# Towards long-tailed, multi-label disease classification from chest X-ray: Overview of the CXR-LT challenge

Gregory Holste[1], Yiliang Zhou[2], Song Wang[1], Ajay Jaiswal[1], Mingquan Lin[2], Sherry Zhuge[3], Yuzhe Yang[4], Dongkyun Kim[5], Trong-Hieu Nguyen-Mau[6], Minh-Triet Tran[6], Jaehyup Jeong[7], Wongi Park[8], Jongbin Ryu[8], Feng Hong[9], Arsh Verma[10], Yosuke Yamagishi[11], Changhyun Kim[12], Hyeryeong Seo[13], Myungjoo Kang[14], Leo Anthony Celi[15,16,17], Zhiyong Lu[18], Ronald M. Summers[19], George Shih[20], Zhangyang Wang[*1], and Yifan Peng[*2]

[1]Department of Electrical and Computer Engineering, The University of Texas at Austin, 78712, Austin, TX USA
[2]Department of Population Health Sciences, Weill Cornell Medicine, 10065, New York, NY USA
[3]School of Information Systems, Carnegie Mellon University, 15213, Pittsburgh, PA USA
[4]Department of Electrical Engineering and Computer Science, Massachussetts Institute of Technology, 02139, Cambridge, MA USA
[5]School of Computer Science, Carnegie Mellon University, 15213, Pittsburgh, PA USA
[6]University of Science, VNU-HCM, 70000, Ho Chi Minh City, Vietnam
[7]KT Research & Development Center, KT Corporation, 06763, Seoul, South Korea
[8]Department of Software and Computer Engineering, Ajou University, 16499, Suwon, South Korea
[9]Cooperative Medianet Innovation Center, Shanghai Jiao Tong University, 200240, Shanghai, China
[10]Wadhwani Institute for Artificial Intelligence, 400079, Mumbai, India
[11]Division of Radiology and Biomedical Engineering, Graduate School of Medicine, The University of Tokyo, 113-0033, Tokyo, Japan
[12]BioMedical AI Team, AIX Future R&D Center, SK Telecom, 04539, Seoul, South Korea
[13]Interdisciplinary Program in AI (IPAI), Seoul National University, 02504, Seoul, South Korea
[14]Department of Mathematical Sciences, Seoul National University, 02504, Seoul, South Korea
[15]Laboratory for Computational Physiology, Massachusetts Institute of Technology, 02139, Cambridge, MA USA
[16]Division of Pulmonary, Critical Care and Sleep Medicine, Beth Israel Deaconess Medical Center, 02215, Boston, MA USA
[17]Department of Biostatistics, Harvard T.H. Chan School of Public Health, 02115, Boston, MA USA
[18]National Center for Biotechnology Information, National Library of Medicine, 20894, Bethesda, MD USA
[19]Clinical Center, National Institutes of Health, 20892, Bethesda, MD USA
[20]Department of Radiology, Weill Cornell Medicine, 10065, New York, NY USA

## Abstract

Many real-world image recognition problems, such as diagnostic medical imaging exams, are "long-tailed" – there are a few common findings followed by many more relatively rare conditions. In chest radiography, diagnosis is both a *long-tailed* and *multi-label* problem, as patients often present with multiple findings simultaneously. While researchers have begun to study the problem of long-tailed learning in medical image recognition, few have studied the interaction of label imbalance and label co-occurrence posed by long-tailed, multi-label disease classification. To engage with the research community on this emerging topic, we conducted an open challenge, **CXR-LT**, on long-tailed, multi-label thorax disease classification from chest X-rays (CXRs). We publicly release a large-scale benchmark dataset of over 350,000 CXRs, each labeled with at least one of 26 clinical findings following a long-tailed distribution. We synthesize common themes of top-performing solutions, providing practical recommendations for long-tailed, multi-label medical image classification. Finally, we use these insights to propose a path forward involving vision-language foundation models for few- and zero-shot disease classification.

**Keywords**: Chest X-ray, Long-tailed learning, Computer-aided diagnosis

---

*Corresponding authors. Email: yip4002@med.cornell.edu, atlaswang@utexas.edu.

# 1  Introduction

Like many diagnostic medical exams, chest X-rays (CXRs) yield a long-tailed distribution of clinical findings. This means that while a small subset of diseases are routinely observed, the majority are quite rare (Zhou et al., 2021). This long-tailed distribution challenges conventional deep learning methods, as they tend to favor common classes and often overlook the infrequent yet crucial classes. In response, several methods (Zhang et al., 2023b) have been proposed lately with a focus on addressing label imbalance in long-tailed medical image recognition tasks (Zhang et al., 2021a; Ju et al., 2021, 2022; Yang et al., 2022). Of note, diagnosing from CXRs is not only a long-tailed problem, but also *multi-label*, since patients often present with multiple disease findings simultaneously. Despite this, only a limited number of studies have incorporated knowledge of label co-occurrence into their learning process (Chen et al., 2020; Wang et al., 2023; Chen et al., 2019a).

