# A clinically accessible small multimodal radiology model and evaluation metric for chest X-ray findings

Juan Manuel Zambrano Chaves [1,2,7], Shih-Cheng Huang [2,7], Yanbo Xu[1,7], Hanwen Xu[3,7], Naoto Usuyama[1,7], Sheng Zhang[1,7], Fei Wang[4], Yujia Xie[1], Mahmoud Khademi[1], Ziyi Yang[1], Hany Awadalla[1], Julia Gong [1], Houdong Hu[1], Jianwei Yang[1], Chunyuan Li[1], Jianfeng Gao[1], Yu Gu[1], Cliff Wong[1], Mu Wei[1], Tristan Naumann [1], Muhao Chen [5], Matthew P. Lungren[1,2,6], Akshay Chaudhari [2], Serena Yeung-Levy [2], Curtis P. Langlotz [2], Sheng Wang [3] ✉ & Hoifung Poon [1] ✉
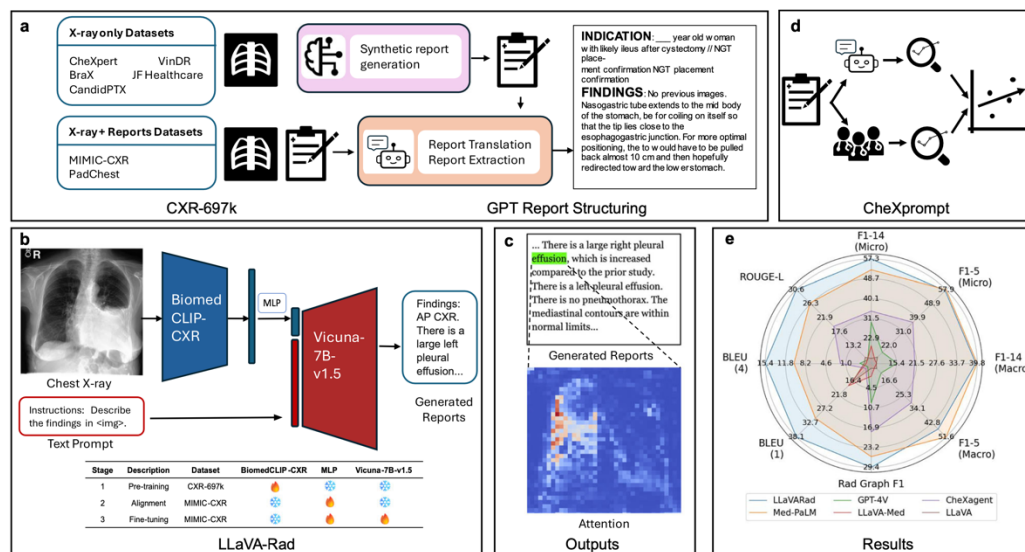
**Fig. 1 | LLaVA-Rad overview. a** To train LLaVA-Rad, we assemble a large dataset with over 697 thousand chest X-ray image-text pairs; GPT-4 is used to synthesize reports from labels, translate reports from Spanish, and process and structure the corresponding radiology reports. **b** We adopt a modular three-stage approach to train LLaVA-Rad, comprised of pre-training, alignment and fine-tuning. **c** A qualitative visualization of the model's attention during its generative process. **d** For evaluation, we also propose a novel factual error scoring approach using GPT-4 and demonstrate its parity with expert evaluation. **e** LLaVA-Rad outperforms much larger generalist and specialized models like GPT-4V and Med-PaLM M on prior standard report evaluation metrics. MLP multi-layer perceptron. The example chest X-ray image in **b** is obtained from ref. 27 with permission for reproduction from the authors.