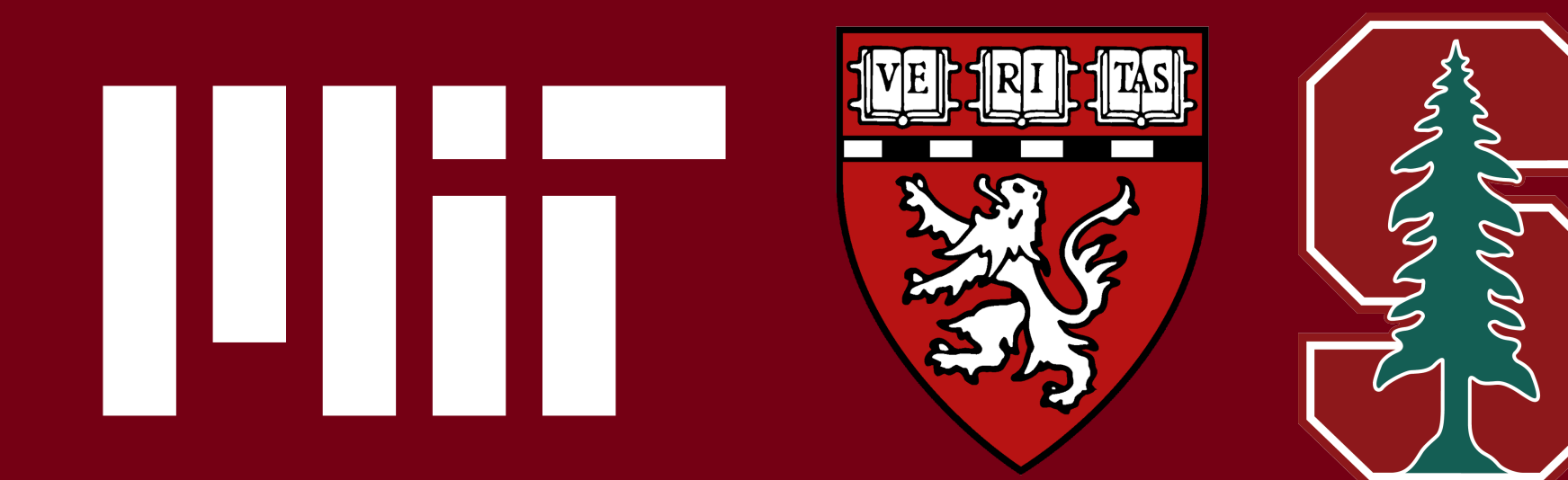




CLIMB: Data Foundations for Large Scale Multimodal Clinical Foundation Models

Wei Dai¹ Peilin Chen¹ Malinda Lu¹ Daniel Li¹ Haowen Wei^{1,2} Hejie Cui³ Paul Pu Liang¹

¹Massachusetts Institute of Technology ²Harvard Medical School ³Stanford University



We collected *CLIMB*, one of the largest multimodal clinical dataset across 15 modalities. Models trained on *CLIMB* achieve up to 32.54% AUC improvements over previous SoTA.

We introduce CLIMB (Clinical Large-scale Integrative Multi-modal Benchmark), a comprehensive dataset unifying 4.51M samples across 44 datasets, totaling 19.01 TB. Multitask pretraining on CLIMB significantly improves performance by up to 32.54% especially for understudied modalities, enables strong few-shot transfer to novel tasks, and enhances multimodal fusion.

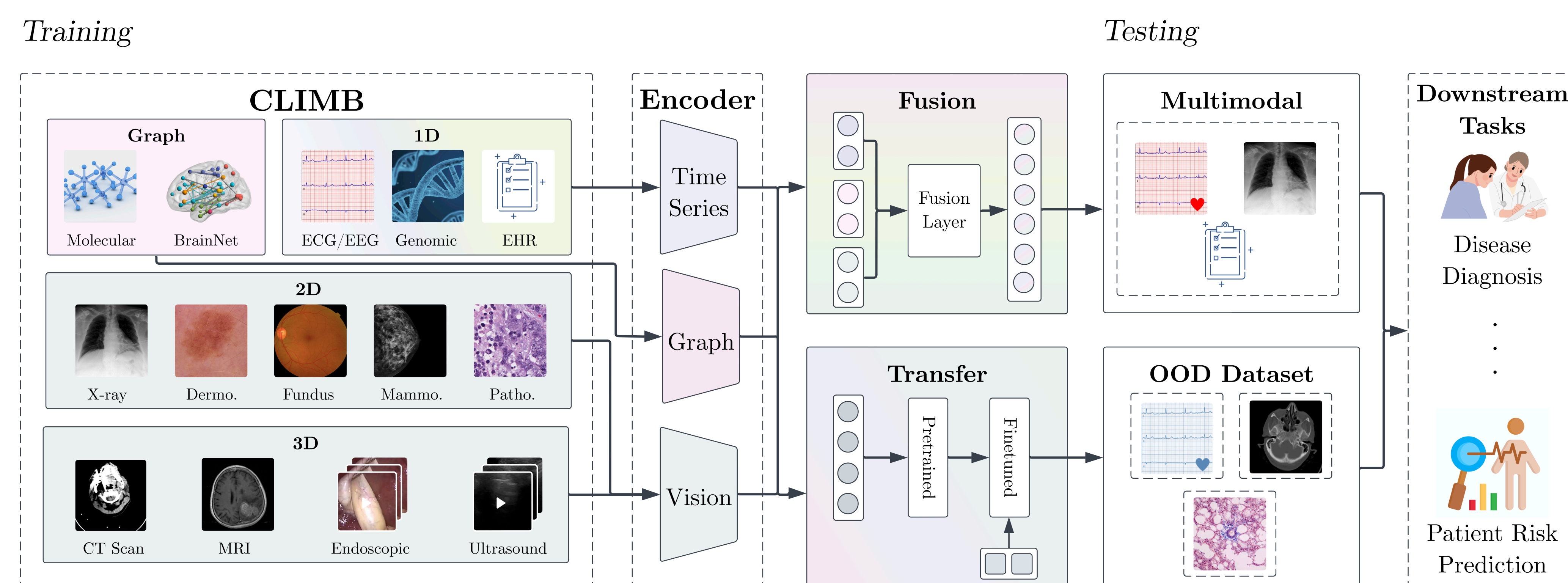
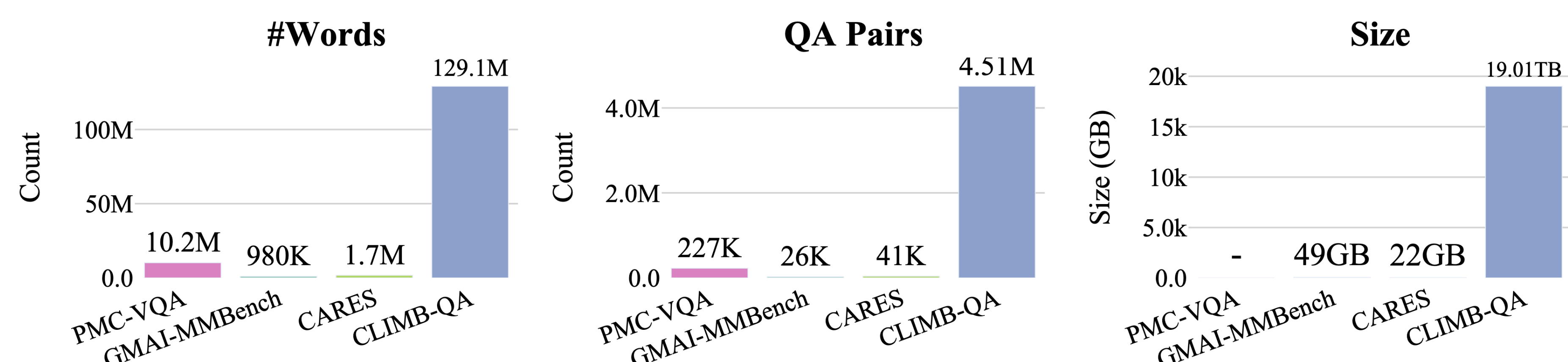


Figure 2. (a) Visualization of CLIMB dataset composition. (b) Focus of dataset collection. (c) Distribution of data collection sites in CLIMB. (d) Example code usage on CLIMB framework. (e) Sample data from CLIMB.

Experimental Results

Multitask pretraining significantly improves performance across tasks, achieving up to **32.54%** AUC improvement in understudied areas.