Owing to the fact that most large-scale image classification benchmarks feature single-label images with a predominately balanced label distribution, we establish a new benchmark for long-tailed, multi-label medical image classification. Specifically, we expanded the MIMIC-CXR (Johnson et al., 2019a) dataset by increasing the set of target disease findings from 14 to 26. This is achieved by introducing 12 new disease findings by parsing the radiology reports associated with each CXR study.

In our effort to engage with the community on this emerging interdisciplinary topic, we have released the data and launched the **CXR-LT** challenge on long-tailed, multi-label thorax disease classification on CXRs. In this paper, we summarize the CXR-LT challenge, consolidate key insights from top-performing solutions, and offer practical perspetive for advancing long-tailed, multi-label medical image classification. Finally, we use our findings to suggest a path forward toward few- and zero-shot disease classification in the long-tailed, multi-label setting by leveraging multimodal foundation models.

Our contributions can be summarized as follows:

1. We have publicly released a large multi-label, long-tailed CXR dataset containing 377,110 images. Each image is labeled with one or multiple labels from a set of 26 disease findings. In additiion, we have provided a "gold standard" subset encompassing human-annotated consensus labels.
2. We conducted **CXR-LT**, a competition for long-tailed, multi-label thorax disease classification on CXRs. We summarize insights from top-performing teams and offer practical recommendations for advancing long-tailed, multi-label medical image classification.
3. Based on the insights from CXR-LT, we propose a methodological path forward for few- and zero-shot generalization to unseen disease findings via multimodal foundation models.

# 2  Methods

## 2.1  Dataset curation

In this section, we detail the data curation process of two datasets: (i) the CXR-LT dataset used in the challenge, and (ii) a manually annotated "gold standard" test set used for additional evaluation of top-performing solutions after the conclusion of the challenge.

### 2.1.1  CXR-LT dataset

The CXR-LT challenge dataset[1] was created by extending the label set of the MIMIC-CXR dataset[2] (Johnson et al., 2019a), resulting in a more challenging, long-tailed label distribution. Following Holste et al. (2022), the radiology reports associated with each CXR study were parsed via RadText (Wang

---

[1]https://physionet.org/content/cxr-lt-iccv-workshop-cvamd/1.1.0/
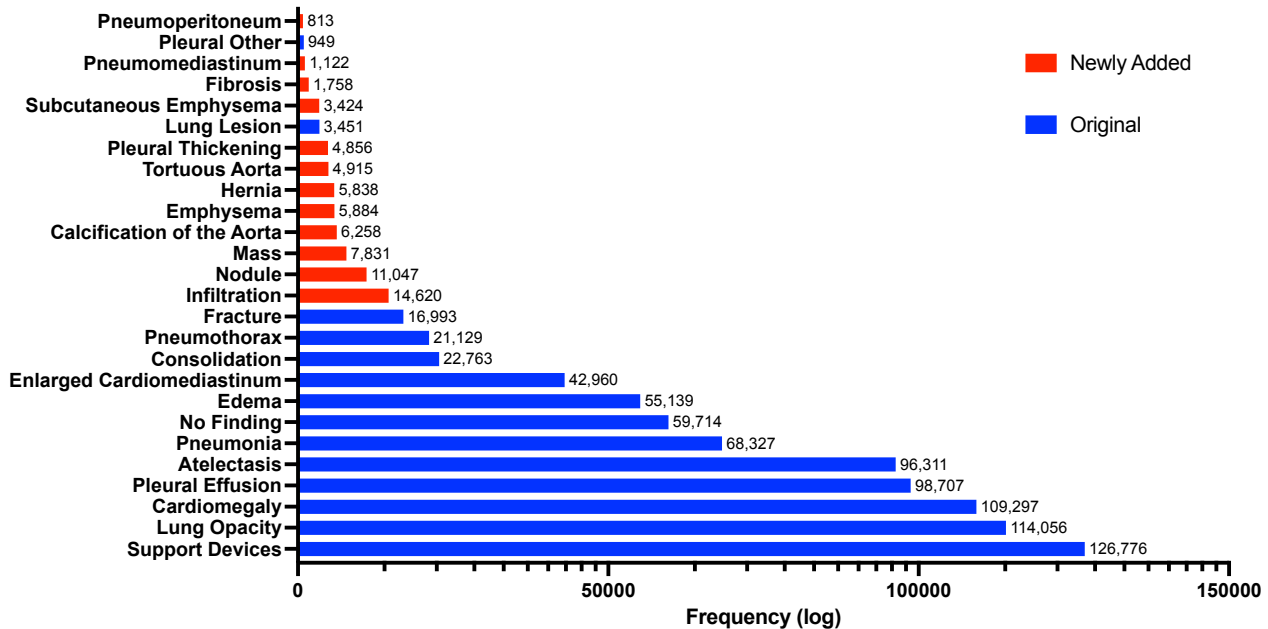[2]https://physionet.org/content/mimic-cxr/2.0.0/

Figure 1: Long-tailed distribution of the CXR-LT 2023 challenge dataset. The dataset was formed by extending the MIMIC-CXR (Johnson et al., 2019a) benchmark to include 12 new clinical findings (red) by parsing radiology reports.

et al., 2022), a radiology text analysis tool, to extract the presence status of *12 new rare disease findings*: (1) Calcification of the Aorta, (2) Emphysema, (3) Fibrosis, (4) Hernia, (5) Infiltration, (6) Mass, (7) Nodule, (8) Pleural Thickening, (9) Pneumomediastinum, (10) Pneumoperitoneum, (11) Subcutaneous Emphysema, and (12) Tortuous Aorta.

The resulting dataset consisted of 377,110 CXRs, each labeled with at least one of 26 disease findings following a long-tailed distribution (Fig. 1). Though MIMIC-CXR contained the images and text reports needed for additional labeling, we used images from the MIMIC-CXR-JPG dataset (Johnson et al., 2019b) in this competition since the preprocessed JPEG images (∼600 GB) would be more accessible to participants than the raw DICOM data (∼4.7 TB) provided in MIMIC-CXR.[3] For the competition, the dataset was randomly split into training (70%), development (10%), and test sets (20%) at the patient level to avoid label leakage. Competition participants would have access to all images, but only have access to labels for the training set.

### 2.1.2 Gold standard test set

While the CXR-LT dataset is large and challenging due to heavy label imbalance and label co-occurrence, it inevitably suffers from label noise much like other datasets with automatically text-mined labels (Abdalla and Fine, 2023). To remedy this, we aimed to construct a "gold standard" set, derived from the challenge test set, with labels that were manually annotated after analyzing the radiology reports. This smaller, but higher quality datasetset would then be used as an auxiliary test set to perform additional evaluations of the top-performing CXR-LT solutions.

To build a gold standard set for evaluation, six annotators manually annotated the presence or absence of the 26 disease findings. A set of 451 radiology reports were randomly sampled from the CXR-LT challenge test set, and each report was annotated by at least two annotators. Prior to annotation, all reports were preprocessed through RadText (Wang et al., 2022) to identify and highlight all relevant disease mentions in the text in order to ease the annotation process. Each annotator was then provided

---

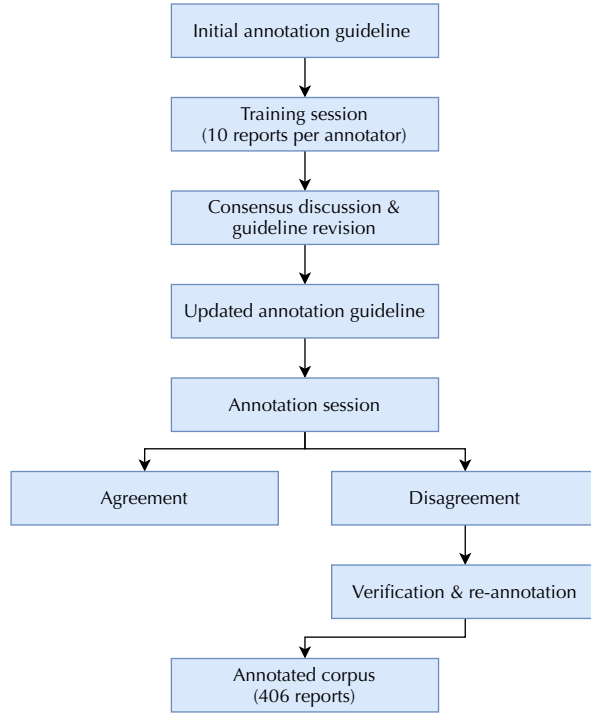[3]https://physionet.org/content/mimic-cxr-jpg/2.0.0/

Figure 2: Flowchart describing CXR-LT gold standard dataset annotation.

with the reports and a list of synonyms for each of the 26 findings. Annotators were asked to select all disease findings that were conclusively affirmed positive in the report. Following MIMIC-CXR, annotators could select "No Finding" if no other findings (except "Support Devices") were present.

Before annotation, a training session was held to align the standards among annotators where each annotator practiced by labeling 10 reports. Any disagreements in this phase were discussed until consensus was reached, leading to the formulation of a shared annotation guideline (Fig. 2). Following the training session, the official annotation process consisted of two rounds: the first round covering 200 reports and the second round covering 251 reports. After each round, individual disease-level disagreements between annotators on a given report were compiled and adjudicated by a third annotator. For the first round, the overall agreement rate was 93.2% and the Cohen's Kappa coefficient was 0.795; for the second round, the agreement rate was 94.9% with a Cohen's kappa of 0.778. After removing reports that were not annotated by at least two readers, the CXR-LT gold standard set consisted of 406 cases. The resulting label distribution of the gold standard set can be found in Supplementary Figure 1.

## 2.2 CXR-LT challenge task

The CXR-LT challenge was formulated as a 26-way multi-label classification problem. Given a CXR, participants were tasked with detecting all disease findings present. If no findings were present, participants could predict "No Finding", with the exception that "No Finding" can co-occur with "Support Devices" as this is not a clinically meaningful *diagnostic* finding. Since this is a multi-label classification problem with severe label imbalance, the primary evaluation metric was mean average precision (mAP), specifically, the "macro-averaged" AP across the 26 classes. While area under the receiver operating characteristic curve (AUROC) is a standard metric employed for related datasets (Wang et al., 2017; Seyyed-Kalantari et al., 2020), AUROC can be heavily inflated in the presence of class imbalance (Fernández et al., 2018; Davis and Goadrich, 2006). Instead, mAP is more appropriate for the long-tailed, multi-label setting since it measures performance across decision thresholds and does not degrade under class imbalance (Rethmeier and Augenstein, 2022). For thoroughness, mean AUROC (mAUROC) and mean F1 score – using a decision threshold of 0.5 for each class – were also calculated.
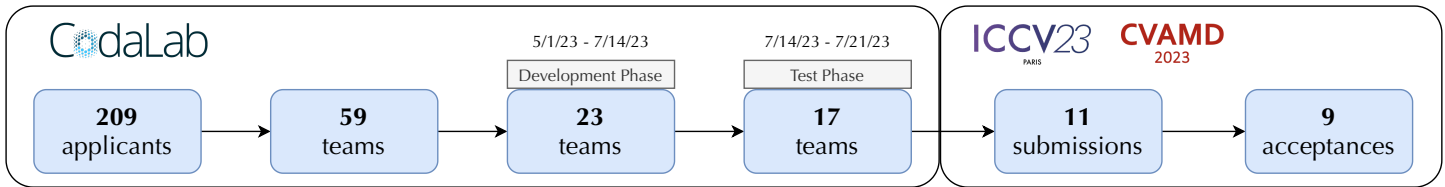
Figure 3: Flowchart describing CXR-LT challenge participation. Over 200 teams applied to participate in the challenge on CodaLab, and 59 teams met registration requirements. Of the 17 teams that participated in the Test Phase, 11 submitted their written solutions for presentation at the ICCV CVAMD 2023 workshop. The top 9 of these submissions were accepted to the workshop and are described in this paper.

The competition was conducted on CodaLab (Pavao et al., 2023). Any registered CodaLab user could apply to participate, but since this competition used MIMIC-CXR-JPG data (Johnson et al., 2019b), which requires credentialing and training through PhysioNet (Goldberger et al., 2000), participants were required to submit proof of PhysioNet credentials to enter. During the Development Phase, registered participants downloaded the labeled training set and (unlabeled) development set, for which they would generate a comma-separated values (CSV) file with predictions to upload. Submissions were then evaluated on the held-out development set and results were updated to a live, public leaderboard. During the Test Phase, test set images (without labels) were released. Participants were asked to submit CSV files with predictions on the much larger, held-out test set and were only given a maximum of 5 successful attempts. For this phase, the leaderboard was kept hidden and the single best-scoring submission (by mAP) by each team was retained. The final Test Phase leaderboard was used to rank participants, primarily by mAP, then by mAUROC in the event of ties.

## 3 Results

### 3.1 CXR-LT challenge participation

The CXR-LT challenge received 209 team applications on CodaLab,[4] of which 59 were approved after providing proof of credentialed access to MIMIC-CXR-JPG (Johnson et al., 2019b). During the Development Phase, 23 teams participated, contributing a total of 525 unique submissions to the public leaderboard. Ultimately, 17 teams participated in the final Test Phase, and 11 of these teams submitted papers describing their challenge solution to the ICCV CVAMD 2023 workshop.[5] The 9 accepted workshop papers, representing the top-performing teams in the CXR-LT challenge, were used for study in this paper (Fig. 3).

### 3.2 Methods of top-performing teams

A summary of top-performing solutions can be found in Table 1, including Test Phase rank, image resolution, backbone architecture used, and other methodological characteristics. Though each solution is described in the paragraphs below, please refer to the paper in each subsection title for full details.

#### 3.2.1 T1 (Kim, 2023)

This team used a two-stage framework that aggregated features across views (e.g., frontal and lateral CXRs). In the first stage, a ConvNeXt-S model (Liu et al., 2022) model was pretrained with Noisy Student (Xie et al., 2020) self-training on the external CXR datasets NIH ChestXRay (Wang et al., 2017),

---

[4] https://codalab.lisn.upsaclay.fr/competitions/12599
[5] https://cvamd2023.github.io/

Table 1: Overview of top-performing CXR-LT challenge solutions. ENS = ensemble; RW = loss reweighting.

| Team | Rank | Image Resolution | Backbone | ENS | RW | Pretraining | Notes |
|---|---|---|---|---|---|---|---|
| T1 | 1 | 1024 | ConvNeXt-S | | ✓ | ImageNet → CheXpert, NIH, VinDr | Two-stage training; cross-view Transformer; ML-Decoder classifier (label as text) |
| T2 | 2 | 512, 768 | EfficientNetV2-S, ConvNeXt-S | ✓ | ✓ | ImageNet | Heavy mosaic augmentation |
| T3 | 3 | 448 | ConvNeXt-B | ✓ | ✓ | ImageNet21K → NIH | Ensemble of "head" and "tail" experts |
| T4 | 4 | 384 | ConvNeXt-B | | ✓ | ImageNet | Custom robust asymmetric loss (RAL) |
| T5 | 5 | 512 | ResNet50* | ✓ | ✓ | ImageNet* | Vision-language modeling (label as text); co-train on NIH, CheXpert |
| T6 | 6 | 448 | ResNeXt101, DenseNet161 | ✓ | | ImageNet → CheXpert, NIH, PadChest | Used synthetic data to augment tail classes |
| T7 | 7 | 224−512 | EfficientNetV2-S | ✓ | | ImageNet, ImageNet21k | Three-stage training with increasing resolution |
| T8 | 8 | 448 | TResNet50 | ✓ | | ImageNet | Heavy CutMix-like augmentation; feature pyramid with deep supervision |
| T9 | 11 | 1024 | ResNet101 | ✓ | | ImageNet | RIDE mixture of experts; LSE pooling; label as text/graph with cross-modal attention |

*T5 additionally used a Transformer text encoder pretrained on PubMedBERT (Gu et al., 2020) and Clinical-T5 (Lehman and Johnson, 2023).

CheXpert (Irvin et al., 2019), and VinDrCXR (Nguyen et al., 2022). Using this backbone as a frozen feature extractor, a Transformer then aggregated multi-view features in a given study. T1 used the ML-Decoder (Ridnik et al., 2023) classification head, which represents the labels as text and performs cross-attention over image (CXR) and text (label) features. Finally, this team utilized a weighted asymmetric loss (Ridnik et al., 2021a) to combat the inter-class imbalance caused by the long-tailed distribution and intra-class imbalance caused by the dominance of negative labels in multi-label classification.

### 3.2.2 T2 (Nguyen-Mau et al., 2023)

This team utilized augmentation, ensemble, and reweighting methods for imbalanced multi-label classification. Specifically, they used an ensemble of EfficientNetV2-S (Tan and Le, 2021) and ConvNeXt-S (Liu et al., 2022) models. Of note, the team made use of heavy "mosaic" augmentation (Bochkovskiy et al., 2020), randomly tiling four CXRs into a single image and using the union of their label sets as ground truth. They used a weighted focal loss (Lin et al., 2017) to handle imbalance, then test-time augmentation and a multi-level ensemble across both model architectures and individual models obtained by stratified cross-validation to improve generalization.

### 3.2.3 T3 (Jeong et al., 2023)

This team proposed an ensemble method based on ConvNeXt-B (Liu et al., 2022) with the CSRA classifier (Zhu and Wu, 2021). After pretraining on the NIH ChestXRay14 (Wang et al., 2017) dataset, T3 trained

three separate models, respectively, on CXR-LT data only from "head" classes, "tail" classes, and all classes; an average of these three models formed the final output. Each model utilized a weighted cross-entropy loss and the Lion optimizer (Chen et al., 2023).

### 3.2.4 T4 (Park et al., 2023)

This team proposed a novel robust asymmetric loss (RAL) for multi-label long-tailed classification. RAL improves upon the popular focal loss (Lin et al., 2017) by including a Hill loss term (Zhang et al., 2021b), which mitigates sensitivity to the negative term of the original focal loss. The team used an ImageNet-pretrained ConvNeXt-B (Liu et al., 2022) with the proposed RAL loss and data augmentation following (Azizi et al., 2021) and (Chen et al., 2019b).

### 3.2.5 T5 (Hong et al., 2023)

This team used a vision-language modeling approach leveraging large pre-trained models. The authors utilized an ImageNet-pretrained ResNet50 (He et al., 2016) and text encoder pre-trained on PubMedBERT (Gu et al., 2020) and Clinical-T5 (Lehman and Johnson, 2023) to extract features from images and label text, respectively. For multi-label classification, they employed a multi-label Transformer query network to aggregate image and text features. To handle imbalance, the team used class-specific loss reweighting informed by validation set performance. They also incorporated external training data (NIH ChestXRay 14 (Wang et al., 2017) and CheXpert (Irvin et al., 2019)), used test-time augmentation, and performed "class-wise" ensembling to improve generalization.

### 3.2.6 T6 (Verma, 2023)

This team used domain-specific pretraining, ensembling, and synthetic data augmentation to improve performance. With an ImageNet-pretrained ResNeXt101 (Xie et al., 2017) and DenseNet101 (Huang et al., 2017), the team further pretrained on the CXR benchmarks NIH ChestXRay14 (Wang et al., 2017), CheXpert (Irvin et al., 2019), and PadChest (Bustos et al., 2020). This team also used RoentGen (Chambon et al., 2022), a multimodal generative model for synthesizing CXRs from natural language, to generate additional CXRs for tail classes in order to combat imbalance.

### 3.2.7 T7 (Yamagishi and Hanaoka, 2023)

This team used a multi-stage training scheme with ensembling and oversampling. For the first stage, an ImageNet21k-pretrained EfficientNetV2-S (Tan and Le, 2021) was trained on $224 \times 224$ resolution images. The weights from this model were then used to train on $320 \times 320$ and $384 \times 384$ images, then $512 \times 512$ images. An ensemble was then formed by averaging the predictions of four models trained on various resolutions, with some models leveraging oversampling of minority classes to mitigate class imbalance. Test-time augmentation and view-based post-processing were also used to boost performance.

### 3.2.8 T8 (Kim et al., 2023)

This team utilized an ImageNet-pretrained TResNet50 (Ridnik et al., 2021b) with heavy augmentation and ensembling. The authors made use of MixUp (Zhang et al., 2017), which linearly combines training images and their labels, and CutMix (Yun et al., 2019), which "cuts and pastes" regions from one training image onto another. They also used a "feature pyramid" approach, extracting pooled features from four layers throughout the network and aggregating these multi-scale features.

Table 2: Final test phase results of the CXR-LT 2023 competition. Presented is average precision (AP) of each team's final model on all 26 classes evaluated on the test set. The best AP for a given class is highlighted in bold.

|  | T1 | T2 | T3 | T4 | T5 | T6 | T7 | T8 | T9 |
|---|---|---|---|---|---|---|---|---|---|
| Atelectasis | **0.622** | 0.609 | 0.611 | 0.607 | 0.606 | 0.610 | 0.602 | 0.595 | 0.546 |
| Calcification of the Aorta | **0.162** | 0.140 | 0.145 | 0.143 | 0.135 | 0.109 | 0.130 | 0.116 | 0.111 |
| Cardiomegaly | 0.661 | 0.652 | 0.652 | 0.648 | 0.653 | 0.652 | **0.668** | 0.640 | 0.581 |
| Consolidation | **0.240** | 0.228 | 0.234 | 0.219 | 0.228 | 0.230 | 0.225 | 0.218 | 0.171 |
| Edema | **0.563** | 0.553 | 0.559 | 0.556 | 0.554 | 0.557 | 0.551 | 0.545 | 0.497 |
| Emphysema | **0.210** | 0.193 | 0.193 | 0.193 | 0.180 | 0.184 | 0.161 | 0.165 | 0.128 |
| Enlarged Cardiomediastinum | 0.186 | 0.184 | 0.184 | 0.184 | **0.186** | 0.185 | 0.183 | 0.177 | 0.140 |
| Fibrosis | **0.167** | 0.163 | 0.157 | 0.153 | 0.154 | 0.154 | 0.116 | 0.132 | 0.120 |
| Fracture | **0.379** | 0.262 | 0.269 | 0.289 | 0.243 | 0.262 | 0.171 | 0.219 | 0.171 |
| Hernia | 0.570 | **0.585** | 0.563 | 0.551 | 0.539 | 0.538 | 0.499 | 0.484 | 0.343 |
| Infiltration | **0.063** | 0.057 | 0.060 | 0.060 | 0.057 | 0.058 | 0.056 | 0.055 | 0.049 |
| Lung Lesion | 0.034 | **0.042** | 0.041 | 0.038 | 0.040 | 0.031 | 0.031 | 0.031 | 0.021 |
| Lung Opacity | **0.617** | 0.597 | 0.603 | 0.597 | 0.594 | 0.598 | 0.590 | 0.584 | 0.529 |
| Mass | **0.250** | 0.224 | 0.213 | 0.206 | 0.200 | 0.227 | 0.187 | 0.167 | 0.112 |
| Nodule | **0.267** | 0.192 | 0.204 | 0.200 | 0.180 | 0.196 | 0.137 | 0.166 | 0.117 |
| Pleural Effusion | **0.843** | 0.829 | 0.831 | 0.832 | 0.830 | 0.805 | 0.822 | 0.822 | 0.781 |
| Pleural Other | **0.070** | 0.037 | 0.043 | 0.040 | 0.039 | 0.016 | 0.059 | 0.042 | 0.007 |
| Pleural Thickening | **0.137** | 0.108 | 0.116 | 0.110 | 0.126 | 0.083 | 0.119 | 0.097 | 0.055 |
| Pneumomediastinum | 0.332 | 0.384 | 0.376 | 0.339 | **0.387** | 0.284 | 0.326 | 0.308 | 0.096 |
| Pneumonia | **0.312** | 0.305 | 0.311 | 0.309 | 0.304 | 0.308 | 0.292 | 0.294 | 0.258 |
| Pneumoperitoneum | **0.324** | 0.316 | 0.261 | 0.283 | 0.303 | 0.237 | 0.262 | 0.235 | 0.155 |
| Pneumothorax | **0.602** | 0.533 | 0.549 | 0.553 | 0.511 | 0.546 | 0.451 | 0.474 | 0.427 |
| Subcutaneous Emphysema | **0.598** | 0.556 | 0.564 | 0.560 | 0.570 | 0.520 | 0.507 | 0.538 | 0.492 |
| Support Devices | **0.918** | 0.906 | 0.916 | 0.913 | 0.910 | 0.910 | 0.894 | 0.903 | 0.887 |
| Tortuous Aorta | **0.066** | 0.061 | 0.060 | 0.060 | 0.058 | 0.056 | 0.063 | 0.053 | 0.045 |
| No Finding | 0.486 | 0.478 | **0.488** | 0.479 | 0.485 | 0.468 | 0.471 | 0.469 | 0.428 |
| Mean | **0.372** | 0.354 | 0.354 | 0.351 | 0.349 | 0.339 | 0.330 | 0.328 | 0.279 |

### 3.2.9 T9 (Seo et al., 2023)

This team built upon ML-GCN (Chen et al., 2019b), a framework for multi-label image classification, which uses GloVe (Pennington et al., 2014) to embed each label as a node within a graph capable of incorporating the co-occurrence patterns of labels. To combat the long-tailed distribution of classes, the authors trained a ResNet101 (He et al., 2016) with class-balanced sampling and the Routing DIverse Experts (RIDE) method (Wang et al., 2020) to diversify the members of their ensemble (Zhang et al., 2023b). They also used log-sum-exp (LSE) pooling (Pinheiro and Collobert, 2015) and a Transformer encoder to attend over image and text features.

## 3.3 Quantitative results

### 3.3.1 CXR-LT Test Phase results

Detailed Test Phase results of 9 top-performing CXR-LT teams can be found in Table 2. T1, the 1st-placed team, reached an mAP of 0.372, considerably outperforming the 2nd-5th-placed teams, who performed

Table 3: Gold standard test set results from CXR-LT 2023 participants. Presented is average precision (AP) of each team's final model on all 26 classes evaluated on our human-annotated gold standard test set. The best AP for a given class is highlighted in bold.

