# RaDialog: Large Vision-Language Models for X-Ray Reporting and Dialog-Driven Assistance

Chantal Pellegrini[1,3], Ege Özsoy[1,3], Benjamin Busam[1],
Benedikt Wiestler[2], Nassir Navab[1], and Matthias Keicher[1]

[1]Computer-Aided Medical Procedures and Augmented Reality, Technical University of Munich, Germany
[2]AI for Image-Guided Diagnosis and Therapy, Technical University of Munich, Germany
[3]Munich Center for Machine Learning (MCML), Garching, Germany.
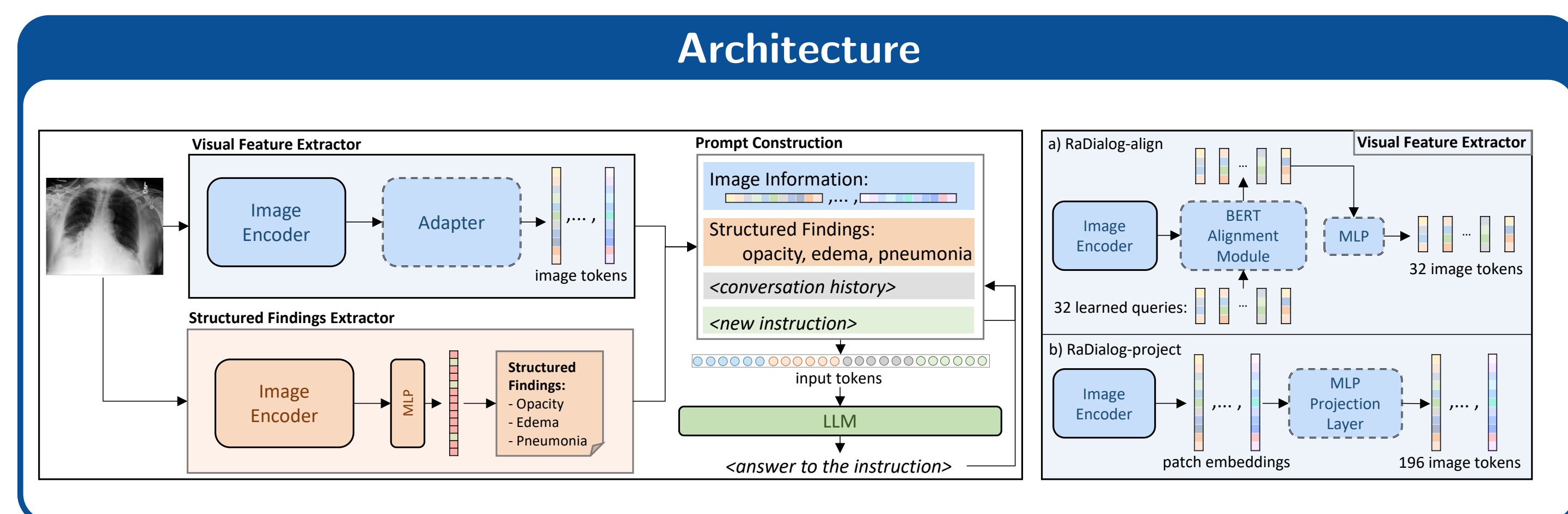
Technische Universität München

## Abstract

Conversational AI tools for generating and discussing accurate radiology reports could transform radiology by **enabling collaborative, human-in-the-loop diagnostic processes**, saving time, and enhancing report quality.

