

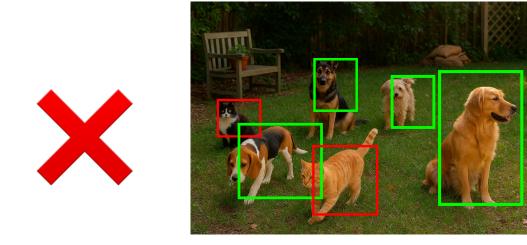
Samsung **R&D Institute** Philippines

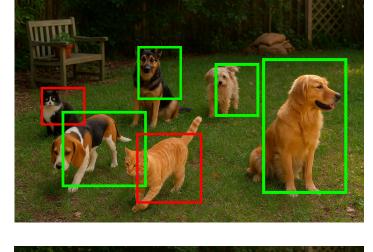
Interpreting Open-Vocabulary Referring Object Detection with Reverse Contrast Attention Drandreb Earl Juanico*, Rowel Atienza, Jeffrey Kenneth Go

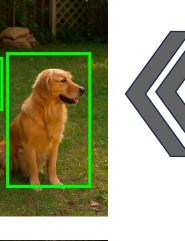


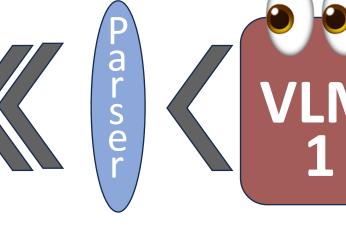


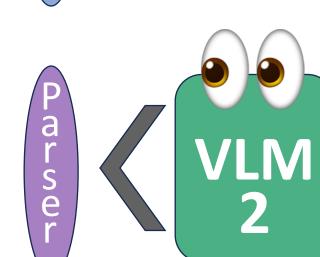
Open-Vocabulary Referring Object Detection







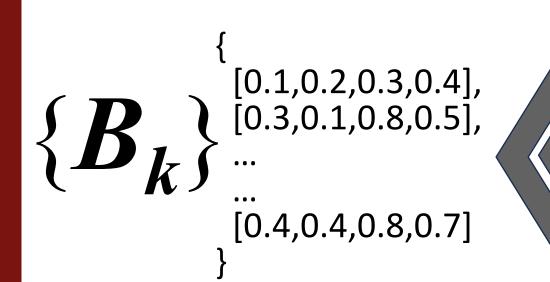


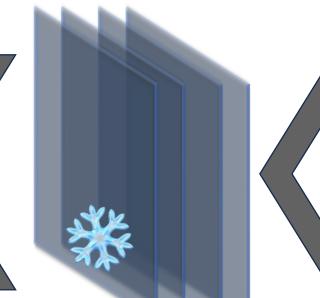


Find me the dogs you see right now.



Inference





LLM

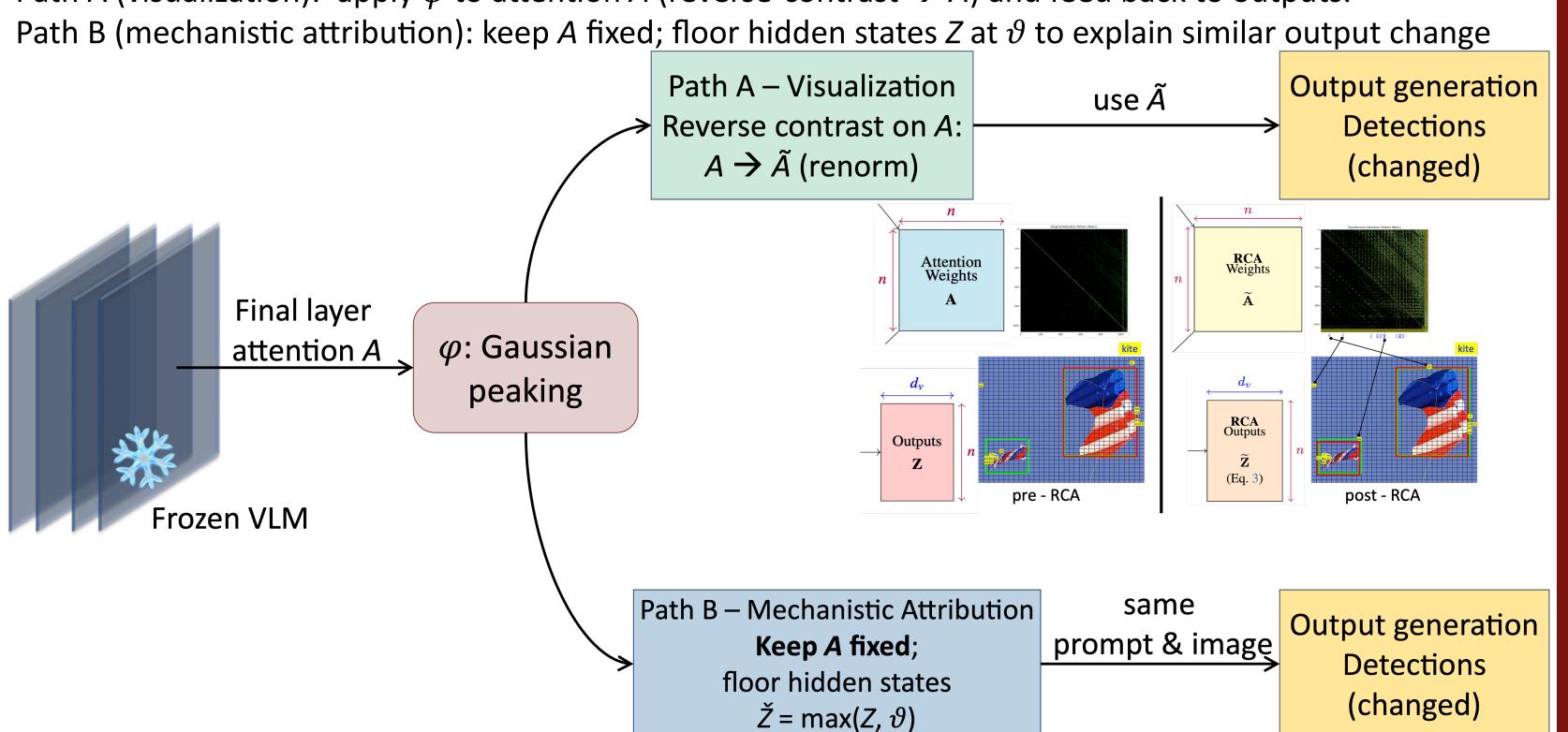


Give the normalized bounding box coordinates in the format [x1, y1, x2, y2] of all instances of {cls} in the image.

