

Explaining Chest X-ray Pathologies in Natural Language

Maxime Kayser^{1*}, Cornelius Emde¹, Oana-Maria Camburu^{2**}, Guy Parsons^{1,3}, Bartłomiej Papież¹, and Thomas Lukasiewicz^{4,1}

¹ University of Oxford, United Kingdom

² University College London, United Kingdom

³ Thames Valley Deanery, Oxford, United Kingdom

⁴ TU Wien, Austria

Abstract. Most deep learning algorithms lack explanations for their predictions, which limits their deployment in clinical practice. Approaches to improve explainability, especially in medical imaging, have often been shown to convey limited information, be overly reassuring, or lack robustness. In this work, we introduce the task of generating natural language explanations (NLEs) to justify predictions made on medical images. NLEs are human-friendly and comprehensive, and enable the training of intrinsically explainable models. To this goal, we introduce MIMIC-NLE, the first, large-scale, medical imaging dataset with NLEs. It contains over 38,000 NLEs, which explain the presence of various thoracic pathologies and chest X-ray findings. We propose a general approach to solve the task and evaluate several architectures on this dataset, including via clinician assessment.

Keywords: Chest X-rays · Natural Language Explanations · XAI

1 Introduction

Deep learning (DL) has become the bedrock of modern computer vision algorithms. However, a major hurdle to adoption and regulatory approval of DL models in medical imaging is the lack of explanations for these models' predictions [21]. The combination of lack of model robustness [28], bias (algorithms are prone to amplifying inequalities that exist in the world) [9, 27], and the high stakes in clinical applications [26, 36] prevent black-box DL algorithms from being used in practice. In this work, we propose natural language explanations (NLEs) as a means to justify the predictions of medical imaging algorithms.

So far, the most commonly used form of explainability in medical imaging is saliency maps, which attribute importance weights to regions in an image. Saliency maps have many shortcomings, including being susceptible to adversarial attacks [8], conveying limited information and being prone to confirmation

* Corresponding author: firstname.lastname@cs.ox.ac.uk

** This work was mostly done while Oana was at the University of Oxford.

bias [1, 3], as well as only telling us *how much* highlighted regions affect the model’s output, and not *why* [7]. NLEs, on the other hand, would be able to fully capture how the evidence in a scan relates to the diagnosis. Furthermore, saliency maps are post-hoc explainers, i.e., they do not constrain the model to learn in an explainable manner. In contrast, self-explaining models have many benefits, including being more robust and having a better prediction performance [35]. Alternative approaches for explainability in medical imaging include latent space disentanglement [30], counterfactual explanations [33], case-based explanations [16], and concept-based explanations [18]. NLEs are a valuable addition to the suite of self-explaining models for medical imaging, as they provide easy-to-understand explanations that are able to communicate complex decision-making processes and mimic the way in which radiologists explain diagnoses [6, 25]. Previous attempts to augment medical image classification with textual information rely on template-generated sentences [6, 22] or sentences from other images based on image similarity [20]. Furthermore, their focus lies mostly on adding descriptive information about a pathology (e.g., its location), instead of explaining the diagnoses. In this work, we leverage the free-text nature of chest X-ray radiology reports, which do not only provide additional details about pathologies, but also the degree of certainty of a diagnosis, as well as justifications of how other observations explain it. Our work builds on the growing work on NLEs approaches in natural language processing [4, 19, 32], natural image understanding [15, 23, 24, 29], as well as task-oriented tasks such as self-driving cars [17] and fact-checking [19].

We propose MIMIC-NLE, the first dataset of NLEs in the medical domain. MIMIC-NLE extends the existing MIMIC-CXR dataset of chest X-rays [14] with diagnoses, evidence labels and NLEs for the diagnoses. We create MIMIC-NLE by using a BERT-based labeler, a set of clinical explanation keywords, and an empirically and clinically validated set of extraction rules. We extracted over 38,000 high-quality NLEs from the over 200,000 radiology reports present in MIMIC-CXR. Our extraction process introduces little noise, on-par or better than for NLE datasets in natural images [15]. Second, we establish an evaluation framework and compare three strong baselines. The evaluation by a clinician validates the feasibility of the task, but also shows that it is a challenging task requiring future research. The code and dataset are publicly available at <https://github.com/maximek3/MIMIC-NLE>.