Few-shot performance on models trained on CLIMB demonstrate significant improvements, achieving up to **29%** improvement.

Multimodal fusion: Single-modality pretraining on CLIMB enhances multimodal learning performance, leading to successful transfer to MIMIC-IV.

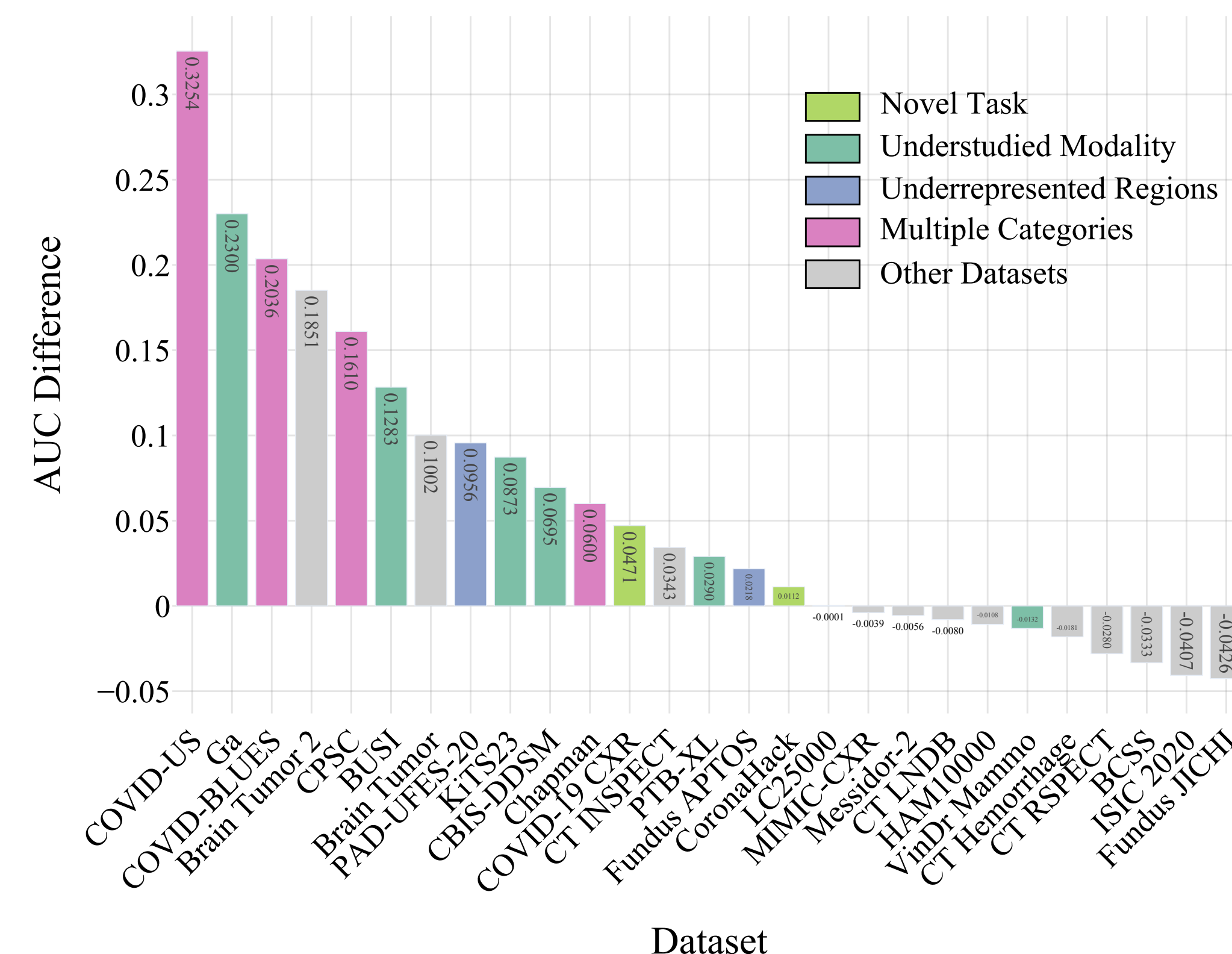


Figure 3. Difference in AUC achieved by the multitask model compared to single-task training.