- **Main Contributions**:
  - Propose a novel **dual-branch architecture** incorporating structured clinical findings with image embeddings.
  - Design an **interactive instruct dataset** to combat the issue of catastrophic forgetting and enable dialog-based human-AI collaboration.
  - Introduce a **context-dropping augmentation** that improves the model's attention to the image information.
  - **In-depth evaluation** on report generation and conversational downstream tasks.
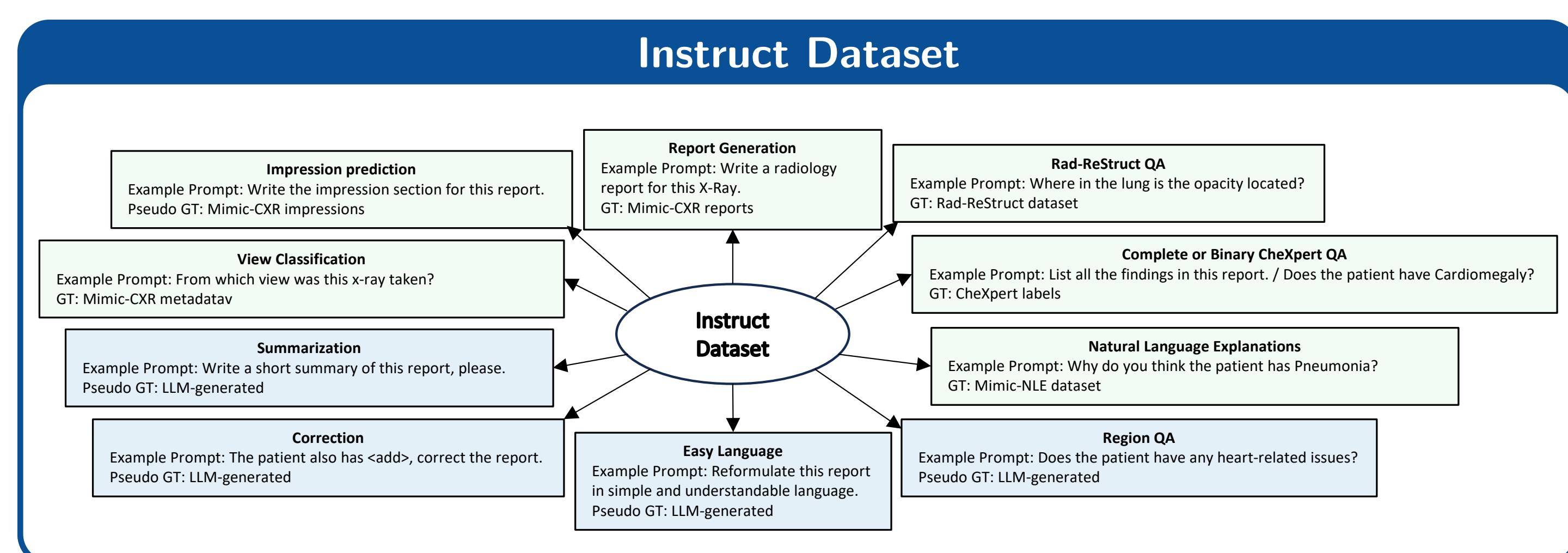
## Method

### Architecture



- **Visual Feature Extractor:** Extracts detailed visual encodings with a domain-specific image encoder and aligns them to the text space with the adapter module:
  - **RaDialog-align:** BERT-based alignment module, pre-trained on x-ray report pairs to extract 32 aligned image tokens
  - **RaDialog-project:** MLP projects all patch embeddings directly to language model tokens, trained end-to-end
- **Structured Findings Extractor:** Builds a structured representation of the main clinical findings in the image using a CLIP-based multi-label classifier, guiding the model on a higher level
- **Prompt Construction:** Image encoding, list of structured findings, conversation history, and instruction are converted into one prompt as input for the LLM
- **Language Model:** The LLM processes the prompt and produces an instruction-specific response. Fine-tuned using LoRA on x-ray report generation and auxiliary interactive downstream tasks.

### Instruct Dataset



- **Catastrophic Forgetting:** Training only on image-report pairs degrades the LLM's performance on other tasks significantly, so a diverse dataset is essential.
- **Instruct Dataset:** New dataset with 580k samples across 10 tasks, each with diverse prompts, to ensure model retains flexibility and conversation abilities:
  - **Core Diagnostic Tasks:** Based on ground-truth data from different datasets, e.g. Findings QA or View Classification.
  - **Replay tasks:** Pseudo ground-truth generated by the non-fine-tuned LLM to prevent forgetting, e.g. Correction or Easy Language
- **Context Dropping Augmentation:** Randomly includes full, partial, or no report in training samples to force the model to adapt to varying levels of report completeness, improving its ability to rely on image features alone for downstream tasks when necessary.

