



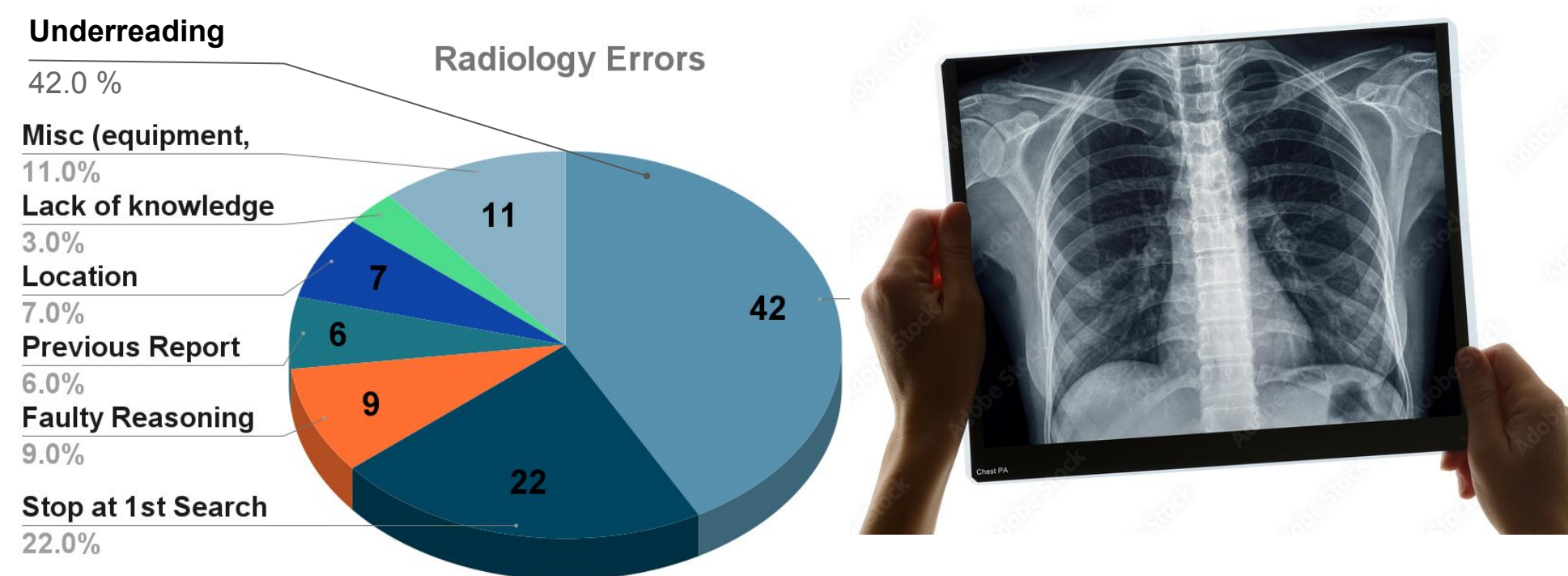
RepsNet: Combining Vision with Language for Automated Medical Reports

Ajay Tanwani, Joelle Barral, Daniel Freedman - Verily and Google Research



Introduction

- 1 billion radiology examinations are performed worldwide with **3-5 % error rate**
- Practitioners spend **5-10 minutes** for writing each report

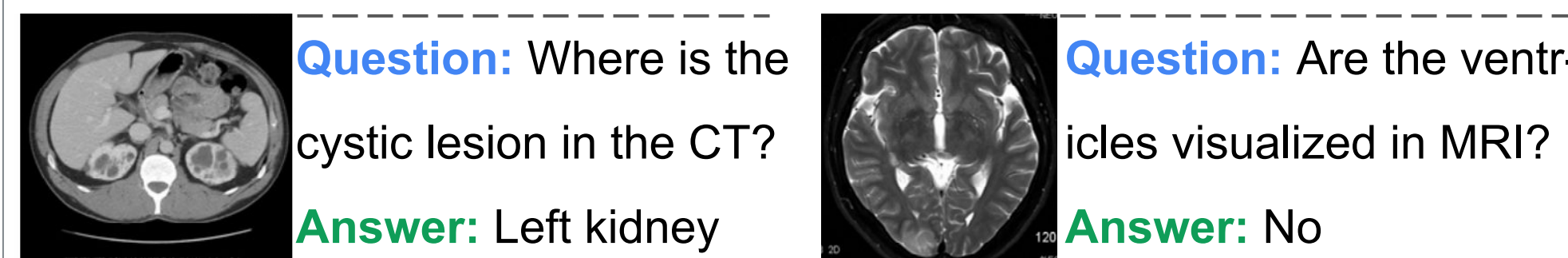


Objective

to improve **report accuracy** and save time in **writing and explaining reports to patients**, while **automating billing codes for describing medical findings**

