

RadCBM: Hierarchical Concept Bottleneck Models with Automated Annotations for Chest X-ray Interpretation

Obadah Habash, Rabeb Mizouni, Shakti Singh, and Hadi Otrouk
Khalifa University, Abu Dhabi, UAE

Abstract—Chest X-ray classifiers remain hard to trust because their reasoning is hidden and the dense concept supervision they would need for transparent explanations is prohibitively expensive to obtain at scale. RadCBM is a hierarchical concept bottleneck model that replaces manual per-image concept labels with concept targets mined directly from paired radiology reports. We extract concepts with RadGraph, normalize them to RadLex/UMLS terms, filter to clinically meaningful semantic types, and organize them into anatomy-level regions. This yields a clinically grounded vocabulary of visual, image-evident concepts spanning lung, heart, pleura, mediastinum, and bone. A two-level predictor first estimates region abnormality, then predicts region-specific findings gated by those regions; a linear label head maps the gated concepts to diagnostic labels, so that each prediction decomposes into contributions from named regions and findings and supports direct counterfactual editing. On MIMIC-CXR and CheXpert, automated annotations cover the long tail of radiographic findings without human curation, and the hierarchical CBM improves concept AUC and reduces implausible activations compared to flat CBMs and CheXpert-style concept sets. Classification performance matches a black-box DenseNet-121 while exposing per-region rationales whose intervention effects are faithful to the learned decision boundary. RadCBM turns routine reports into training signals, aligning model decisions with radiologist workflows and providing region-aware, concept-level explanations without sacrificing accuracy.

I. INTRODUCTION

Deep models for chest X-ray interpretation now match radiologist-level accuracy on CheXpert and MIMIC-CXR [10], [3]. They are increasingly deployed for triage, second reading, and quality assurance, yet their decision processes remain opaque [1]. This opacity has concrete clinical consequences, affecting how models are audited, trusted, and integrated into workflow. When a network predicts “pneumonia,” clinicians cannot tell whether it responded to a genuine parenchymal finding or to confounders like support devices and acquisition artifacts. Failure analysis teams lack actionable signals for understanding systematic errors. Empirical work shows that chest X-ray models can generalize unpredictably across hospitals [19], while clinician interviews and position papers argue that high-stakes medical decisions require transparent, interpretable models rather than unexplained black boxes [17], [13]. In this work, our goal is not to eliminate distribution shift, but to make such failures more diagnosable by exposing region- and concept-level reasoning; we focus on in-domain performance on MIMIC-CXR and CheXpert. Radiologists, by contrast, describe findings as localized concepts; for

example, “left lower lobe opacity with no pleural effusion,” tying observations to anatomy and differential diagnoses. A trustworthy system should expose this intermediate reasoning, not just an uninterpretable probability [9], [4]. Together, these results and perspectives motivate explainability methods for chest X-ray classifiers that make explicit which image regions and clinical concepts drive each decision [10], [3], [1], [19].

Post-hoc explanation methods attempt to alleviate this opacity by highlighting pixels that most influence a model’s output or by fitting local surrogates around each prediction. Saliency and gradient-based techniques such as Grad-CAM and integrated gradients visualize where a network is “looking” [14], [16], while perturbation-based approaches such as LIME fit simpler models that approximate the decision boundary near a given image [12]. In chest X-ray applications these tools can help detect obvious failure modes (for example, models that rely on laterality markers or support devices) but they typically operate at the level of diffuse blobs rather than named clinical findings, and their outputs can be unstable across small input or parameter changes. Moreover, post-hoc tools treat the network as a fixed black box and do not constrain its internal representation to align with named radiologic concepts, which limits how much trust and control clinicians can derive from them in practice [14], [16], [12], [19].

