# VividMed: Vision Language Model with Versatile Visual Grounding for Medicine

**Lingxiao Luo\*, Bingda Tang\*, Xuanzhong Chen, Rong Han, and Ting Chen**
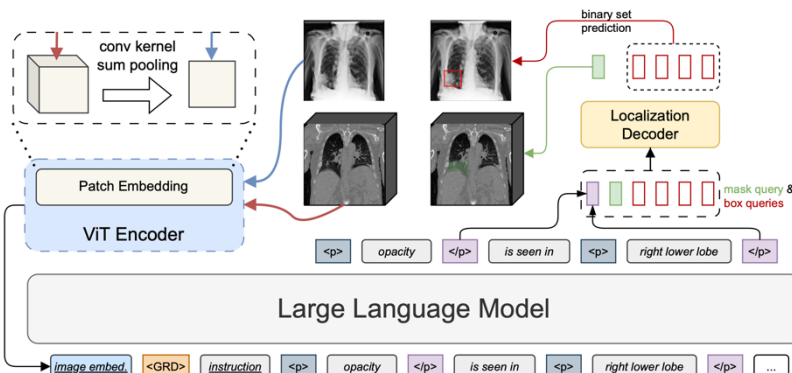
Tsinghua University

Figure 1: The architecture of VividMed, which is built upon a base VLM (left and lower) and a promptable localization module (upper right). The model identifies key phrases for grounding by enclosing them with bracket tokens, and the hidden states of the closed bracket token is used for prompting the localization module. The model accepts both 2D and 3D images as input by adaptively adjusting weights in the patch embedding layer. The vision encoder of the localization module is omitted for clarity.
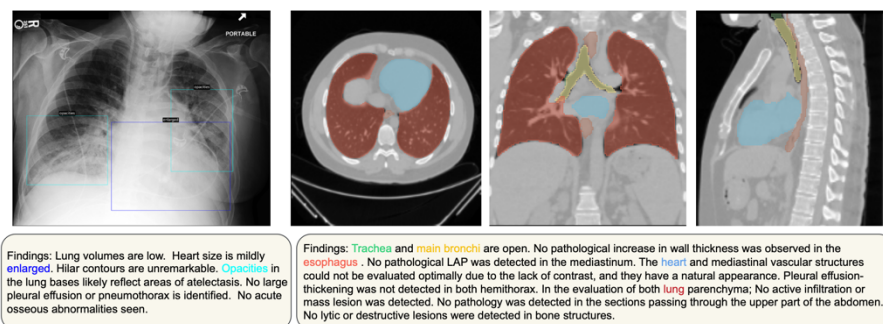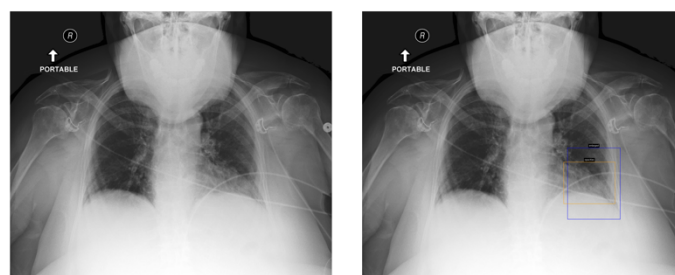


Figure 2: Selected qualitative results for grounded report generation, zoom in for better view. Impressions are omitted for clarity.



Figure 3: In this example, the model wrongly identifies cardiomegaly and gives an unusual visual grounding result, which may remind the radiologist in clinical practice.