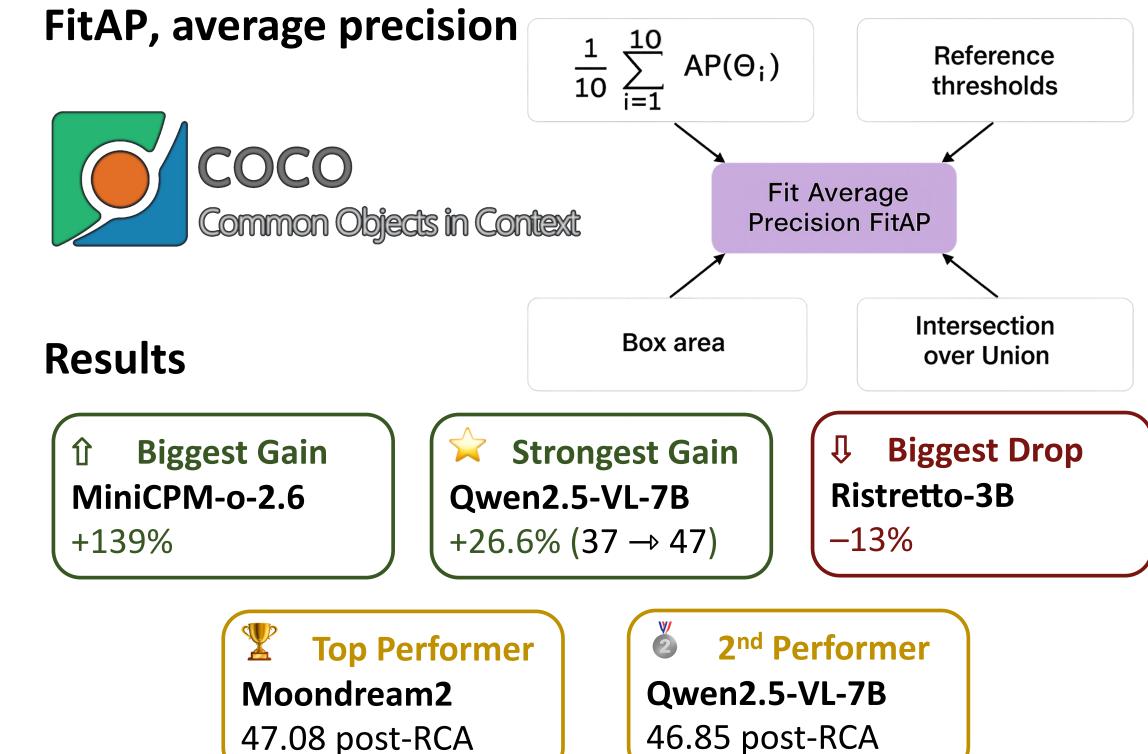
35B

Reverse Contrast Attention

RCA: Two linked interpretations from one transformer Path A (visualization): apply φ to attention A (reverse-contrast $\rightarrow \tilde{A}$) and feed back to outputs.