## Acknowledgements

## Results

### Report Generation



**Comparison to SOTA domain-specific VLLMs:**

| Method | | MIMIC-CXR | | | | IU-Xray (OOD) | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | FT | CE | BS | R-L | FT | CE | BS | R-L |
| LLaVA-Med | ✗ | 10.7 | 0.19 | 15.1 | ✗ | 5.0 | 0.20 | 15.8 |
| Rad-FM | ✓ | 15.4 | 0.22 | 15.6 | ✗ | 5.9 | 0.20 | 13.8 |
| XrayGPT | ✓ | 19.3 | 0.33 | 22.0 | ✓ | 9.9 | 0.39 | 25.7 |
| LLM-CXR | ✓ | 21.1 | - | - | - | - | - | - |
| CheXagent | ✓ | 22.2 | 0.36 | 25.9 | ✓ | 14.1 | 0.51 | 34.6 |
| R2GenGPT | ✓ | 24.7 | 0.36 | **27.6** | - | - | - | - |
| RaDialog_align-rep | ✓ | 39.4 | **0.40** | 27.1 | ✗ | 22.6 | **0.47** | **31.0** |
| RaDialog_align-ins | ✓ | 38.6 | 0.39 | 27.0 | ✗ | 22.9 | 0.46 | 30.2 |
| RaDialog_project-rep | ✓ | **39.7** | 0.36 | 25.6 | ✗ | 23.0 | 0.45 | 29.6 |
| RaDialog_project-ins | ✓ | 39.2 | 0.37 | 26.7 | ✗ | **23.1** | 0.45 | 30.4 |

**Comparison to SOTA report generation methods:**

| Method | CE | BS | B-4 | MTR | R-L |
| --- | --- | --- | --- | --- | --- |
| R2Gen (Chen et al., 2020) | 27.6 | 0.27 | 10.3 | 14.2 | 27.7 |
| Kiut (Wang et al., 2022) | 32.1 | - | 11.3 | 16.0 | 28.5 |
| COMG (Gu et al., 2024a) | 34.5 | - | 10.4 | 13.7 | 27.9 |
| HKRG (Wang et al., 2025) | 33.9 | - | **14.3** | **16.7** | **31.0** |
| ORID (Gu et al., 2024b) | 35.2 | - | 11.6 | 15.0 | 28.4 |
| MPO (Xiao et al., 2024) | 35.3 | - | 13.9 | 16.2 | 30.9 |
| RaDialog-project-instruct | **39.2** | **0.37** | 9.4 | 14.2 | 26.7 |

- Our architecture design leads to fine-grained image understanding and improved clinical correctness compared to prior methods.
- Further, RaDialog is the first LVLM to maintain SOTA report generation performance and outperforms all other methods significantly.

### Conversational Downstream Tasks



**Performance on Downstream Tasks:**

| Task | Metric | Comp. Method | Comp. Result | RaDialog |
| --- | --- | --- | --- | --- |
| Report Correction | CE improvement | XrayGPT | 10.0 | **33.4** |
| Finding Prediction | F1 | XrayGPT | 20.6 | **39.7** |
| Impression Generation | R-L | CheXagent | 40.3 | **45.8** |
| View Classification | Accuracy | CheXagent | **97.5** | 95.9 |

- Our instruct dataset and training setup successfully enables dialog-based assistance by avoiding catastrophic forgetting and teaching domain-specific conversational tasks.
- In a human evaluation comparing RaDialog against XRayGPT, the radiologist preferred RaDialog in 84% of cases for report generation and 71% in conversational tasks.

Code: Paper: