

Knowledge-grounded Adaptation Strategy for Vision-language Models: Building Unique Case-set for Screening Mammograms for Residents Training

Aisha Urooj Khan¹, John Garrett², Tyler Bradshaw², Lonie Salkowski², Jiwoong Jason Jeong³, Amara Tariq¹, and Imon Banerjee^{1,3}

¹ Department of Radiology, Mayo Clinic

² Department of Radiology, UW Madison School of Medicine and Public Health

³ School of Computing and Augmented Intelligence, Arizona State University

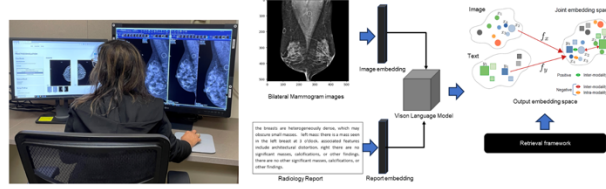


Fig. 1: Multimodal learning for screening mammogram: (a) a session with radiology resident for the case review; (b) framework generating joint embedding space for bilateral mammogram and free-text radiology reports. Illustration of joint embedding space (right) is adapted from CrossCLR [22].

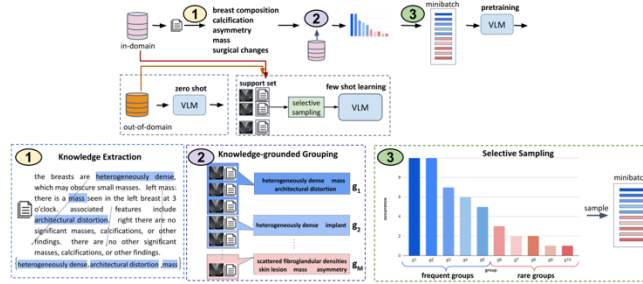


Fig. 2: Workflow for adapting the VLM with the proposed selective sampling to learn joint representation aware of fine-grained knowledge. The pretrained model is tested on out-of-domain data for zero shot evaluation. For few shot learning, support set is obtained from the training data to fine-tune model.

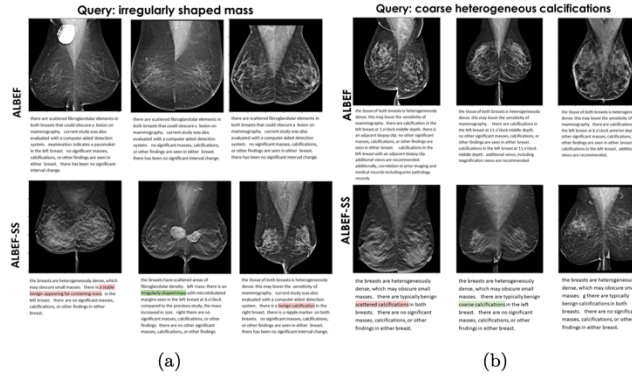


Fig. 3: Qualitative results for Retrieval model. An example with highlighted green words is marked relevant by the radiologist for case build. Concepts highlighted with the pink show not exact but related findings in the image-report pair.