#### PTCOG-AO2025-ABS-0091

Vision-language Model for Enhanced Organ Segmentation in Proton/Heavy Ion Radiotherapy

Chenbin Liu\*, Yihao Zhao¹, Cuiyun Yuan¹, Ying Liang¹, Yang Li¹, Chunxia Li¹, Man Zhao¹, Jun Hu¹, Ningze Zhong¹, Wei Liu²

\*'<sup>1</sup> Radiation Oncology, Chinese Academy of Medical Sciences, Cancer Hospital, Shenzhen Center, China, <sup>2</sup> Radiation Oncology, Mayo Clinic Arizona, United States of America

# **Objectives**

Accurate delineation of organs-at-risk is critical for proton and heavy ion radiotherapy planning, as precise dose delivery depends heavily on robust segmentation of anatomical structures. Current automated segmentation methods often underutilize prior anatomical knowledge and expert textual descriptions, limiting their accuracy and clinical applicability. This study aims to develop and validate a vision-language model that leverages detailed organ descriptions to improve segmentation performance specifically for proton/heavy ion therapy workflows.

#### Methods

We develop Med-VLM, a novel medical image segmentation framework that integrates transformer-based vision encoders (ViT-base) with text encoders (BioBERT) using low-rank adaptation (LoRA) and a query-based feature mixer. The model utilizes expert-provided descriptions of organ location and spatial relationships as input, alongside CT image slices. During training, image and text encoder weights remain frozen, while adapters, the feature mixer, and a mask decoder are fine-tuned. The mask decoder employs dilated convolutions and a bi-directional transformer to generate precise segmentation masks. Training was conducted on three large medical image datasets (FLARE, SegTHOR, and MSD) comprising 47,000 images and 100,000 high-quality masks. Segmentation performance was evaluated using the Dice Similarity Coefficient (DSC), 95th percentile Hausdorff Distance (HD95), and Average Surface Distance (ASD).

## Results

Our model demonstrated superior segmentation accuracy compared to state-of-the-art methods (nnUnet, SAM, MedSAM, Lvit) across all datasets. On the combined test set, Med-VLM achieved DSC 0.98±0.03 (vs. 0.95±0.03 for nnUnet), HD95 21.86±39.95 (vs. 45.13±62.61 for nnUnet), and ASD 4.56±4.67 (vs. 13.11±13.94 for nnUnet). The incorporation of authoritative organ descriptions significantly improved boundary delineation, especially in regions critical for proton/heavy ion therapy planning. Ablation studies confirmed the importance of detailed textual input and the effectiveness of LoRA and BioBERT-based text encoding.

### **Conclusions**

Our model represents a significant advance in automatic organ segmentation for proton and heavy ion radiotherapy, leveraging expert textual descriptions to achieve state-of-the-art accuracy and boundary delineation. This approach has strong potential to streamline and enhance clinical workflows in advanced radiotherapy, supporting more precise and reliable dose delivery to target volumes while minimizing exposure to organs-at-risk.