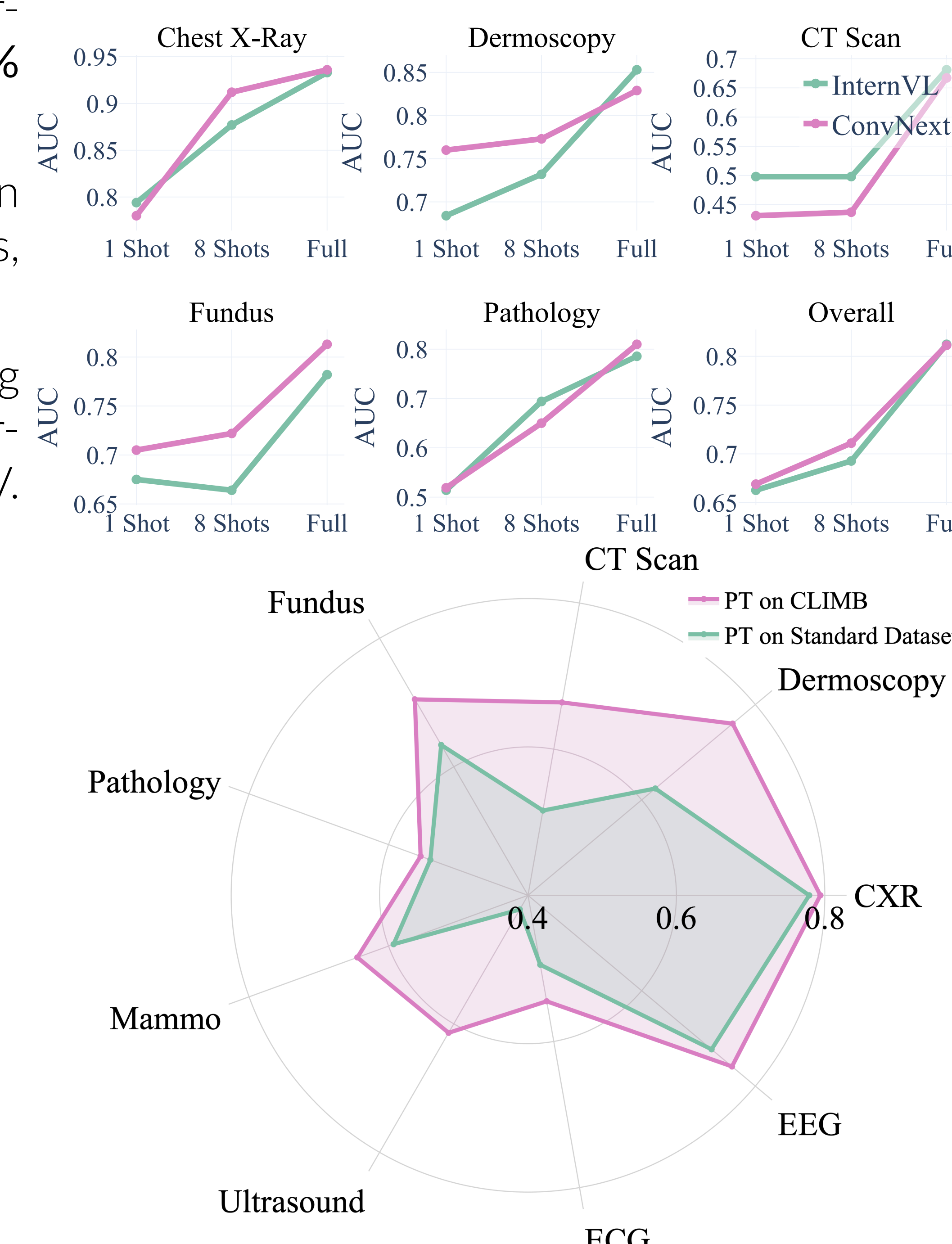
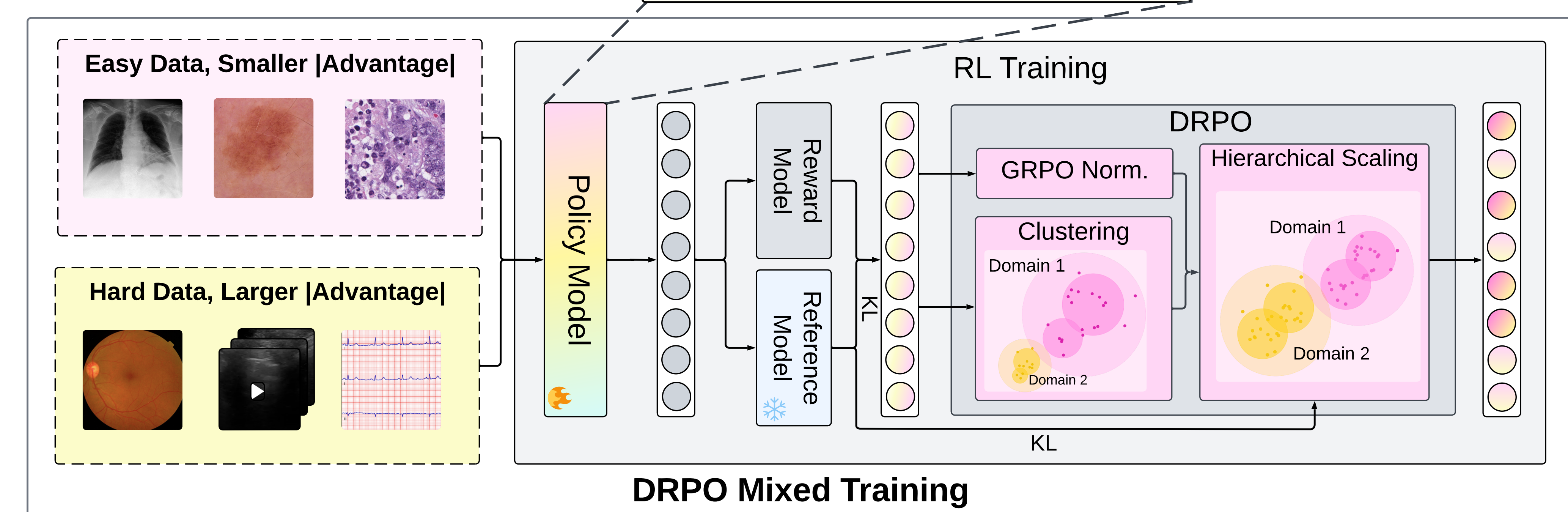
