

A Vision-Language Foundation Model to Enhance Efficiency of Chest X-ray Interpretation

Zhihong Chen^{1,2,*}, Maya Varma^{1,2,3,*}, Justin Xu^{1,2,4}, Magdalini Paschali^{1,2}, Dave Van Veen^{1,5}, Andrew Johnston², Alaa Youssef^{1,2}, Louis Blankemeier^{1,5}, Christian Bluethgen^{1,6}, Stephan Altmayer², Jeya Maria Jose Valanarasu^{1,3}, Mohamed Siddig Eltayeb Muneer², Eduardo Pontes Reis^{1,2}, Joseph Paul Cohen¹, Cameron Olsen², Tanishq Mathew Abraham⁷, Emily B. Tsai², Christopher F. Beaulieu², Jenia Jitsev^{8,9}, Sergios Gatidis^{1,2}, Jean-Benoit Delbrouck^{1,2}, Akshay S. Chaudhari^{1,2,10}, Curtis P. Langlotz^{1,2,10,11}

¹Stanford Center for Artificial Intelligence in Medicine and Imaging, Stanford University, Palo Alto, CA, USA.

²Department of Radiology, Stanford University, Stanford, CA, USA. ³Department of Computer Science, Stanford University, Stanford, CA, USA. ⁴Big Data Institute, University of Oxford, Oxford, UK. ⁵Department of Electrical Engineering, Stanford University, Stanford, CA, USA. ⁶Department of Radiology, University Hospital Zurich, Zürich, Switzerland. ⁷Stability AI, London, UK. ⁸Jülich Supercomputing Centre, Jülich, Germany. ⁹LAION, Germany.

¹⁰Department of Biomedical Data Science, Stanford University, Stanford, CA, USA. ¹¹Department of Medicine, Stanford University, Stanford, CA, USA. Corresponding to: {zhihongc,mvarma2,jbde,akshaysc,langlotz}@stanford.edu

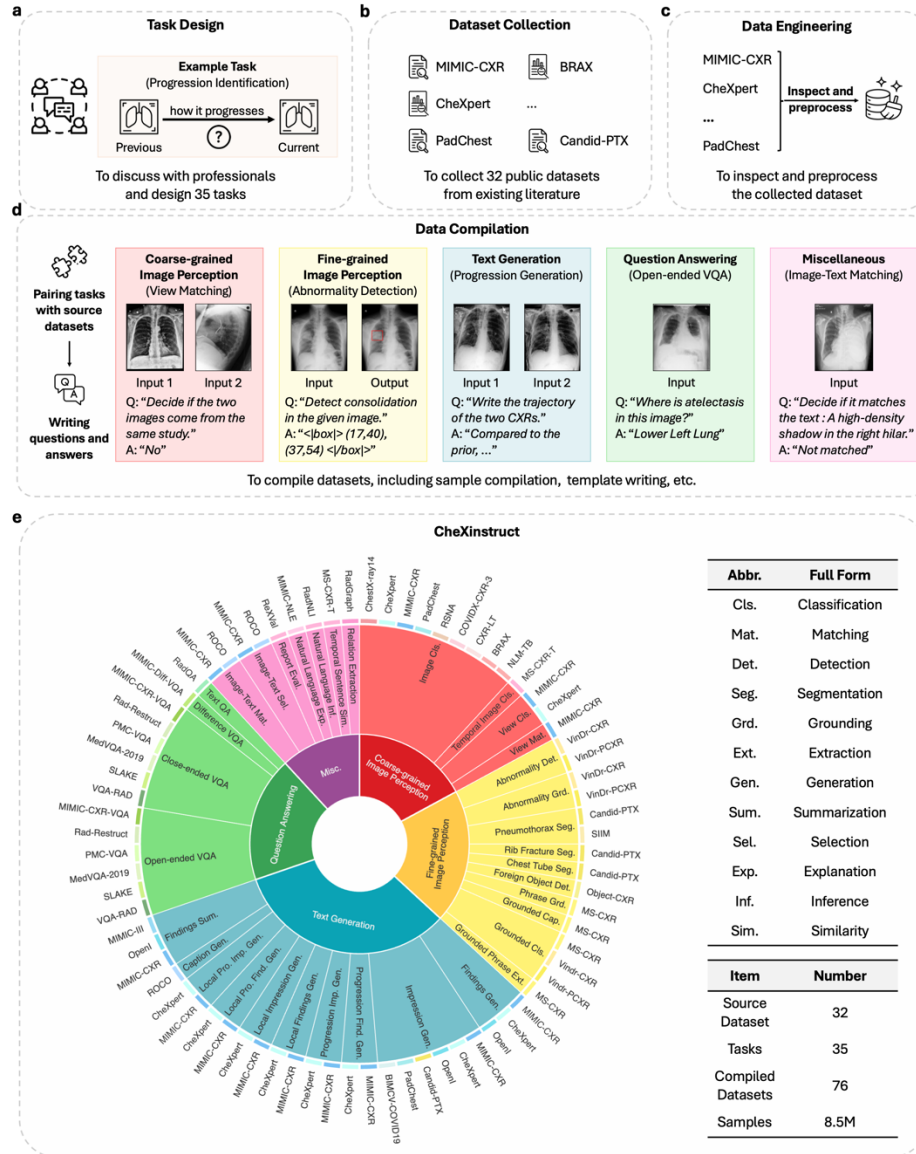


Figure 1 | Curation of CheXInstruct. a, Identification of CXR interpretation tasks. We defined 35 tasks that users are likely to perform with CXR FMs. b, Source dataset collection. To create training data samples for each of our defined tasks, we collected 32 public datasets. c, Data engineering. We performed both manual quality control and automated data engineering to preprocess collected source data. d, CheXInstruct compilation. We used the preprocessed datasets to generate training samples for each of our 35 defined tasks. e, Overview of CheXInstruct with data statistics.

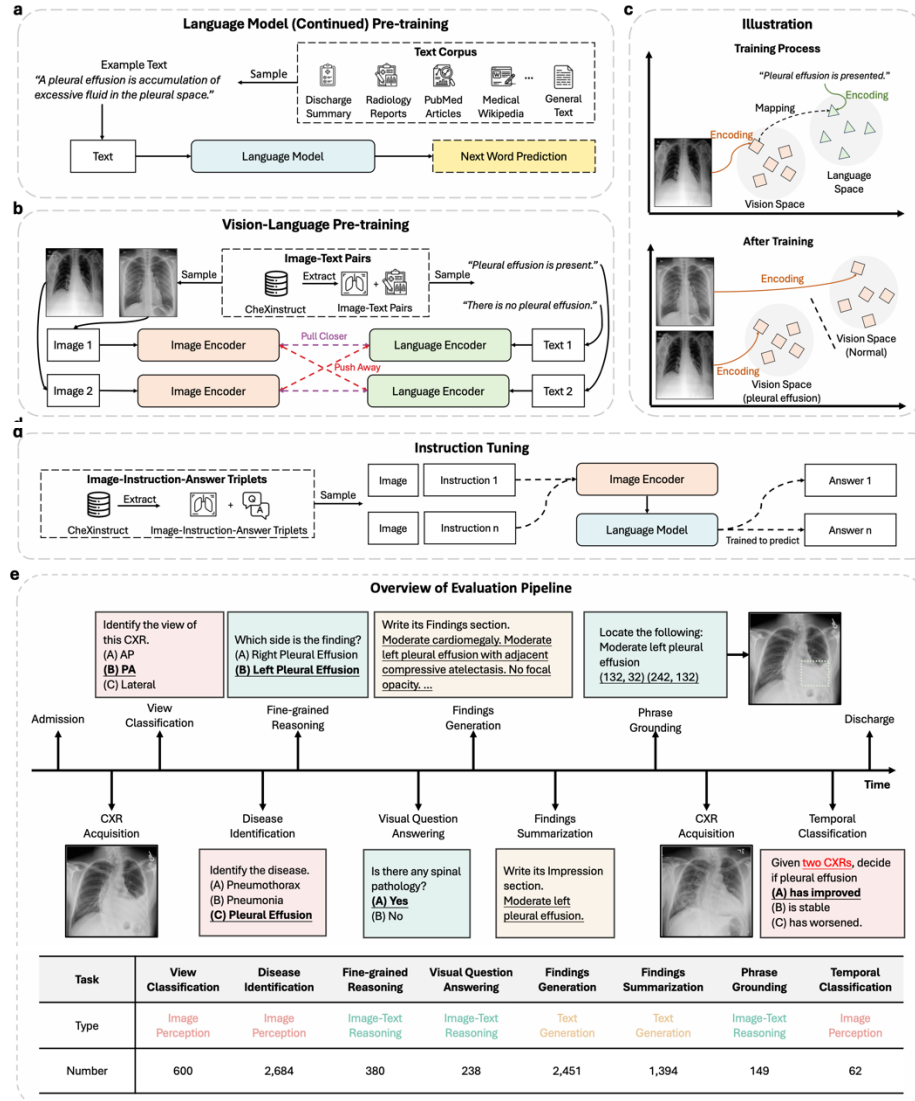


Figure 2 | Training and evaluating CheXagent. a, To develop CheXagent, we first trained a language model on clinical text. b, We then trained an image encoder to learn useful visual representations of imaging findings by leveraging paired text. c, This procedure enabled the visual encoder to capture semantic meaning with respect to key findings within its latent representation space. d, Finally, we jointly trained the image encoder and language model on data triplets from CheXInstruct, providing CheXagent with the capability to respond to user instructions. e, We constructed eight evaluation tasks to assess image perception, reasoning, and text generation capabilities.