VQA-Rad Results

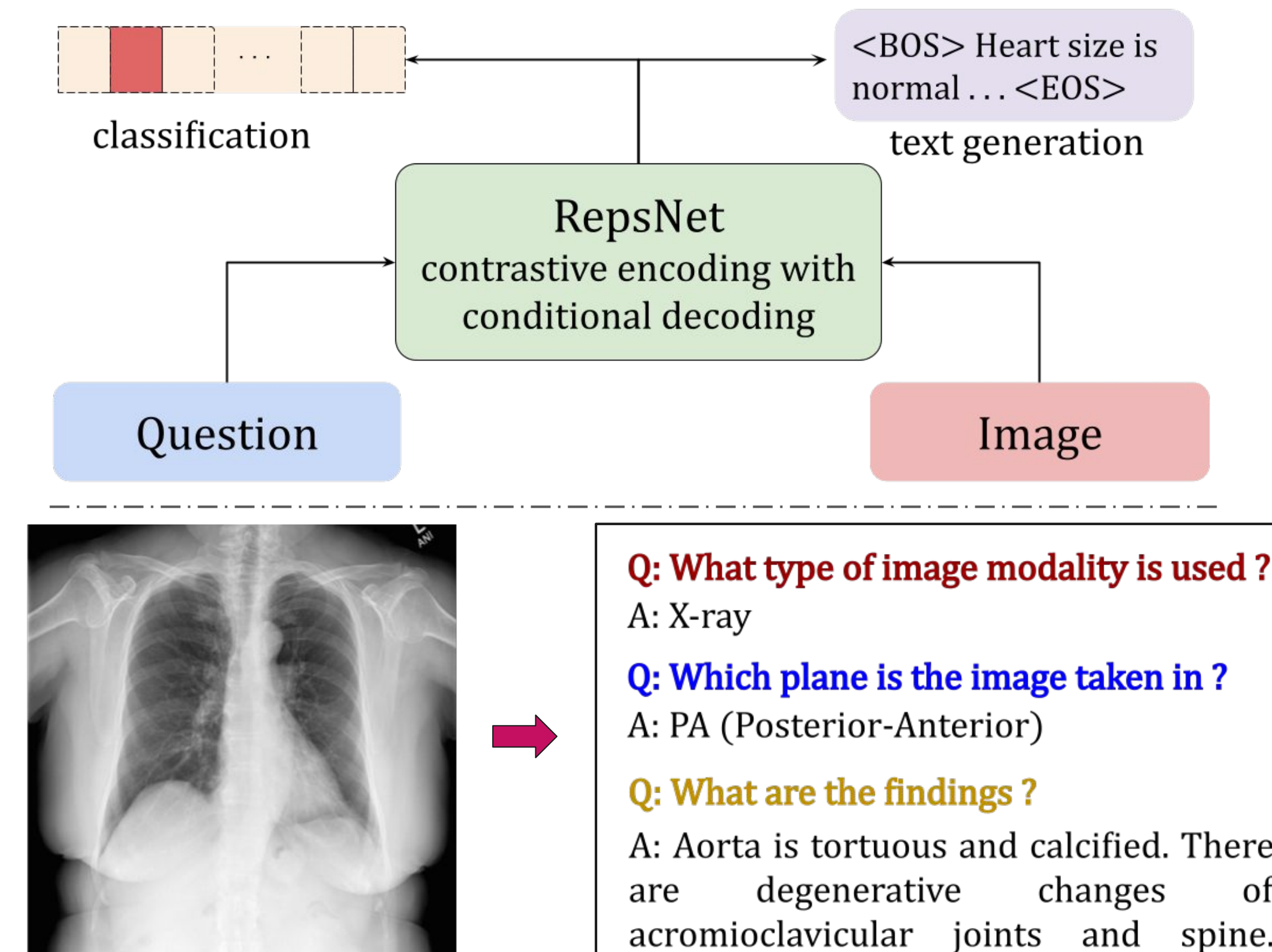
- **State-of-the-art classification accuracy** by adapting RepsNet on small **VQA-Rad datasets**



	2018	2019	All
MEVF [5]	66.10	—	—
MMQ [11]	67.00	—	—
QCR [41]	69.65	—	—
CLEF [1]	—	62.40	—
CRPD [21]	72.70	—	—
RepsNet-0	81.08	67.57	63.69
RepsNet-5	83.55	79.83	71.93
RepsNet-10	87.05	81.17	80.37

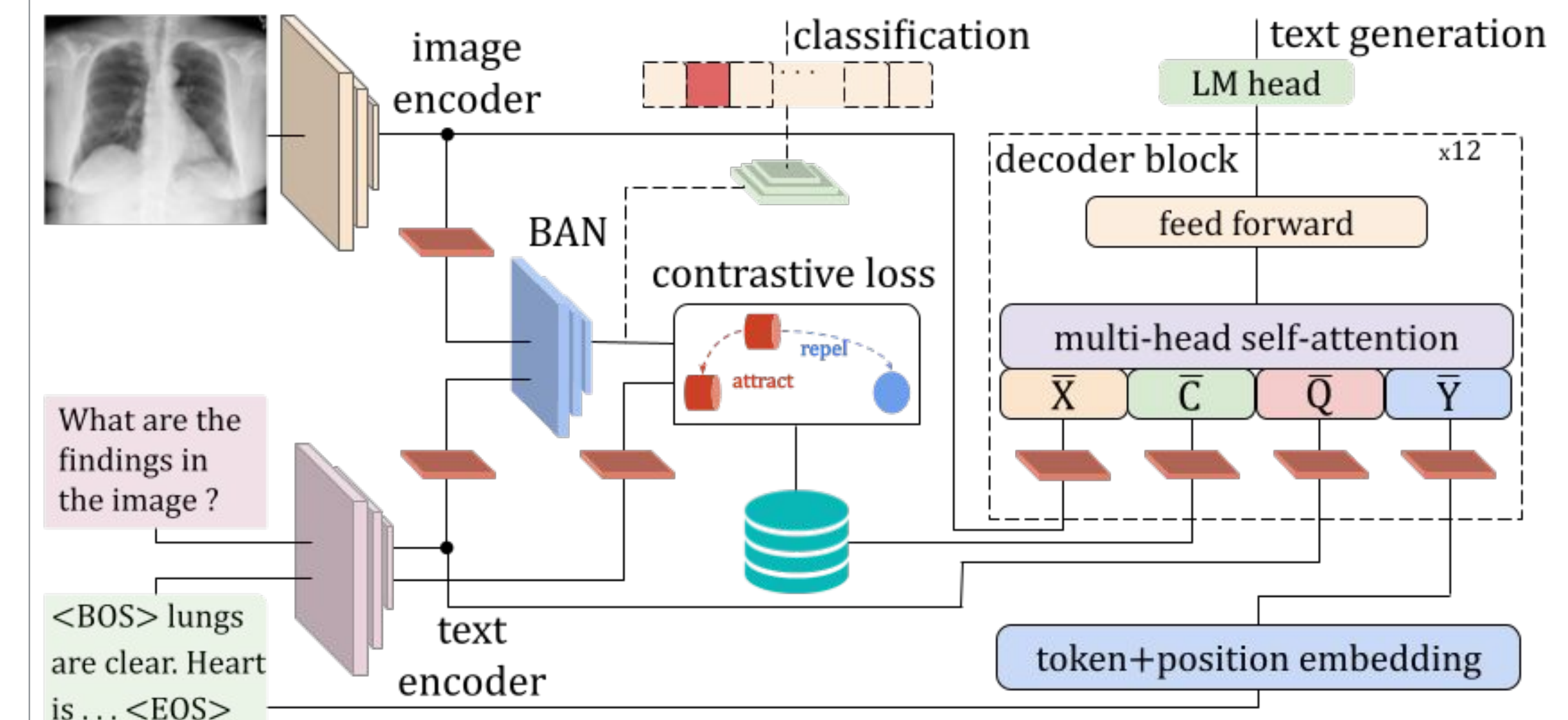
RepsNet Main Idea

- RepsNet fills medical reports by **classification** or **text generation** via **visual question answering**



RepsNet Model

- RepsNet adapts **pre-trained vision and language models**
- **Joint pre-training** of images and text in the **encoding phase by contrastive learning**
- **Augment language model with image and prior context of nearest neighbour answers** in the decoding phase



IU-XRay Results

- Joint pre-training of images and medical findings in the encoding phase, and conditional decoding of GPT-2 language model on image and nearest neighbour reports gives **state-of-the-art BLUE scores on IU-Xray dataset**

	B1	B2	B3	B4
Co-Att [16]	0.45	0.29	0.20	0.15
HRGR [20]	0.44	0.30	0.21	0.15
CMAS [15]	0.46	0.30	0.21	0.15
Mem-T [7]	0.47	0.30	0.22	0.16
VTI [25]	0.49	0.36	0.29	0.15
PPKED [22]	0.48	0.31	0.22	0.17
RepsNet	0.58	0.44	0.32	0.27

Qualitative Results & Next Steps

- RepsNet shows **strong alignment with ground-truth** in describing medical findings
- Ongoing work: **Grounding medical knowledge + video report generation + application to gastroenterology**