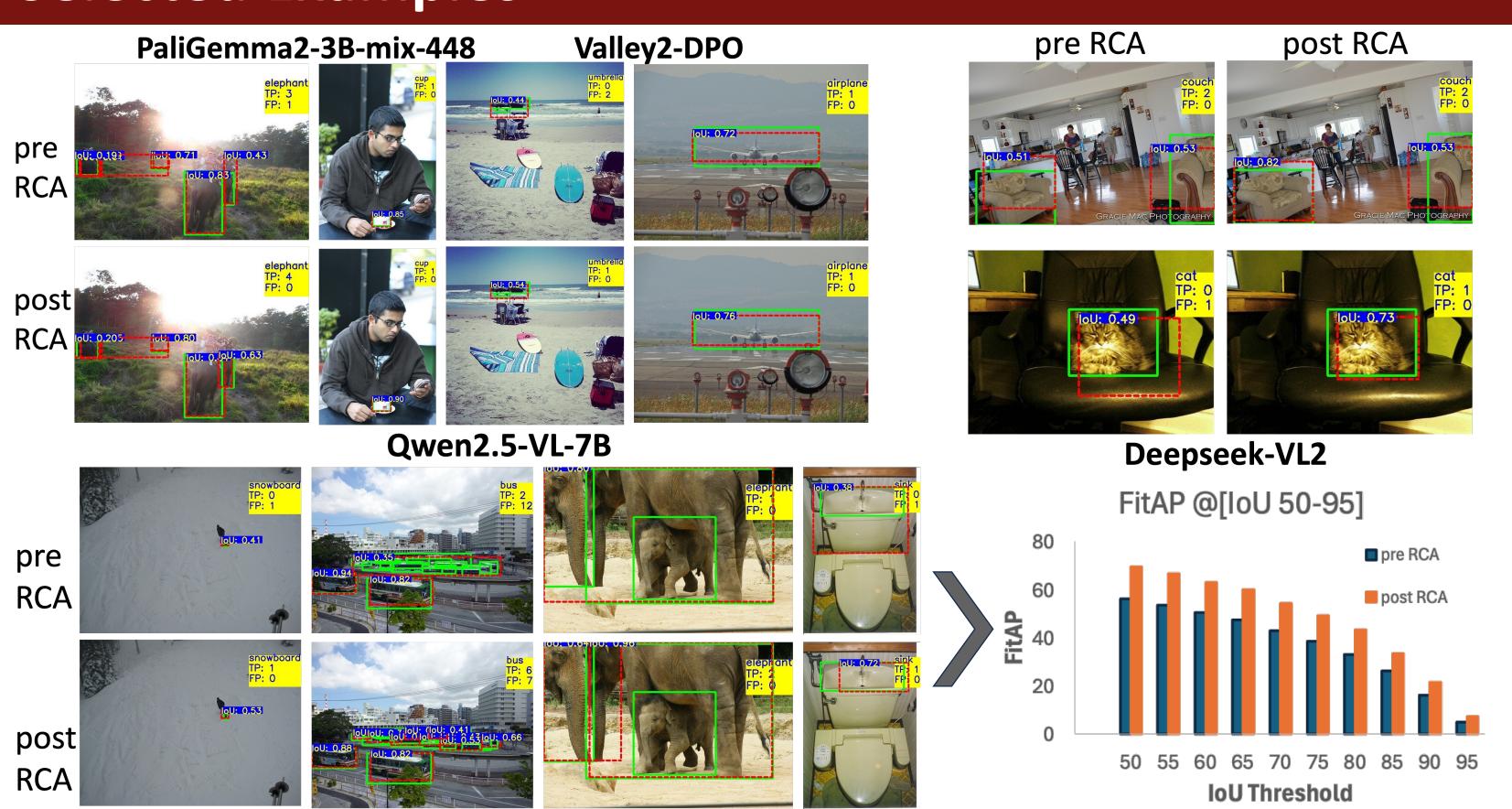


Evaluation

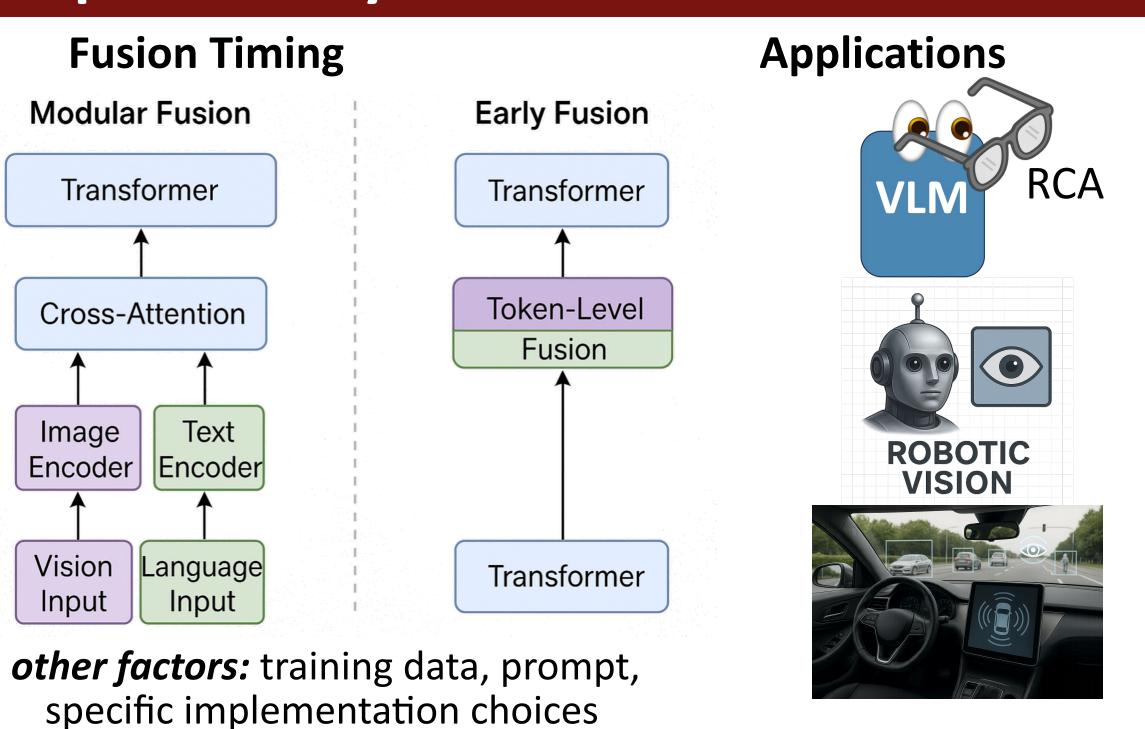


11 of 15 selected VLM gained post-RCA

Selected Examples



Explainability



VLM Selection

Best models via:



OpenCompass 司南

Checkpoints available via:



Hugging Face



Acknowledgement: This work was supported by Samsung Research Philippines. DEJ acknowledges the support from the DOST-SEI, ERDT, and ICCV Broadening Participation