2 MIMIC-NLE

Gathering NLEs is expensive, especially when radiologic expertise is required. To our knowledge, there is currently no NLE dataset for medical imaging. To address this, we show that it is possible to automatically distill NLEs from radiology reports, as radiologists typically explain their findings in these reports. We leverage the radiology reports from MIMIC-CXR [14], a publicly available chest X-ray dataset with 227,827 radiology reports. By applying various filters, labelers, and label hierarchies, we extract 38,003 *image-NLE* pairs, or 44,935

image-diagnosis-NLE triplets (as some NLEs explain multiple diagnoses). The extraction process is summarized in Fig. 1. Our filters consist mainly of removing sentences that contain anonymized data, or provide explanations based on patient history or technical details of the scan. More details are provided in the appendix. Furthermore, we only consider frontal, i.e., anteroposterior (AP) and posteroanterior (PA) scans, as these are most commonly used for diagnosis in routine clinical pathways and are most likely to contain the visual information required to generate NLEs.

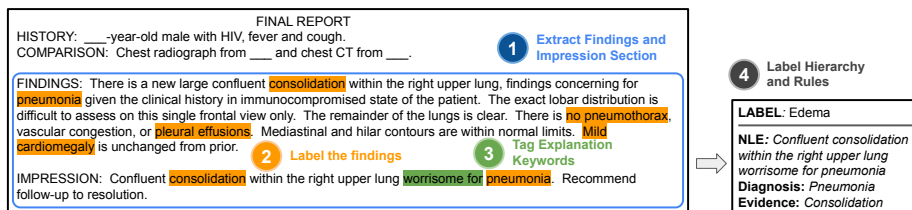


Fig. 1. The steps required to extract NLEs from raw radiology reports. We first extract the Findings and Impression sections, which contain the descriptive part of the report. Next, we identify the labels referred to in each sentence and the sentences that contain explanation keywords. Based on this information, we leverage the rules described in Table 1 to extract valid NLEs, as well as their diagnosis and evidence labels.

First, based on exploring the reports and discussions with clinicians, we observe that a small selection of phrases, such as “compatible with” or “worrisome for” (see full list in Appendix A), are very accurate identifiers of sentences where a potential pathology is explained by observations made on the scan (i.e., an *explanatory* sentence). Next, we make use of the CheXbert labeler [34], which can extract 14 different chest X-ray labels from clinical text, to identify the findings mentioned in each sentence. Thus, for each sentence in the MIMIC-CXR radiology reports, we know whether it is an explanatory sentence, which labels it refers to, and if the labels have a negative, positive, or uncertain (i.e., they are *maybe* present) mention. As the goal of NLEs is to explain predictions, we need to establish which of the labels mentioned in an NLE are being explained and which are part of the evidence. To be able to determine the evidence relationships present in NLEs, we need to restrict ourselves to a limited set of label combinations. For example, if an explanatory sentence in a radiology report indicates the presence of both *Atelectasis* and *Consolidation*, it is not obvious whether *Consolidation* is evidence for *Atelectasis* (e.g., “Right upper lobe new consolidation is compatible with atelectasis with possibly superimposed aspiration.”) or whether they are both explained by a different finding (e.g., “A persistent left retrocardiac density is again seen reflecting left lower lobe atelectasis or consolidation.”). However, for other label combinations, such as *Consolidation* and *Pneumonia*, the evidence relationship is usually clear, i.e., *Consolidation* is the evidence for *Pneumonia*.

We therefore propose the evidence graph in Fig. 2, which depicts the label relationships that have a known and high-confidence evidence relationship (i.e., their co-occurrence in a sentence lets us deduce an evidence relationship with high probability). The graph was constructed with an external radiologist and by empirically validating the co-occurrences of these labels in MIMIC-CXR.

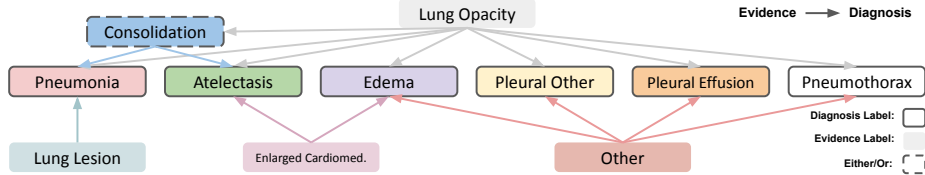


Fig. 2. Our evidence graph that visualizes which of the labels can act as evidence for which diagnosis labels.

Based on the evidence graph and manual inspection of every label combination that appears at least 25 times in an explanatory sentence (121 combinations in total), we established 12 mutually exclusive rules that mark a sentence as a valid NLE. The rules take into account the label combination, their uncertainty label (i.e., *Uncertain* or *Positive*), and the presence of explanation keywords. The rules are defined in Table 1. For each sentence in MIMIC-CXR, we use these rules to determine whether the label combination referred to in the sentence makes it a valid NLE and to determine which of the labels are part of the evidence and the diagnosis. Some label combinations, such as *Consolidation* and *Pneumonia*, are considered a valid NLE even if they do not contain an explanation keyword (as the keywords are not exhaustive). However, for label combinations with less frequent evidence relationships, such as *Enlarged Cardiomed.* and *Edema*, we also require the presence of an explanation keyword. It is worth noting that we focus on explanations for positive and uncertain cases only, as strictly negative findings will generally not require case-specific explanations.

Based on these rules, we obtain 38,003 NLEs. One of the authors evaluated a subset of 100 NLEs and found an accuracy of 92% (compared to between 49% to 91% for natural images NLE datasets [15]). We deemed an NLE as correct if it correctly uses the extracted evidence labels to explain the extracted diagnosis labels. The main reasons for incorrect NLEs were failures of the CheXbert labeler. Thus, with improved labelers, the accuracy could be further improved. Using the train, dev, and test splits in MIMIC-CXR [14], we get split sizes of 37,016, 273, and 714 for MIMIC-NLE, respectively.

3 Models

We propose self-explaining models that learn to detect lung conditions and explain their reasoning in natural language. The learned image representations are

constrained by mapping them to language that explains the evidence backing a diagnosis. Our approach is illustrated in Fig. 1.

We denote the vision model, i.e., the image classification model, as task model M_T . In our case, $M_T(x) = Y$, where x is a radiographic scan, and Y is the prediction vector $Y \in \mathbb{R}^{n_{\text{unc}} \times n_{\text{path}}}$, with $n_{\text{unc}} = 3$ being the number of certainty levels, and $n_{\text{path}} = 10$ being the number of pathologies. This follows the *U-MultiClass* approach from Irvin et al. [13], i.e., for each pathology we classify the image as *negative*, *uncertain*, or *positive*.

Table 1. This table denotes all the included label combinations for NLEs, including which of the labels are being explained and which are the evidence. The column “*kw req.*” specifies which label combinations additionally require the presence of an explanation keyword to be considered an NLE. “*Other / misc.*” refers to evidence that has not been picked up by the CheXbert labeler. If not denoted by U or P , all labels can be either positive or uncertain. A^U and B^U are the sets A and B , where all labels are given as uncertain. $\mathcal{P}_{\geq 2}(A^U)$ is the power set of A^U , where each set has at least two labels (i.e., any combination of at least two labels from A^U).

MIMIC-NLE Label Combinations		
Evidence	Diagnosis Label(s)	kw req.
<i>Other / misc.</i>	$d \in A = \{\text{Pleural Eff., Edema, Pleural Other, Pneumoth.}\}$	yes
<i>Other / misc.</i>	$s \in \mathcal{P}_{\geq 2}(A^U)$	yes
Lung Opacity	$d \in B = A \cup \{\text{Pneumonia, Atelectasis}\}$	no
Lung Opacity	$s \in \mathcal{P}_{\geq 2}(B^U)$	no
Lung Opacity	Consolidation	no
Consolidation	Pneumonia	no
{Lung Op., Cons.}	Pneumonia	no
Lung Lesion	Pneumonia	yes
Lung Opacity	$\{\text{Atelectasis}^P, \text{Pneumonia}^U\}$	no
Consolidation	$\{\text{Atelectasis}^U, \text{Pneumonia}^U\}$	no
Enlarged Card.	Edema	yes
Enlarged Card.	Atelectasis	yes

To our knowledge, this is the first application of NLEs to multi-label classification. We address this by generating an NLE for every label that was predicted as *uncertain* or *positive* and is considered a diagnosis label in our evidence graph. Given a set of pathologies P (see Fig. 2), we generate an NLE for every pathology $p_j \in P$. We denote the explanation generator as M_E . For every pathology p_j , we condition the NLE generation on the M_T prediction, i.e., we have $M_E(x_{\text{REP}}, Y, p_j) = e_j$, where x_{REP} is the learned representation of the image x , and e_j is the NLE that explains the classification of p_j . By backpropagating the loss of M_E through M_T , M_T is constrained to learn representations that embed the correct reasoning for each diagnosis. During training, we condition M_E on the ground-truth (GT) pathology p_j and prediction vector Y .

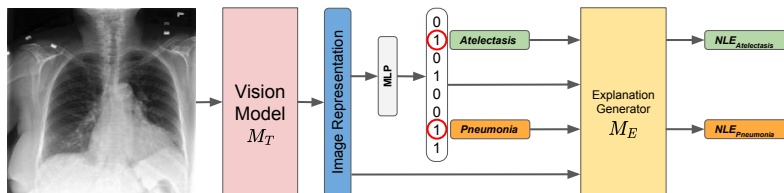


Fig. 3. The model pipeline to provide an NLE for a prediction.

4 Experimental Setup

We evaluate three baselines on our MIMIC-NLE dataset. We propose both automatic and expert-based evaluations.

Baselines. First, we adapt two state-of-the-art chest X-ray captioning architectures to follow the NLE generation approach outlined in Fig. 3. We re-implement TieNet [37], following their paper and a publicly available re-implementation⁵ and RATCHET [11], a recent transformer-based captioning method. For each approach, we use DenseNet-121 [12] as the vision model M_T . DenseNet-121 is a convolutional neural networks (CNN) widely used for chest X-ray classification [13]. For TieNet, contrary to the original model, we do not predict the labels from the learned text representations, but condition the NLE generation on the labels. TieNet’s decoder, an LSTM [10], is conditioned on the pooled features of the image x_{REP} via an attention layer, on the prediction vector Y by adding its embedding to x_{REP} , and on the diagnosis by initializing the hidden state of the decoder with an encoding of the predicted diagnosis p_j in textual form. RATCHET considers 49×1024 feature maps as x_{REP} , and the embedding of Y is modeled as an additional, 50-th, feature map. We condition on p_j in the same way as in TieNet, by initializing the hidden state of the decoder.

We also propose an additional baseline referred to as DPT (**D**enseNet-121 + **GPT**-2), which is inspired by state-of-the-art NLE approaches for natural image understanding [15, 24]. DPT leverages a DenseNet-121 as M_T and a vision-enabled GPT-2 language model [31] as M_E . An attention layer enforces M_E to focus on the same regions in the scan as M_T . We get attention-weighted $7 \times 7 \times 1024$ feature maps from the last layer of DenseNet-121. We consider the NLE generation as a sequence-to-sequence translation from the 49 feature maps (to which we add the relative position of each feature map as a position encoding), a token for Y , and the diagnosis p_j in text form. For fair comparison, we use the same GPT-2 vocabulary for all models and initialize all word embeddings with the pre-trained GPT-2 weights. All models use DenseNet-121 from TorchXRyVision [5], which was pre-trained on CheXpert [13]. All text generation is done via greedy decoding.

⁵ <https://github.com/farrell1236/RATCHET/tree/tienet>

Model Training. We use a maximum sequence length of 38 during training, which corresponds to the 98th percentile of the training set. We use the Adam optimizer with weight decay for the transformer models and without weight decay for TieNet (as it failed otherwise). We also use a linear scheduler with warmup for the learning rate. For each model, we experiment with different learning rates and batch sizes as hyperparameters. Using the dev set, we obtain the best hyperparameters as follows: for DPT, a learning rate of 5×10^{-4} and batch size of 16; for RATCHET, a batch size 16 and learning rate 5×10^{-5} ; for TieNet, a batch size of 32 and learning rate of 1×10^{-3} . We selected these based on the product of the task score S_T , CLEV score, and the average of the BERTScore [38] and METEOR [2] score, all outlined below.

Evaluation Metrics. We evaluate models both on their ability to solve the task, i.e., image multi-label classification, and their ability to explain how they solved the task, i.e., providing NLEs. The task score S_T is given by the weighted AUC score, where we consider *uncertain* and *positive* as one class (following Irvin et al. [13]). For the NLEs scores, we only consider the NLEs that explain correctly predicted labels, as in [4, 15]. Previous works have shown that automated natural language generation (NLG) metrics are underperforming for NLEs, as the same answers can be explained in different syntactic forms and even different semantic meanings [15]. We therefore propose the CLEV (CLinical EVIDence) score, which verifies whether an NLE refers to the right clinical evidence. For this, we leverage the CheXbert labeler, which extracts the evidence labels referred to in the GT and generated NLEs. For example, if the GT NLE mentions *Lung Opacity* and *Consolidation* as evidence for *Pneumonia*, we expect the generated NLE to contain the same findings. The CLEV score is the accuracy over all the generated NLEs, i.e., what share of them contains exactly the same evidence labels as the corresponding GT. We also provide the NLG metrics of BLEU, Rouge, CIDEr, SPICE, BERTScore, and METEOR as in [15]. For BERTScore, we initialize the weights with a clinical text pretrained BERT model.⁶ It is worth noting that, out of the given suite of NLG metrics, BERTScore and METEOR have previously been shown to have the highest (although still low in absolute value) empirical correlation to human judgment in natural image tasks [15].

Clinical Evaluation. Given the difficulties of evaluating NLEs with automatic NLG metrics, we also provide an assessment of the NLEs by a clinician. They were presented with 50 X-ray scans, a diagnosis to be explained, and four different NLEs: the GT and one for each of our three models. The NLEs are shuffled for every image, and the clinician does not know which is which. They are then asked to judge how well each NLE explains the diagnosis given the image on a Likert scale of 1 (*very bad*) to 5 (*very good*).

5 Results and Discussion

Table 5 contains the AUC task score, the evaluation conducted by a clinician, and the automatic NLG metrics. We also provide AUC results for DenseNet-121 only,

⁶ https://huggingface.co/emilyalsentzer/Bio_ClinicalBERT

Table 2. The S_T score, clinical evaluation, and NLG scores for our baselines on the MIMIC-NLE test set. \geq GT reflects the share of generated NLEs that received a rating on-par or better than the GT. Clin.Sc. reflects the average rating of 1 (lowest) to 5 (highest) that was given to the NLEs by a clinician. R-L refers to Rouge-L, and Bn to the n -gram BLEU scores. Best results are in bold. As we only evaluate NLEs for correctly predicted diagnoses, our NLG metrics cover 534, 560, and 490 explanations for RATCHET, TieNet, and DPT, respectively.

	AUC	\geq GT	Clin.Sc.	CLEV	BERTS.	MET.	B1	B4	R-L	CIDEr	SPICE
GT	-	-	3.20	-	-	-	-	-	-	-	-
DenseNet-121	65.2	-	-	-	-	-	-	-	-	-	-
RATCHET	66.4	48%	2.90	74.7	77.6	14.1	22.5	4.7	22.2	37.9	20.0
TieNet	64.6	40%	2.60	78.0	78.0	12.4	17.3	3.5	19.4	33.9	17.2
DPT	62.5	48%	2.66	74.9	77.3	11.3	17.5	2.4	15.4	17.4	13.7

i.e., without an M_E module, which shows that providing NLEs can improve task performance, e.g., for RATCHET. We observed that the main reason why GT NLEs obtain an absolute rating score of 3.2/5 is inter-annotator disagreement between our clinician and the author of the reports (sometimes due to lack of patient context given in this scenario). Another reason is that some of the GT NLEs refer to a change in pathology with respect to a previous study, which is not something that can be assessed from the image. We also observe that the CLEV score neither correlates to expert evaluation nor NLG metrics. One explanation could be that the evidence labels in MIMIC-NLE are highly imbalanced, i.e., predominantly *Lung Opacity*. Therefore, as the CLEV score does not take into account much of the diversity that is inherent in our NLEs, such as the location of findings, and their size and appearance, a model that generates generic NLEs that make reference to *Lung Opacity* will yield a good CLEV score, as was the case for TieNet. Overall, the NLG metrics are on-par with the NLG metrics for report generation [11]. Hence, the difficulty of generating longer texts (reports) seems to be offset by the degree of difficulty of specific, but shorter NLEs. The results also indicate that the NLG metrics are generally poor at reflecting expert judgment. More precisely, BERTScore, which showed the highest correlation with human judgment for natural images [15], has poor indicative qualities for medical NLEs. One reason could be that automatic NLG metrics generally put equal emphasis on most words, while certain keywords, such as the location of a finding, can make an NLE clinically wrong, but only contribute little to the NLG score. While the GPT-2 based architecture of DPT proved very efficient for natural images [15], its performance is less convincing on medical images. A reason could be that GPT-2 is too large and relies on embedded commonsense knowledge to generate NLEs, which is less helpful on highly specific medical text. Example NLE generations are provided in the appendix.

6 Summary and Outlook

In this work, we introduced MIMIC-NLE, the first dataset of NLEs in the medical domain. We proposed and validated three baselines on MIMIC-NLE. Providing NLEs for medical imaging is a challenging and worthwhile task that is far from being solved, and we hope our contribution paves the way for future work. Open tasks include providing robust automatic metrics for NLEs and introducing more medical NLE datasets and better performing NLE models.

Acknowledgments

We thank Sarim Ather for useful discussions and feedback. M.K. is supported by the EPSRC Center for Doctoral Training in Health Data Science (EP/S02428X/1), and by Elsevier BV. This work has been partially funded by the ERC (853489—DEXIM) and by the DFG (2064/1—Project number 390727645). This work has also been supported by the Alan Turing Institute under the EPSRC grant EP/N510129/1, by the AXA Research Fund, the ESRC grant “Unlocking the Potential of AI for English Law”, the EPSRC grant EP/R013667/1, and by the EU TAILOR grant. We also acknowledge the use of GPU computing support by Scan Computers International Ltd. BWP acknowledges a Nuffield Department of Population Health Research Fellowship. OMC acknowledges a Leverhulme Early Career Fellowship.

Bibliography

- [1] Adebayo, J., Gilmer, J., Muelly, M., Goodfellow, I., Hardt, M., Kim, B.: Sanity checks for saliency maps. In: *NeurIPS* (2018)
- [2] Banerjee, S., Lavie, A.: METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In: *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization* (2005)
- [3] Bornstein, A.M.: Is artificial intelligence permanently inscrutable? (2016)
- [4] Camburu, O.M., Rocktäschel, T., Lukasiewicz, T., Blunsom, P.: e-SNLI: Natural language inference with natural language explanations. In: *NeurIPS* (2018)
- [5] Cohen, J.P., et al.: TorchXRyVision: A library of chest X-ray datasets and models. *arXiv preprint arXiv:2111.00595* (2020)
- [6] Gale, W., et al.: Producing radiologist-quality reports for interpretable artificial intelligence. *arXiv preprint arXiv:1806.00340* (2018)
- [7] Ghassemi, M., Oakden-Rayner, L., Beam, A.L.: The false hope of current approaches to explainable artificial intelligence in health care. *The Lancet Digital Health* **3**(11) (2021)
- [8] Gu, J., Tresp, V.: Saliency methods for explaining adversarial attacks. *arXiv preprint arXiv: 1908.08413* (2019)
- [9] Hajian, S., Bonchi, F., Castillo, C.: Algorithmic bias: From discrimination discovery to fairness-aware data mining. In: *SIGKDD* (2016)
- [10] Hochreiter, S., Schmidhuber, J.: Long short-term memory. *Neural Computation* **9**(8) (1997)
- [11] Hou, B., Kaissis, G., Summers, R.M., Kainz, B.: RATCHET: Medical transformer for chest X-ray diagnosis and reporting. In: *MICCAI* (2021)
- [12] Huang, G., Liu, Z., van der Maaten, L., Weinberger, K.Q.: Densely connected convolutional networks. In: *CVPR* (2017)
- [13] Irvin, J., et al.: CheXpert: A large chest radiograph dataset with uncertainty labels and expert comparison **33**(01) (2019)
- [14] Johnson, A.E.W., et al.: MIMIC-CXR-JPG, a large publicly available database of labeled chest radiographs. *arXiv:1901.07042* (2019)
- [15] Kayser, M., Camburu, O.M., Salewski, L., Emde, C., Do, V., Akata, Z., Lukasiewicz, T.: e-ViL: A dataset and benchmark for natural language explanations in vision-language tasks. In: *ICCV* (2021)
- [16] Kim, E., Kim, S., Seo, M., Yoon, S.: XProtoNet: Diagnosis in chest radiography with global and local explanations. In: *CVPR* (2021)
- [17] Kim, J., Rohrbach, A., Darrell, T., Canny, J., Akata, Z.: Textual Explanations for Self-Driving Vehicles. In: *ECCV* (2018)
- [18] Koh, P.W., Nguyen, T., Tang, Y.S., Mussmann, S., Pierson, E., Kim, B., Liang, P.: Concept bottleneck models. In: *ICML* (2020)
- [19] Kotonya, N., Toni, F.: Explainable automated fact-checking for public health claims. In: *EMNLP* (2020)

- [20] Kougia, V., et al.: RTEEX: A novel framework for ranking, tagging, and explanatory diagnostic captioning of radiography exams. *Journal of the American Medical Informatics Association* **28**(8) (2021)
- [21] Langlotz, C.P., et al.: A roadmap for foundational research on artificial intelligence in medical imaging: From the 2018 NIH/RSNA/ACR/The Academy Workshop. *Radiology* **291**(3) (2019)
- [22] Lee, H., Kim, S.T., Ro, Y.M.: Generation of multimodal justification using visual word constraint model for explainable computer-aided diagnosis. In: *Interpretability of Machine Intelligence in Medical Image Computing and Multimodal Learning for Clinical Decision Support* (2019)
- [23] Majumder, B.P., Camburu, O.M., Lukasiewicz, T., McAuley, J.: Knowledge-grounded self-rationalization via extractive and natural language explanations. In: *ICML* (2022)
- [24] Marasović, A., et al.: Natural language rationales with full-stack visual reasoning: From pixels to semantic frames to commonsense graphs. In: *EMNLP Findings* (2020)
- [25] Miller, T.: Explanation in artificial intelligence: Insights from the social sciences. *Artificial Intelligence* **267** (2019)
- [26] Mozaffari-Kermani, M., et al.: Systematic Poisoning Attacks on and Defenses for Machine Learning in Healthcare. *IEEE Journal of Biomedical and Health Informatics* **19**(6) (2015)
- [27] Obermeyer, Z., et al.: Dissecting racial bias in an algorithm used to manage the health of populations. *Science* **366**(6464) (2019)
- [28] Papernot, N., McDaniel, P., Goodfellow, I.: Transferability in machine learning: From phenomena to black-box attacks using adversarial samples. *arXiv preprint arXiv:1605.07277* (2016)
- [29] Park, D.H., et al.: Multimodal explanations: Justifying decisions and pointing to the evidence. In: *ICCV* (2018)
- [30] Puyol-Antón, E., et al.: Interpretable deep models for cardiac resynchronization therapy response prediction. In: *MICCAI* (2020)
- [31] Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I.: Language Models are Unsupervised Multitask Learners. *openai.com* (2019)
- [32] Rajani, N.F., McCann, B., Xiong, C., Socher, R.: Explain yourself! Leveraging language models for commonsense reasoning. In: *ACL* (2019)
- [33] Schutte, K., Moindrot, O., Hérent, P., Schiratti, J.B., Jégou, S.: Using StyleGAN for visual interpretability of deep learning models on medical images. *arXiv preprint arXiv:2101.07563* (2021)
- [34] Smit, A., et al.: Combining automatic labelers and expert annotations for accurate radiology report labeling using BERT. In: *EMNLP* (2020)
- [35] Stacey, J., Belinkov, Y., Rei, M.: Natural language inference with a human touch: Using human explanations to guide model attention. *AAAI* (2022)
- [36] Vayena, E., Blasimme, A., Cohen, I.G.: Machine learning in medicine: Addressing ethical challenges. *PLoS medicine* **15**(11) (2018)
- [37] Wang, X., et al.: TieNet: Text-image embedding network for common thorax disease classification and reporting in chest X-rays. In: *CVPR* (2018)
- [38] Zhang, T., Kishore, V., Wu, F., Weinberger, K.Q., Artzi, Y.: BERTScore: Evaluating text generation with BERT. In: *ICLR* (2019)

A Supplementary

Table 3. Our different filters that were applied to extract image-diagnosis-NLE triplets from MIMIC-CXR. # of sentences corresponds to unique image-sentence pairs. “Non-descriptive aspects” are aspects that cannot be derived from the image itself. Duplicates are sentences that are from the same report and mention the same labels.

Filter	# of sentences
Extract Findings and Impression sections	1,383,533
Remove sentences with anonymized data	1,304,465
Remove sentences referring to non-descriptive aspects: Patient history: “prior”, “compare”, “change”, “deteriorat”, “increase”, “decrease”, “previous”, “patient” Recommendations: “recommend”, “perform”, “follow” Technical: “CT”, “technique”, “ position”, “exam”, “assess”, “view”, “imag”	1,007,002
Filter by rules from Table 1	43,612
Remove duplicates	39,094
Remove studies without AP or PA images	38,003

Table 4. The distribution of diagnosis labels in all our NLEs. Diagnosis label combinations are ordered by occurrence. # of sentences corresponds to unique image-sentence pairs. The table is displayed in two halves.