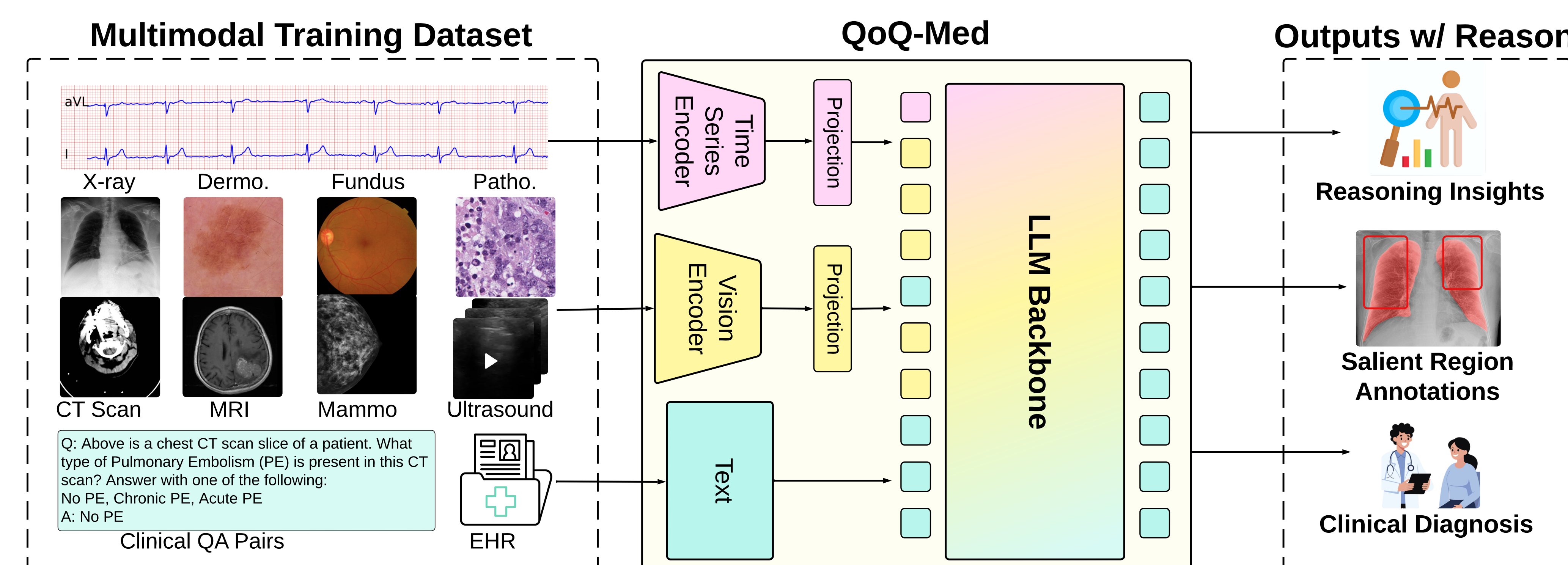


Figure 4. Few-shot performance of models across different pretraining (PT) datasets.

Below describes our preprint

QoQ-Med: Building Multimodal Clinical Foundation Models with Domain-Aware GRPO Training



Challenge: In multimodal training, frequent modalities and easy samples dominates training.
Idea: Hierarchically clusters samples during training, upscaling rewards from rare and hard domains.

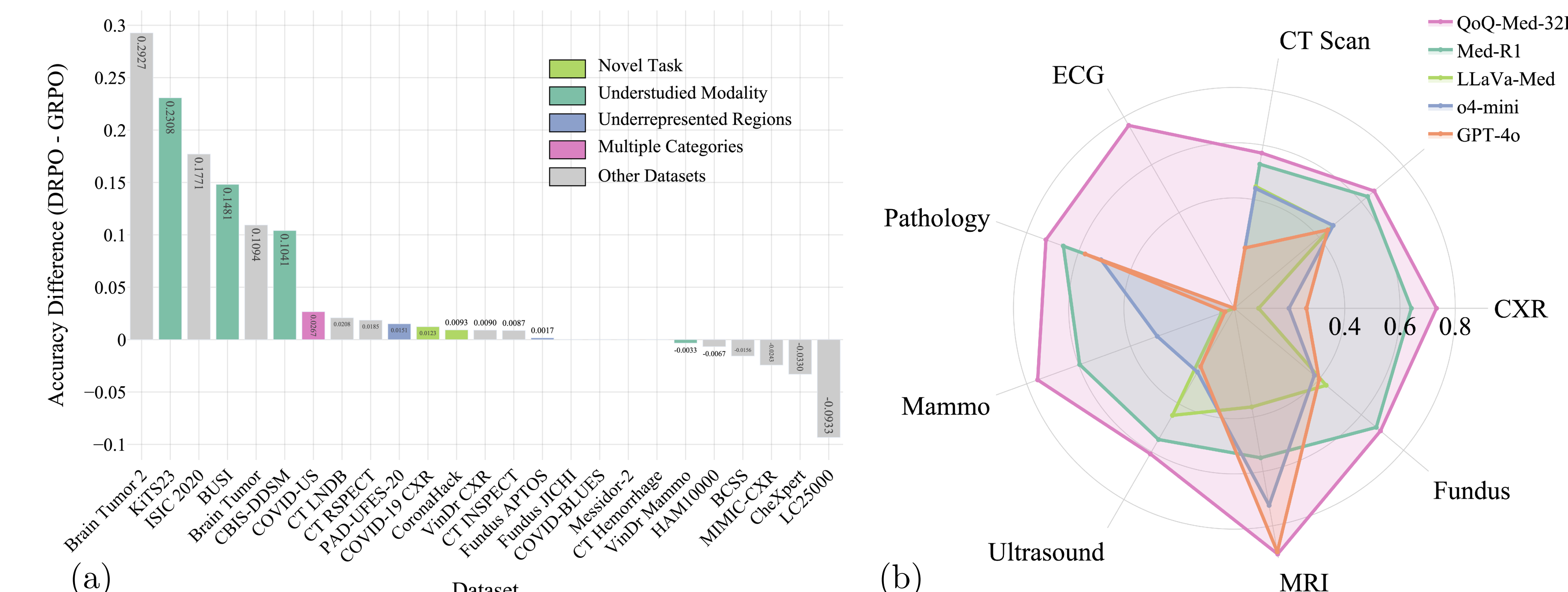
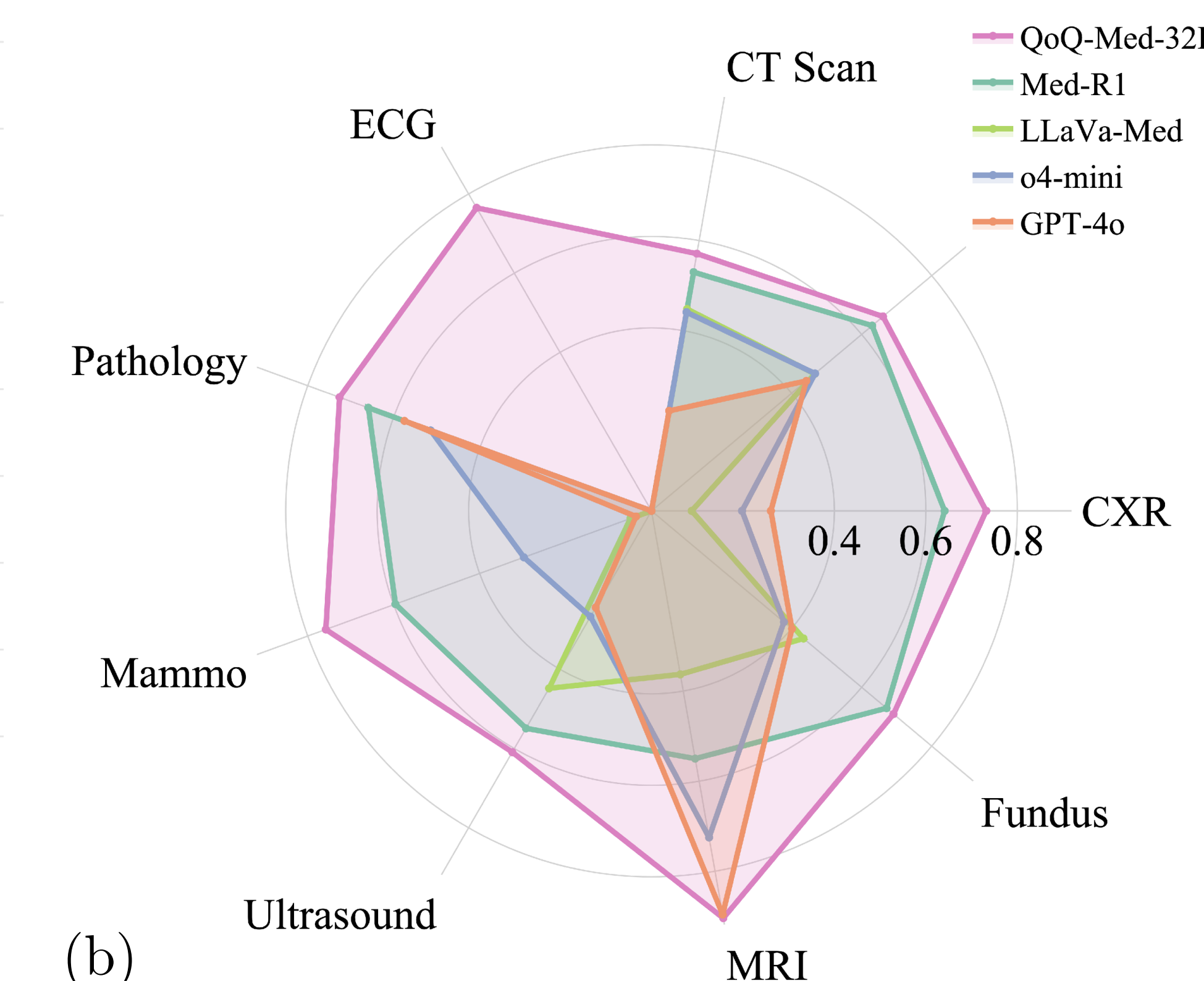


Figure 5. (a) DRPO performance vs GRPO.



(b) QoQ-Med vs other open/closed source models

* We start from current pretrained SoTA, then train the model on CLIMB with the following tasks:

