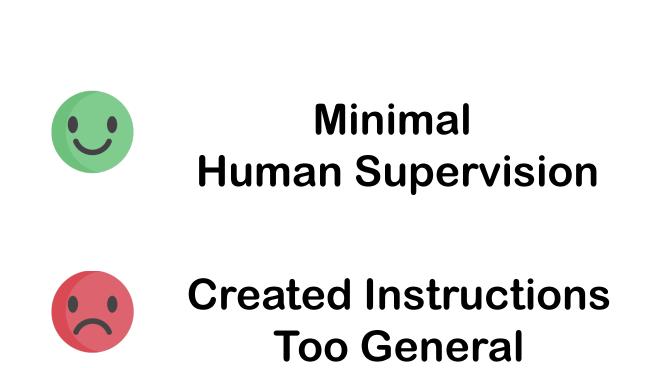


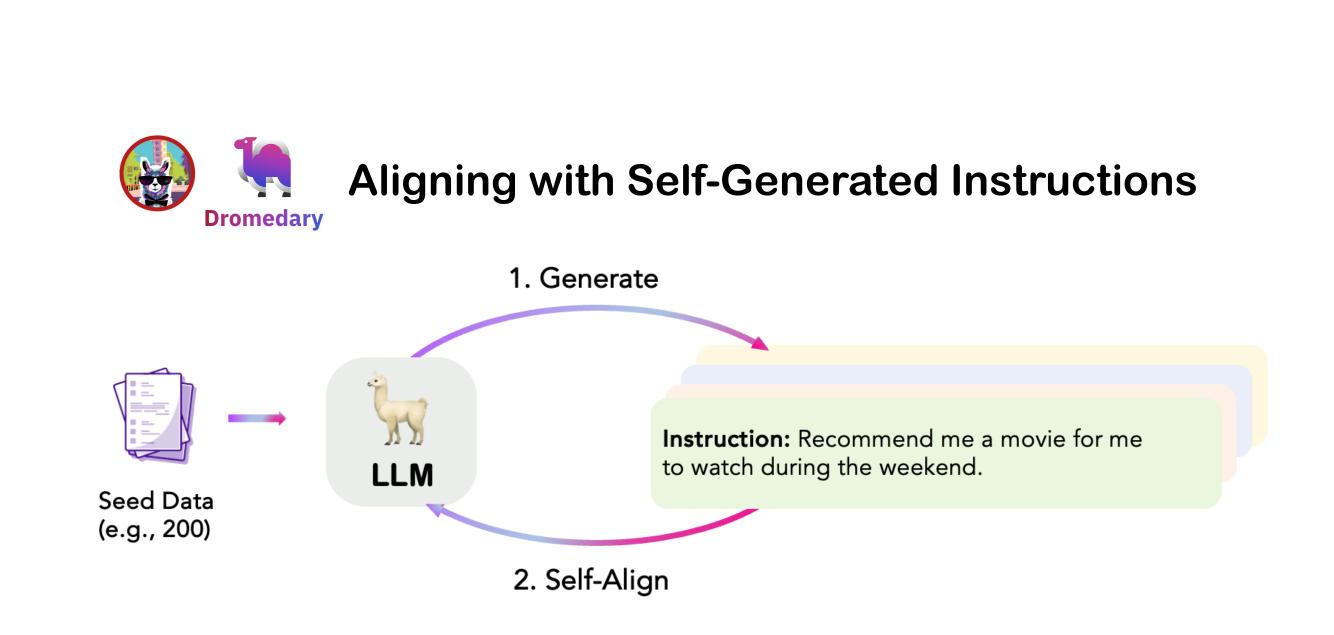


#### Instruction-Tuning & Self-Alignment

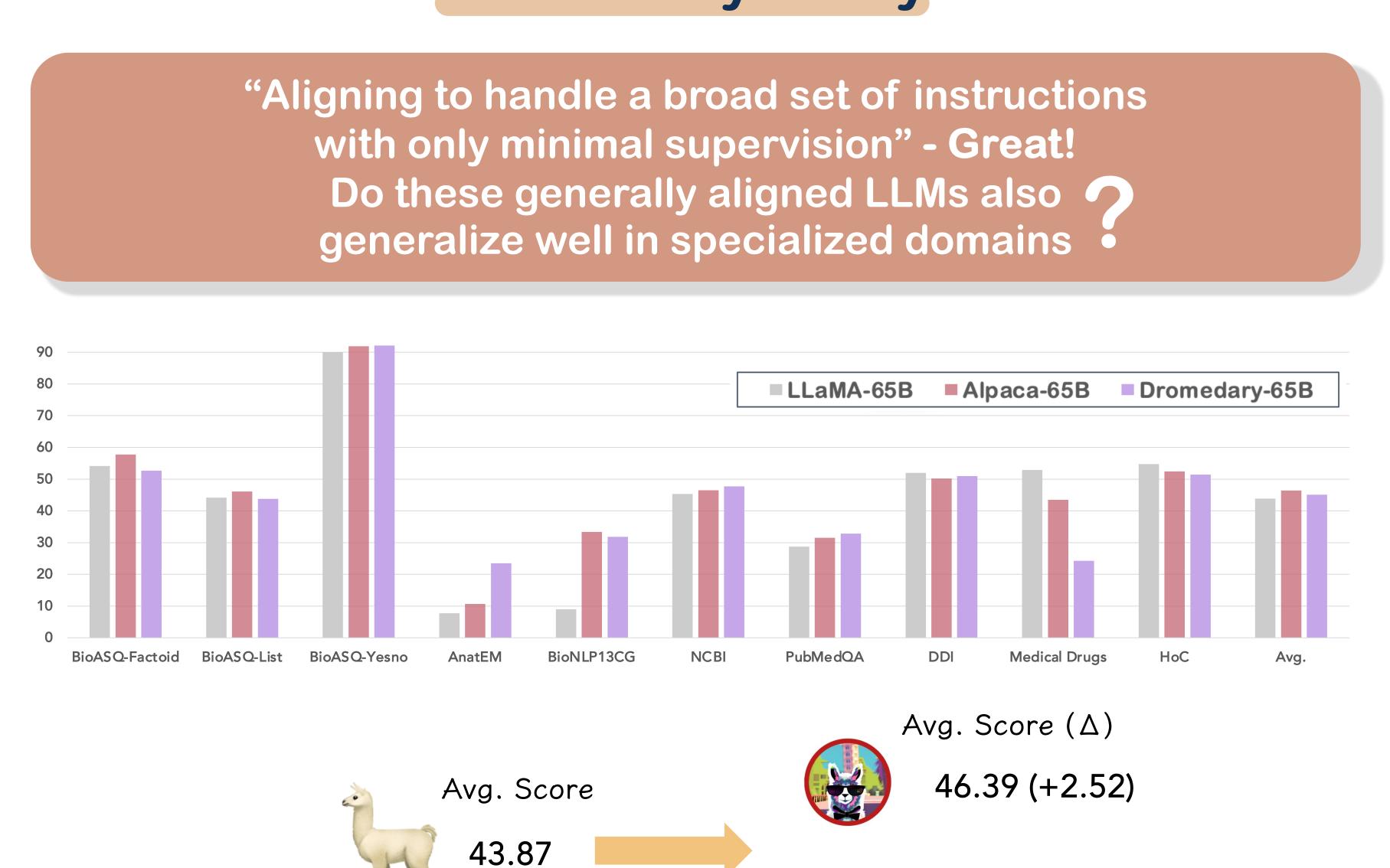
#### Multi-Task Tuning with Human Annotated Instructions







## **Preliminary Study**

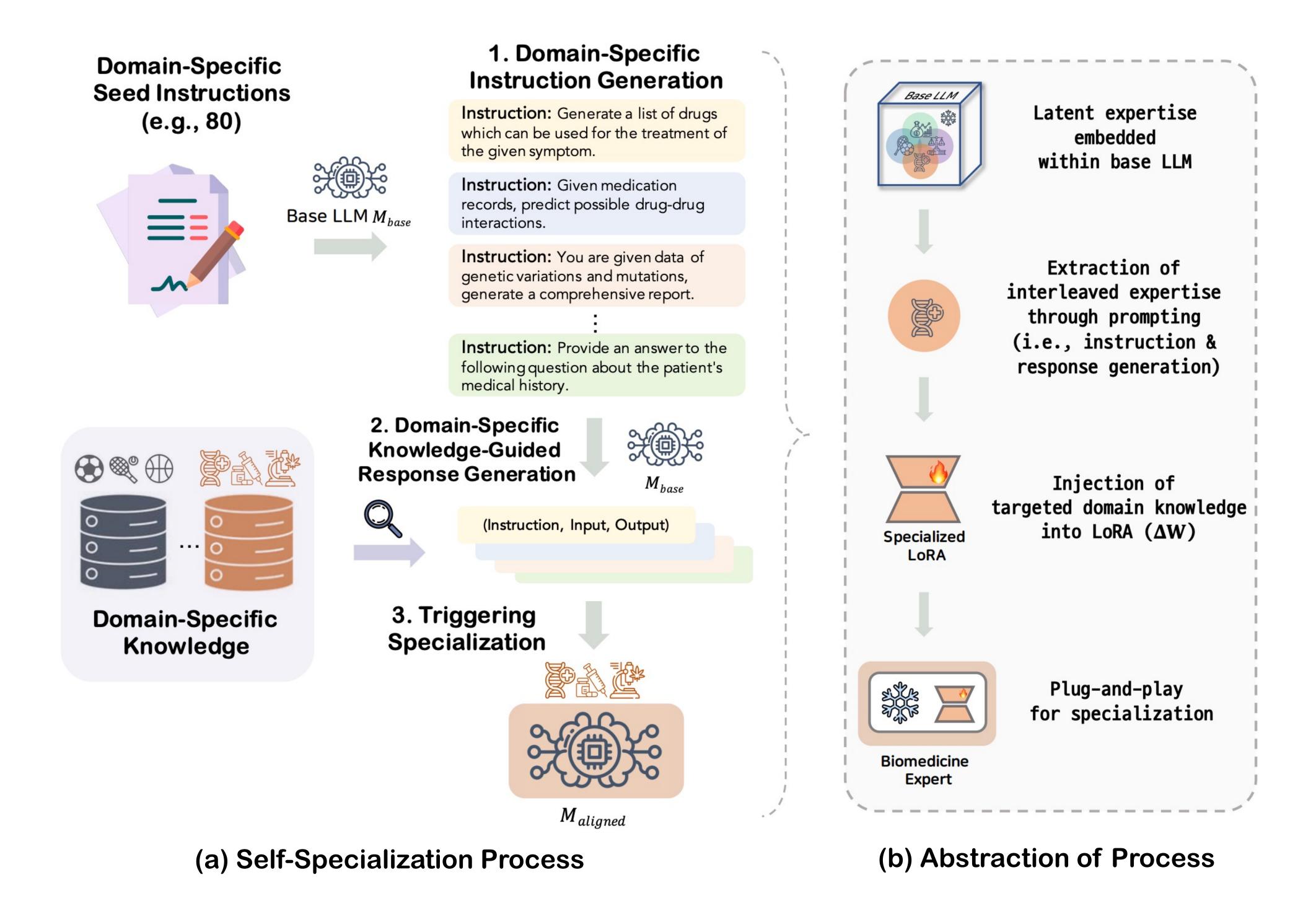


Only a slight advantage over the base model, although they are aligned to handle a broad set of instructions

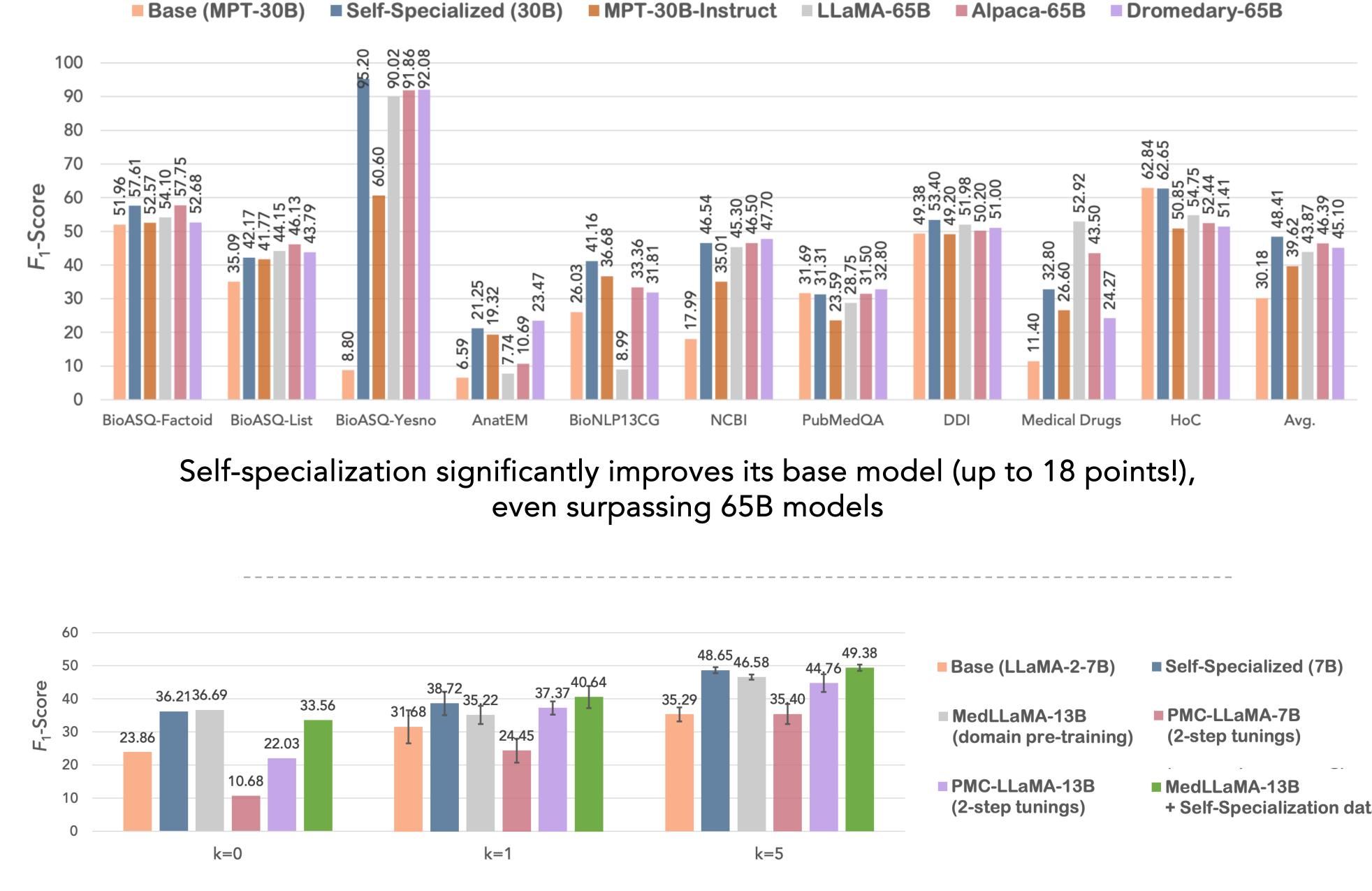
**Dromedary** 

Self-Align

# Self-Specialization for Uncovering Domain Expertise

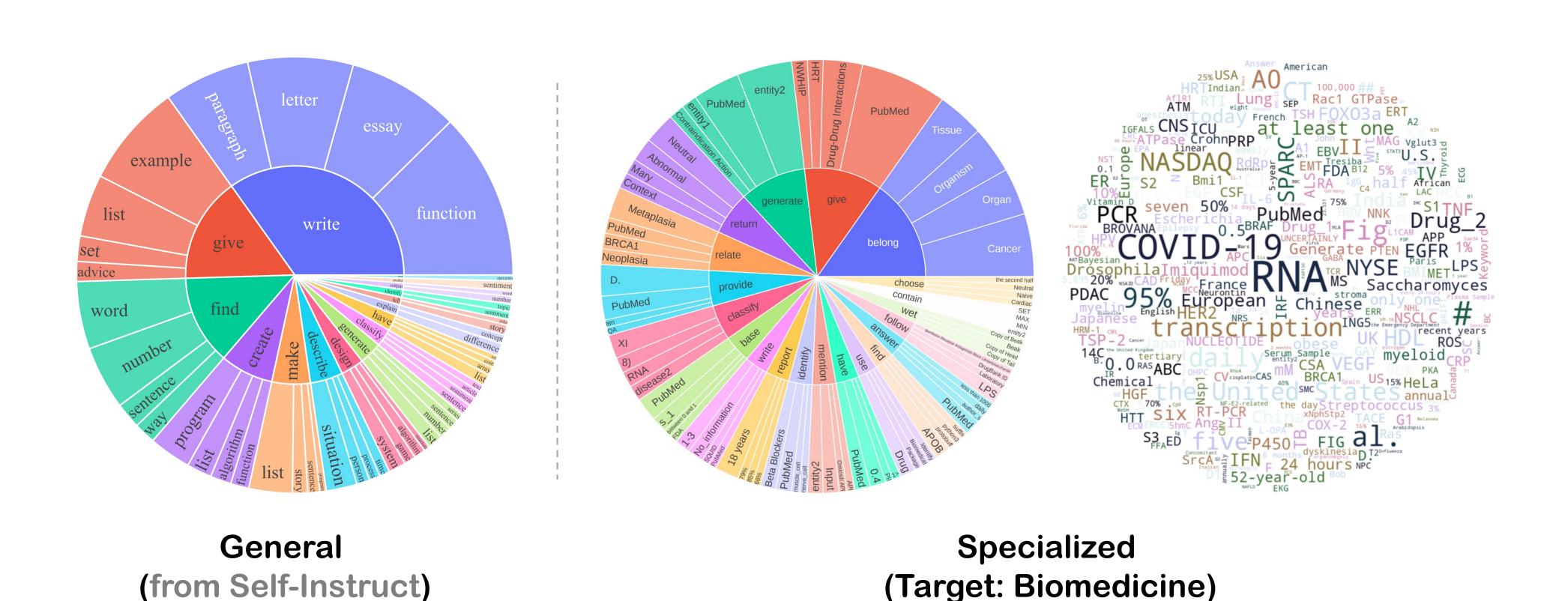


#### Results in Biomedical Domain



Self-specialization is on par or more effective than extensive domain pre-training + complementary

## Generated Data through Self-Specialization



Self-Specialization

Carving out latent expertise

Biomedicine Expert

Plug-and-play

Sports

Finance

Law

Plug-and-play

#### Key Takeaways

Q. Can we self-align LLMs with an expert domain like biomedicine with limited supervision?

- 1. Benchmarking of General-Purpose Aligned Models Highlighting the intrinsic challenge of encoding vast general knowledge into a finite set of parameters
- 2. Exploring a Lightweight Solution, Self-Specialization

  Targeted self-alignment to uncover latent expertise within LLMs with minimal supervision
- 3. Findings
  - Remarkable effectiveness in biomedical and financial domains
- Highly efficient and practical: Tuning with QLoRA on single A100 (using 5K generated data, ~3 hrs)