# Benchmarking general and retina-specific foundation models with quality-aware conditioning for fundus phenotyping

*Yosef Solewicz, Technion - Israel Institute of Technology*

We present a unified benchmark for retinal fundus analysis that compares frozen and fine-tuned vision foundation models across five datasets (UKBB, HYFundus, BRSET, ODIR, Papila). Evaluation is structured as partly in-domain (datasets used for training/validation) and partly out-of-domain (held-out external datasets), enabling a direct assessment of transfer robustness.

The benchmark includes both medical and general backbones (MedSigLIP, RETFound variants, and DINO-family models) and evaluates age regression, sex classification, and eye laterality prediction. We first test frozen representations with lightweight task heads, then study adaptation strategies from parameter-efficient tuning to full fine-tuning. Across datasets, fine-tuning is the main source of improvement, especially for age prediction and for backbones that are weaker in the frozen setting. Full fine-tuning generally narrows performance gaps between models and improves OOD robustness. In contrast, gains for sex and eye are usually smaller because frozen baselines are already strong, with near-saturated eye laterality in several settings.

We also tested quality-aware conditioning modules at embedding, score, and input levels. These produced task- and dataset-specific gains, but they were less consistent and typically smaller than the gains from backbone fine-tuning. Overall, our results identify fine-tuning as the most effective adaptation strategy in this benchmark.