

MedCLIP-SAM: Bridging Text and Images for Universal Medical Segmentation

Taha Koleilat^{1*}, Hojat Asgariandehkordi¹, Hassan Rivaz¹, and Yiming Xiao²

¹ Department of Electrical and Computer Engineering, Concordia University,
Montreal, Canada

{taha.koleilat,hojat.asgariandehkordi,hassan.rivaz}@concordia.ca

² Department of Computer Science and Software Engineering, Concordia University,
Montreal, Canada
yiming.xiao@concordia.ca

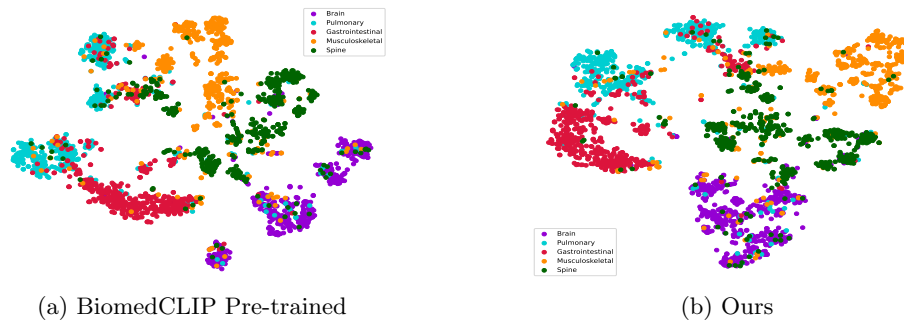


Fig. 1: Visualization of 2D t-SNE image embeddings based on the anatomical categories of the medical diagnosis, computed based on the MedPix dataset, which was used to finetune the BiomedCLIP model.

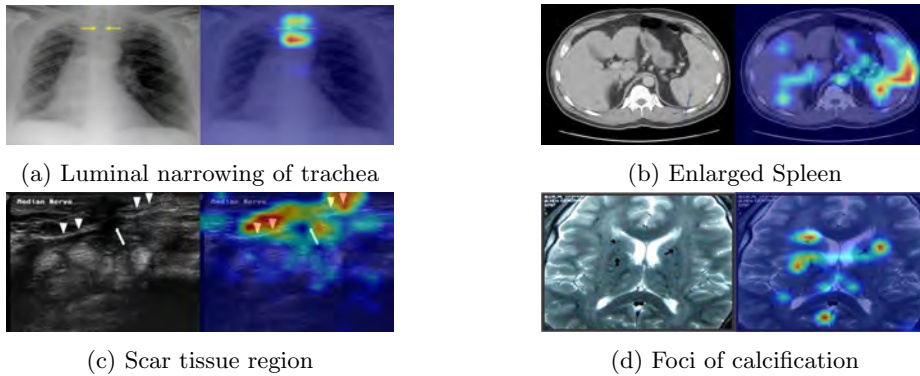


Fig. 2: Sample attention heatmaps with gScoreCAM after providing the text prompt to localize the region of abnormality in the radiology scans taken from the ROCO validation set. Arrows show the anatomical/pathological regions.

* Corresponding author.