Diagnosis Labels	# of sentences	Diagnosis Labels	# of sentences
Atelectasis	10,616	Atel., Edema, Pneumonia	104
Pneumonia	9,032	Atelectasis, Edema	103
Edema	5,098	Pl. Eff., Pneumothorax	65
Atelectasis, Pneumonia	4,773	Pleural Effusion, Pneumonia	55
Pleural Effusion	4,585	Edema, Pleural Effusion	31
Consolidation	916	Atelectasis, Pleural Other	17
Pleural Other	846	Atel., Edema, Pl. Eff.	9
Edema, Pneumonia	623	Edema, Pl. Eff., Pneumonia	8
Atelectasis, Pleural Effusion	397	Atel., Edema, Pl. Eff., Pneum.	6
Pneumothorax	311	Edema, Pleural Other	6
Pleural Effusion, Pleural Other	195	Pleural Other, Pneumonia	4
Atel., Pl. Eff., Pneumonia	190	Other	13

Table 5. The distribution of evidence labels in all our NLEs. Evidence label combinations are ordered by occurrence. # of sentences corresponds to unique image-sentence pairs. The table is displayed in two halves.

Evidence Labels	# of sentences	Evidence Labels	# of sentences
Lung Opacity	29,115	Consolidation, Lung Opacity	692
Other / misc.	5,842	Enlarged Cardiomeastinum	141
Consolidation	2,102	Lung Lesion	111

Table 6. The occurrence of different explanation keywords in the unfiltered sentences from MIMIC-CXR. The table is displayed in two halves.

Keyword	Count	Keyword	Count
suggest	19,787	relate	5,241
reflect	17,878	may represent	5,141
due	17,771	potentially	3,160
consistent with	14,654	worrisome for	1,685
concerning for	10,243	indicate	1,568
compatible with	6,025	account	1,193
likely represent	5,239	suspicious for	684

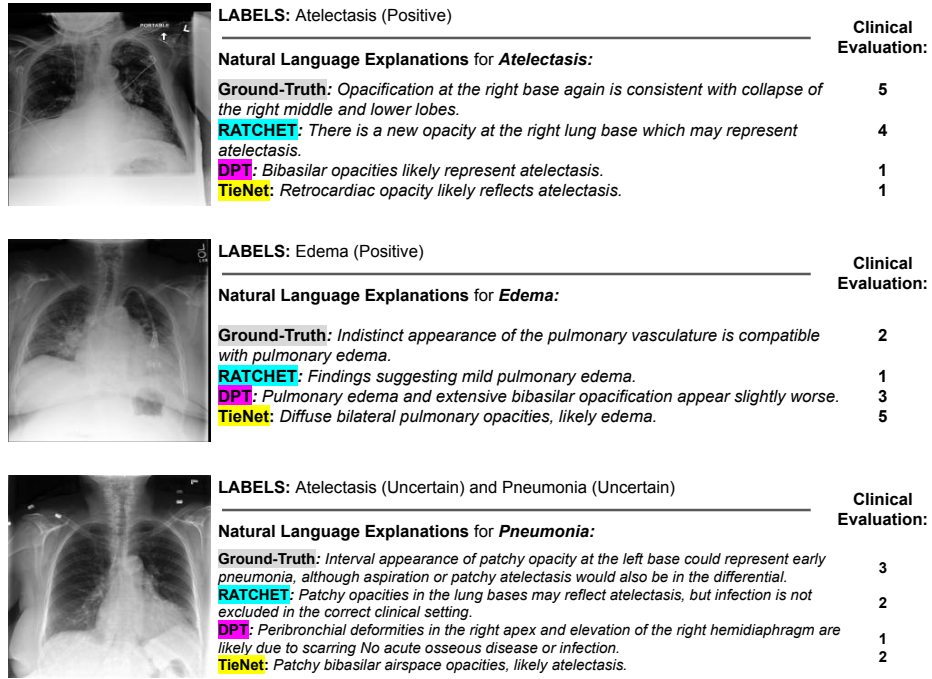


Fig. 4. The GT NLEs and three model-generated NLEs explaining three different diagnosis labels on three different scans. The clinical evaluation is given on a Likert scale, where 5 is the highest, and 1 is the lowest score.