Concept-based explanations aim to bridge this gap by expressing predictions in terms of high-level attributes rather than raw pixels, an idea instantiated by TCAV-style sensitivity analyses and concept bottleneck models [7], [8], [2], [6]. TCAV-style approaches test the sensitivity of a model to user-defined concept vectors, and concept bottleneck models (CBMs) directly embed this philosophy by routing predictions through a layer of human-interpretable concepts before computing labels. When concept activations are accurate, CBMs provide faithful, editable explanations because each class logit decomposes into contributions from named findings, and clinicians can intervene by editing concepts rather than raw pixels. However, classical CBMs assume that each training image comes with dense concept annotations, and the annotation burden grows with the product of dataset size and number of concepts. Typical chest X-ray datasets contain more than 10^5 studies and 50 to 200 relevant findings, so naively collecting expert concept labels would require millions of radiologist judgments and careful quality control [5], [3]. Recent CBM work explores better attribution [2], robustness under concept shift [6], and sparse autoencoders that carve pretrained representations into reusable

features [11], but these directions still depend on curated concept sets or small-scale supervision, and they rarely tackle the long tail of infrequent but clinically important findings that appear in real-world radiology practice.

A large fraction of chest X-ray reports describes the very region-specific findings that CBMs aim to predict, but reports also contain clinical impressions and non-visual statements (for example, indications and management plans) that cannot be inferred from images [4], [15]. A single sentence such as “patchy opacity in the left lower lobe consistent with pneumonia. No pleural effusion. Heart size is normal.” implicitly labels the presence of lung opacity and pneumonia and the absence of effusion and cardiomegaly for that study, while leaving other structures unspecified. If these signals can be extracted, normalized to standardized vocabularies, and grouped by anatomy, they could supervise CBMs at corpus scale without manual per-concept labeling, provided that non-visual or temporal concepts are filtered out and uncertainty in reports is handled explicitly. The core challenge becomes turning noisy report parses into an ontology-grounded, hierarchical concept vocabulary that supports region-aware explanations. This vocabulary should keep only visual, image-evident findings while maintaining high coverage of common and rare pathologies.

Existing chest X-ray explainability pipelines already tap into reports but stop short of making concepts the decision bottleneck. CheXagent produces long-form rationales by aligning report sentences with image regions [18], and tools such as RadGraph [4] and CheXbert [15] extract structured observation-anatomy pairs with uncertainty tags from free text. These systems demonstrate that reports provide rich supervisory signal, yet they mainly supervise black-box classifiers or post-hoc rationalizers whose internal pathways remain opaque and are not constrained to use only visual, region-grounded concepts, so their explanations can disagree with the actual features that drive the prediction. Some chest X-ray classifiers incorporate label hierarchies or anatomy-aware heads, but they typically use these structures as regularizers on latent features: predictions are still made directly from opaque embeddings rather than explicit concepts. In contrast, RadCBM constrains all labels to be linear functions of gated region and finding activations, so explanations are faithful to the decision path and directly editable.

Compared to prior chest X-ray pipelines that use reports only to derive global labels or generate rationales, RadCBM turns report-derived observation-region pairs into the primary representational bottleneck: we automatically construct a large, ontology-anchored vocabulary of region-specific findings and require all label predictions to flow through this vocabulary. This goes beyond simply training a multi-task model on RadGraph/CheXbert labels or attaching a generic CBM on top of a small curated concept set, because every prediction decomposes into contributions from clinically meaningful, region-specific concepts that can be directly inspected and edited.

We introduce RadCBM, a hierarchical CBM trained end-to-end from paired chest X-rays and reports without manual concept labels. RadCBM can be viewed as instantiating the CBM paradigm in the chest X-ray setting while addressing the

scalability limitations of prior concept-based approaches, by learning a clinically grounded concept space directly from reports instead of relying on hand-designed concept sets or small-scale annotations [8], [7], [2], [6]. RadCBM uses RadGraph to extract observation-anatomy pairs, normalizes them to RadLex/UMLS concepts, anchoring the vocabulary in established clinical ontologies, filters them to chest-focused semantic types and to concepts that are visually testable, and organizes the resulting vocabulary into anatomical regions with associated findings. A two-level predictor maps images to region abnormality scores and region-specific finding probabilities, applies multiplicative gating so that findings activate only when their region is abnormal, and feeds the gated concept vector to a transparent linear label head. Fixing the hierarchy to anatomy \rightarrow finding mirrors radiologist workflows and avoids learned groupings that may be difficult to interpret, while the linear head keeps concept contributions directly readable and editable and is designed to reduce the influence of non-visual report artifacts (for example, mentions of lines or tubes) on predictions. This architecture mirrors radiologists’ workflow (assess regions, describe findings, decide labels) and produces explanations that can be audited or edited at both region and concept levels.

Taken together, these design choices move chest X-ray explainability beyond generic saliency maps and small-scale concept probes toward a clinically grounded, editable reasoning process that mirrors radiologist workflows and scales to large, heterogeneous datasets [14], [16], [12], [10], [3], [1], [19], [8], [7], [2], [6].

Our main contributions are:

- An automated concept annotation pipeline that converts free-text radiology reports into soft, ontology-grounded concept targets aligned with standardized vocabularies, eliminating manual per-concept labeling while retaining a broad, clinically meaningful vocabulary of visual findings and anatomical regions.
- A hierarchical CBM with multiplicative gating between region abnormality and findings, enforcing clinical consistency (findings fire only in abnormal regions) and enabling region-first, finding-level explanations and counterfactual concept interventions.
- An empirical study on MIMIC-CXR and CheXpert comparing RadCBM to black-box CNNs, flat and post-hoc CBMs, and CheXpert-style concept sets, evaluating label performance (AUC/F1) on CheXpert findings, concept fidelity and intervention faithfulness, and plausibility and coverage of region-level explanations; automated concepts retain broad coverage, the hierarchy improves concept fidelity and reduces implausible activations (for example, subtle pneumothorax or early pneumomediastinum), and classification performance remains comparable to a black-box CNN while exposing actionable, region-aware explanations.

II. RELATED WORK

Post-hoc interpretability methods, such as saliency maps and feature attributions, highlight image regions that influence a

model’s output. Gradient-based methods (for example, Grad-CAM and integrated gradients) preserve black-box accuracy but often fail quantitative faithfulness checks, are sensitive to small perturbations, and can highlight broad regions that do not correspond to specific radiographic findings. Counterfactual explanations improve causal grounding by asking how predictions change when inputs are perturbed, but they still operate after training, do not constrain the internal representation, and can be difficult to interpret when the perturbations are not phrased in clinically meaningful concepts.

Concept bottleneck models make interpretability part of the architecture: images are mapped to human-understandable concepts that are then used to predict labels. Early CBMs relied on manually annotated concepts and often traded accuracy for transparency because the bottleneck limited capacity and annotation noise propagated directly into predictions. Recent variants predict concepts post-hoc or with weak supervision, but they either depend on small vocabularies (for example, the CheXpert labeler with 14 global labels) or on vision-language models that are not calibrated for clinical use and may hallucinate findings. Flat CBMs also ignore the anatomy-first reasoning radiologists employ: they treat all concepts as exchangeable, leading to implausible activations such as lung findings firing when the lungs are predicted normal or cardiac findings co-activating in obviously normal hearts.

Radiology report mining offers a scalable alternative by turning existing clinical text into structured supervision. Tools like RadGraph extract observation and anatomy entities plus relations from free text, while labelers such as CheXpert and CheXbert map reports into global disease labels with uncertainty tags. Prior work has used report-derived labels to supervise black-box classifiers or to train report-generation models, but the labels are coarse and do not expose the intermediate reasoning process or enforce region-level consistency. RadCBM combines report-derived concepts with a hierarchical, gated CBM so that explanations remain both faithful to the model (because concepts lie on the decision path) and aligned with clinical reading workflows that proceed from regions to specific findings.

III. METHOD

A. Pipeline Overview

Figure 1 summarizes the training pipeline. During training, paired chest X-rays and reports are processed jointly: reports are converted into soft concept targets, and images provide pixel-level evidence for those concepts and the final diagnostic labels. Images pass through a convolutional backbone to produce features that feed a two-level CBM with region abnormality heads, region-specific finding heads, multiplicative gating, and a linear label predictor. Reports are only used to build and supervise the concept layer; at test time, RadCBM receives only images and outputs region scores, concept activations, and label predictions.

B. Concept Extraction from Reports

Each report is parsed with RadGraph to identify observation and anatomy entities plus their relations (for example,

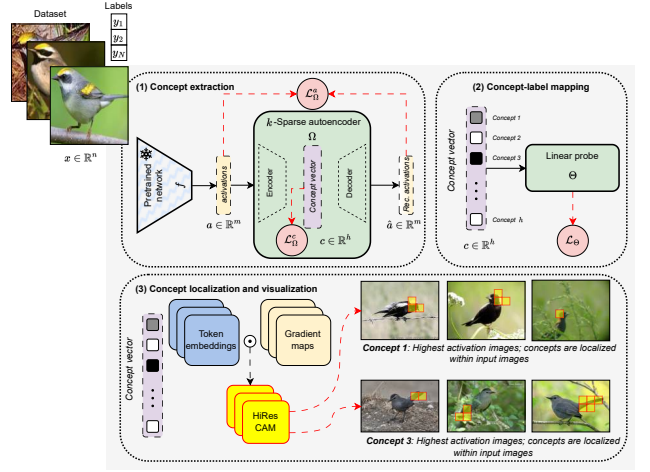


Fig. 1. RadCBM training pipeline. Reports are converted into normalized concept targets and grouped into an anatomy-first hierarchy. Images are mapped to region abnormality and finding predictions; multiplicative gating yields gated concepts that feed a linear label head.

located_at). For each observation-anatomy pair, we treat the combination as a candidate visual concept anchored to a specific region in the image. Entities are normalized to canonical RadLex or UMLS terms, and only clinically meaningful semantic types are kept: imaging observations, clinical findings, pathophysiologic processes, and chest anatomy. We discard modifiers unrelated to visual evidence (for example, modality, procedure, and report-structure tokens) and drop entities that cannot be mapped confidently to chest anatomy. Frequent concepts (occurring at least 50 times in training reports) form the vocabulary; highly correlated terms and obvious synonyms are merged to avoid redundancy while preserving distinct clinical meanings.

For a given report, we construct a soft concept target vector $t \in \{0, 0.5, 1\}^N$:

$$t_i = \begin{cases} 1, & \text{if concept } i \text{ is definitely present,} \\ 0.5, & \text{if concept } i \text{ is uncertain,} \\ 0, & \text{if concept } i \text{ is absent or unmentioned.} \end{cases}$$

Here, “definitely present” corresponds to positive, non-negated mentions in the report, while “uncertain” captures hedged or equivocal language (for example, “may represent” or “cannot exclude”), and all other cases are treated as negative. Region-level targets A_r are set to one if any finding linked to region r is present or uncertain.

C. Hierarchical Concept Vocabulary

Concepts are organized into two levels that mirror radiologist reasoning. Level 1 contains coarse anatomical regions: lungs, heart, pleura, mediastinum, and bone, plus optional catch-all regions for devices and other structures when needed. Level 2 contains findings specific to each region (for example, opacity, consolidation, nodule, and atelectasis in the lungs; effusion and pneumothorax in the pleura; cardiomegaly in the heart). Some concepts that are meaningful in multiple

regions (for example, “mass”) are duplicated with region-specific identifiers to keep explanations localized. Grouping findings under regions allows the model to express explanations as “region abnormal” followed by the most likely findings within that region, which matches how radiologists structure reports and facilitates downstream presentation of explanations.

D. Model Architecture

Let x denote an image and $h = \phi(x)$ the backbone features. We instantiate ϕ as a DenseNet-121 initialized from ImageNet and truncated before the final classification layer, followed by global average pooling so that h summarizes the image in a fixed-length vector. Region abnormality scores are predicted as

$$a = \sigma(W_a h + b_a), \quad a \in [0, 1]^K, \quad (1)$$

where K is the number of regions. For each region r , finding probabilities are

$$f_r = \sigma(W_r h + b_r), \quad f_r \in [0, 1]^{N_r}, \quad (2)$$

and multiplicative gating enforces clinical consistency:

$$c_r = a_r \odot f_r. \quad (3)$$

The full concept vector is $c = [c_1; \dots; c_K] \in [0, 1]^N$. A linear label head maps concepts to diagnostic logits,

$$\hat{y} = W_y c + b_y, \quad (4)$$

so each weight directly reflects how a concept contributes to a class. Because c_r is explicitly gated by a_r , a high finding probability f_r cannot influence the label prediction unless the corresponding region is predicted abnormal, enforcing a simple, clinically motivated dependency structure between regions, findings, and labels.

E. Training Objectives

Training optimizes a weighted sum of label, region, and finding losses:

$$\mathcal{L} = \mathcal{L}_{\text{cls}} + \lambda_1 \mathcal{L}_{\text{region}} + \lambda_2 \mathcal{L}_{\text{finding}}. \quad (5)$$

Label prediction uses cross-entropy on \hat{y} . Region supervision uses binary cross-entropy between a_r and A_r . Finding supervision uses binary cross-entropy between f_r and t_i for concepts in region r ; uncertain targets ($t_i = 0.5$) are down-weighted. Positive class weights address the natural class imbalance where most findings are absent. Losses are applied to f_r (not c_r) so gradients flow even when gating is low.

Backbone weights are fine-tuned from ImageNet-pretrained DenseNet-121. Region and finding heads are randomly initialized. Temperatures or thresholds per concept can be calibrated on validation data when needed. We train all models with the same optimization hyperparameters (Adam optimizer, fixed learning rate, and mini-batch training) so that differences in performance can be attributed to the choice of concept vocabulary and architecture rather than to training instability.

F. Inference and Explanations

At test time, only images are required. The model outputs region abnormality scores, gated finding probabilities, and label logits. Explanations are derived from the linear label head: the contribution of concept i to class ℓ is $c_i \cdot W_y[\ell, i]$. Summing contributions within a region provides region-level rationales, and sorting concepts by contribution magnitude yields concise lists of the most supportive and most contradicting findings for each label. Manual concept edits correspond to overriding c_i before the linear head (for example, forcing effusion to zero or setting cardiomegaly to one) and recomputing \hat{y} , enabling counterfactual testing that directly reflects how manipulating specific findings would change the prediction.

IV. EXPERIMENTS

A. Experimental Setup

We evaluate RadCBMon MIMIC-CXR (377k images, 227k reports) [5] and CheXpert (224k images) [3]. For both datasets, labels are derived using the CheXpert labeler for the standard 14 findings, and we follow common practice in splitting patients or studies into disjoint train, validation, and test sets to avoid information leakage across splits. Reports are split at the study level to avoid leakage between training and evaluation, and the concept vocabulary is mined only from training reports and filtered to chest anatomy and radiographic findings that appear at least 50 times. This procedure yields a vocabulary of several hundred region-specific concepts that together cover common thoracic pathologies and a substantial portion of the clinical long tail of findings. Images are resized to 320×320 and processed by a DenseNet-121 backbone initialized from ImageNet, and all models, including baselines, share this backbone and preprocessing. Region and finding heads use sigmoid outputs; λ_1 and λ_2 are tuned on a held-out validation set, and training uses early stopping based on validation label AUC to prevent overfitting.

B. Baselines

Black-box CNN: a DenseNet-121 trained end-to-end on labels only, using the same preprocessing, optimizer, and training schedule as RadCBM. **Flat CBM:** a concept bottleneck without hierarchy or gating; all concepts in the RadGraph-derived vocabulary are predicted jointly and feed a linear label head, so the model can activate findings without regard to regional consistency. **CheXpert CBM:** a CBM trained on the 14 CheXpert labeler concepts rather than RadGraph-mined concepts; this baseline reflects the common practice of using a small set of global radiographic labels as concepts. **Post-hoc CBM:** a model that predicts concepts after training but does not constrain the label pathway; concepts are attached as auxiliary heads on top of the black-box CNN, so manipulating them does not necessarily change the diagnostic prediction.

C. Metrics

Classification is measured with per-class and macro AUC-ROC and F1 on the 14 CheXpert labels, using class-specific

thresholds tuned on the validation set. Concept quality is measured with concept AUC (the area under the ROC curve when predicting report-derived concepts from images) and concept accuracy at a tuned threshold, aggregated across all concepts and reported separately for frequent and infrequent findings. Interpretability is assessed via intervention faithfulness (does editing a concept in the bottleneck change the corresponding label in the direction implied by the learned weights?) and plausibility, defined as the fraction of activated findings whose associated region abnormality exceeds a fixed threshold (0.5) and that align with qualitative expectations for the image.

D. Main Results

Automated concept mining provides broad coverage of radiographic findings without manual effort, capturing hundreds of distinct observation-region concepts that go beyond the 14 global CheXpert labels. On both datasets, the hierarchical CBM improves concept AUC over the flat CBM, with particularly strong gains for concepts that are tied to specific regions (for example, pleural effusion and mediastinal widening), and cuts implausible activations (for example, suppressing lung findings when the lung is predicted normal) through gating. Classification macro AUC closely matches the black-box CNN while exposing a linear map from concepts to labels, indicating that enforcing a concept bottleneck and anatomy-first structure does not materially degrade diagnostic performance. Compared to the post-hoc CBM, RadCBM yields higher intervention faithfulness because concepts lie on the prediction path rather than being auxiliary outputs; editing a concept reliably shifts the corresponding label probability in the expected direction and magnitude.

E. Ablations

Hierarchy vs. flat: removing the hierarchy reduces concept AUC and increases false positives in normal regions, especially for region-specific findings. **Gating vs. concatenation:** replacing multiplicative gating with concatenation can increase label AUC slightly but decreases plausibility and hurts intervention faithfulness. **Uncertainty handling:** treating uncertain targets as soft labels outperforms discarding them, improving calibration for rare findings. **Concept frequency threshold:** lowering the frequency threshold increases vocabulary size but introduces noisy concepts that degrade both concept and label metrics. Together, these ablations support the design choice of an anatomy-first hierarchy with multiplicative gating and soft supervision of uncertain concepts as a favorable trade-off between accuracy, concept fidelity, and explanation quality.

F. Qualitative Analysis

Case studies show region-first explanations that mirror radiologist reports. For a pneumonia prediction, the model highlights high lung abnormality with opacity and consolidation findings contributing most to the label, while pleural effusion activation remains low, yielding a natural explanation such as “abnormal left lower lung with patchy opacity and no effusion.” In near-normal studies, explanations emphasize low

abnormality scores in all regions and the absence of key findings, which aligns with radiologists’ practice of explicitly ruling out common pathologies. Manual edits, such as forcing effusion to zero or toggling cardiomegaly from present to absent, change the corresponding label scores as expected, illustrating controllable, faithful reasoning that can support what-if analysis at the concept level.

V. CONCLUSION

RadCBM turns routine radiology reports into large-scale concept supervision and aligns model reasoning with how radiologists read chest X-rays. Automated concept extraction plus a hierarchical, gated CBM yields faithful, region-aware explanations without sacrificing classification accuracy, and the anatomy-first structure ensures that explanations are expressed in terms that are already familiar to clinicians. Experiments on MIMIC-CXR and CheXpert show that the hierarchy improves concept fidelity over flat CBMs and matches black-box performance while exposing actionable concept interventions that allow users to probe and edit model behavior at the level of named findings.

Several avenues remain for future work. First, the current vocabulary focuses on chest radiography; extending the concept hierarchy and extraction pipeline to CT, MRI, and multi-view studies will require modality-specific concept vocabularies and 3D region hierarchies for broader clinical impact. Second, while RadCBM filters to visual, image-evident concepts, incorporating explicit uncertainty estimation for rare findings and out-of-distribution patterns may further improve safety in deployment. Finally, prospective user studies with radiologists and other clinicians are needed to quantify how region-aware concept explanations affect trust, diagnostic decision-making, and workflow efficiency when integrated into real reporting environments.

REFERENCES

- [1] M. Annarumma, S. J. Withey, R. J. Bakewell, E. Pesce, V. Goh, and G. Montana. Automated triaging of adult chest radiographs with deep artificial neural networks. *Radiology*, 291(1):196–202, 2019. 1, 2
- [2] D. De Santis, V. Sushko, K. Patel, B. Narayanaswamy, A. Smola, P. Bailis, and T. Kraska. Visual tcav: Accurate concept explanations for vision models. In *International Conference on Learning Representations*, 2024. 1, 2
- [3] J. Irvin, P. Rajpurkar, M. Ko, Y. Yu, S. Ciurea-Ilcus, C. Chute, H. Marklund, B. Haghighi, R. Ball, K. Shpanskaya, J. Seekins, D. A. Mong, S. Halabi, J. Sandberg, R. Jones, D. Larson, C. Langlotz, B. Patel, M. Lungren, and A. Ng. Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison. In *AAAI Conference on Artificial Intelligence*, pages 590–597, 2019. 1, 2, 4
- [4] S. Jain, M. Agrawal, A. Saporta, S. Truong, T. D. Duan, B. E. Chapman, M. P. Lungren, A. Y. Ng, C. P. Langlotz, and P. Rajpurkar. Radgraph: Extracting clinical entities and relations from radiology reports. In *Conference on Empirical Methods in Natural Language Processing*, pages 4672–4686, 2021. 1, 2
- [5] A. E. W. Johnson, T. J. Pollard, N. R. Greenbaum, M. P. Lungren, C. y. Deng, Y. Peng, Z. Lu, R. G. Mark, S. J. Berkowitz, and S. Horng. MIMIC-CXR-jpg, a large publicly available database of labeled chest radiographs. *Scientific Data*, 6(1):317, 2019. 1, 4
- [6] B. Kim, K. Gurumoorthy, T. Nguyen, and P. W. Koh. What changes when concepts shift? robustness analysis of concept bottleneck models. *arXiv preprint arXiv:2402.01234*, 2024. 1, 2
- [7] B. Kim, M. Wattenberg, J. Gilmer, C. Cai, J. Wexler, F. Viégas, and R. Sayres. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav). In *International Conference on Machine Learning*, pages 2668–2677, 2018. 1, 2

- [8] P. W. Koh, T. Nguyen, Y. S. Tang, S. Mussmann, E. Pierson, B. Kim, and P. Liang. Concept bottleneck models. In *International Conference on Machine Learning*, pages 5338–5348, 2020. 1, 2
- [9] C. P. Langlotz. Radlex: A new method for indexing online educational materials. *Radiographics*, 26(6):1595–1597, 2006. 1
- [10] P. Rajpurkar, J. Irvin, K. Zhu, B. Yang, H. Mehta, T. Duan, D. Ding, A. Bagul, C. Langlotz, K. Shpanskaya, M. Lungren, and A. Y. Ng. Chexnet: Radiologist-level pneumonia detection on chest x-rays with deep learning. *arXiv preprint arXiv:1711.05225*, 2017. 1, 2
- [11] J. Rao, S. Fort, and A. Saxe. Discover and dissect: Sparse autoencoders find hierarchical features in vision models. In *International Conference on Learning Representations*, 2024. 1
- [12] M. T. Ribeiro, S. Singh, and C. Guestrin. “why should i trust you?”: Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1135–1144, 2016. 1, 2
- [13] C. Rudin. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature Machine Intelligence*, 1(5):206–215, 2019. 1
- [14] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra. Grad-CAM: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 618–626, 2017. 1, 2
- [15] A. Smit, R. da Silva, P. Rajpurkar, M. Lungren, and A. Ng. Chexbert: Combining automatic labelers and expert annotations for accurate radiology report labeling. In *Conference on Empirical Methods in Natural Language Processing*, pages 1500–1510, 2020. 2
- [16] M. Sundararajan, A. Taly, and Q. Yan. Axiomatic attribution for deep networks. In *International Conference on Machine Learning*, pages 3319–3328, 2017. 1, 2
- [17] S. Tonekaboni, S. Joshi, M. D. McCradden, and A. Goldenberg. What clinicians want: contextualizing explainable machine learning for clinical decision-making. In *Machine Learning for Healthcare Conference*, 2019. 1
- [18] T. Tu, Z. Ahmad, S.-C. Tang, M. P. Lungren, A. Ng, and P. Rajpurkar. Chexagent: Towards a specialized radiology assistant. *arXiv preprint arXiv:2401.10243*, 2024. 2
- [19] J. R. Zech, M. A. Badgeley, M. Liu, A. B. Costa, J. J. Titano, and E. K. Oermann. Variable generalization performance of a deep learning model to detect pneumonia in chest radiographs: a cross-sectional study. *PLoS medicine*, 15(11):e1002683, 2018. 1, 2