| | T1 | T2 | T3 | T4 | T5 | T6 | T7 | T8 | T9 |
|---|---|---|---|---|---|---|---|---|---|
| Atelectasis | 0.465 | 0.481 | 0.494 | 0.453 | **0.500** | 0.464 | 0.444 | 0.455 | 0.449 |
| Calcification of the Aorta | 0.658 | 0.609 | **0.688** | 0.662 | 0.621 | 0.578 | 0.544 | 0.613 | 0.541 |
| Cardiomegaly | 0.696 | 0.718 | 0.718 | 0.704 | 0.720 | 0.732 | **0.754** | 0.718 | 0.664 |
| Consolidation | 0.415 | 0.449 | 0.471 | 0.474 | 0.436 | **0.476** | 0.411 | 0.426 | 0.406 |
| Edema | 0.600 | 0.590 | **0.601** | 0.572 | 0.589 | 0.587 | 0.571 | 0.587 | 0.540 |
| Emphysema | 0.298 | 0.356 | 0.362 | 0.309 | 0.359 | 0.388 | **0.394** | 0.279 | 0.292 |
| Enlarged Cardiomediastinum | 0.351 | **0.370** | 0.349 | 0.349 | 0.337 | 0.341 | 0.328 | 0.314 | 0.337 |
| Fibrosis | 0.417 | 0.491 | 0.460 | 0.458 | 0.487 | **0.497** | 0.397 | 0.437 | 0.428 |
| Fracture | **0.583** | 0.455 | 0.535 | 0.494 | 0.524 | 0.501 | 0.385 | 0.366 | 0.389 |
| Hernia | 0.759 | **0.808** | 0.804 | 0.766 | 0.722 | 0.723 | 0.714 | 0.708 | 0.514 |
| Infiltration | 0.065 | 0.049 | 0.059 | 0.048 | 0.049 | 0.046 | 0.053 | 0.080 | **0.095** |
| Lung Lesion | 0.028 | **0.071** | 0.033 | 0.042 | 0.030 | 0.032 | 0.066 | 0.040 | 0.044 |
| Lung Opacity | 0.642 | 0.651 | **0.678** | 0.656 | 0.656 | 0.650 | 0.655 | 0.652 | 0.623 |
| Mass | 0.410 | **0.477** | 0.389 | 0.376 | 0.369 | 0.412 | 0.321 | 0.363 | 0.128 |
| Nodule | **0.435** | 0.325 | 0.296 | 0.300 | 0.272 | 0.359 | 0.362 | 0.257 | 0.212 |
| Pleural Effusion | **0.856** | 0.835 | 0.836 | 0.842 | 0.831 | 0.808 | 0.822 | 0.837 | 0.804 |
| Pleural Other | **0.385** | 0.212 | 0.224 | 0.259 | 0.230 | 0.121 | 0.249 | 0.245 | 0.138 |
| Pleural Thickening | **0.386** | 0.249 | 0.249 | 0.204 | 0.244 | 0.190 | 0.251 | 0.210 | 0.132 |
| Pneumomediastinum | 0.771 | 0.830 | 0.795 | 0.800 | 0.823 | 0.757 | **0.837** | 0.776 | 0.569 |
| Pneumonia | 0.114 | **0.131** | 0.127 | 0.123 | 0.111 | 0.111 | 0.114 | 0.103 | 0.096 |
| Pneumoperitoneum | 0.586 | 0.539 | **0.632** | 0.539 | 0.557 | 0.443 | 0.499 | 0.526 | 0.487 |
| Pneumothorax | **0.675** | 0.649 | 0.701 | 0.641 | 0.598 | 0.663 | 0.562 | 0.568 | 0.587 |
| Subcutaneous Emphysema | 0.845 | 0.825 | 0.852 | **0.887** | 0.874 | 0.777 | 0.774 | 0.813 | 0.837 |
| Support Devices | 0.952 | 0.950 | 0.957 | **0.960** | 0.956 | 0.951 | 0.932 | 0.946 | 0.944 |
| Tortuous Aorta | 0.362 | 0.329 | 0.322 | 0.309 | 0.336 | **0.386** | 0.265 | 0.299 | 0.280 |
| No Finding | 0.731 | **0.841** | 0.834 | 0.729 | 0.852 | 0.828 | 0.815 | 0.807 | 0.766 |
| Mean | **0.519** | 0.511 | 0.518 | 0.498 | 0.503 | 0.493 | 0.481 | 0.478 | 0.435 |

similarly with mAP ranging from 0.349 to 0.354; further, T1 achieved best performance on 20 out of 26 classes. Maximum per-class AP ranged widely from 0.063 (Infiltration) to 0.918 (Support Devices), owing primarily to the challenges posed by label imbalance and noise. Additional Test Phase results by AUROC can be found in Supplementary Table 1.

### 3.3.2 Gold standard test set

As described in Section 2.1.2, a "gold standard" test set of 406 newly labeled CXRs was curated for additional evaluation on a small subset of the challenge test set with higher-quality labels. Detailed results of top-performing teams on the gold standard test set can be found in Table 3. Owing to improved label quality, AP values were generally higher in the gold standard set than the original challenge test set, with certain classes experience large changes in performance – for example, the maximum AP jumps from 0.162 to 0.688 for Calcification of the Aorta and from 0.598 to 0.887 for Subcutaneous Emphysema. Despite this, the overall correspondence and ranking of team results remained consistent between the official CXR-LT challenge test set and the gold standard set (Fig. 4; $R^2 = 0.958$, $r = 0.979$, $P = 4.7 \times 10^{-6}$). Additional gold standard set results by AUROC can be found in Supplementary Table 2.
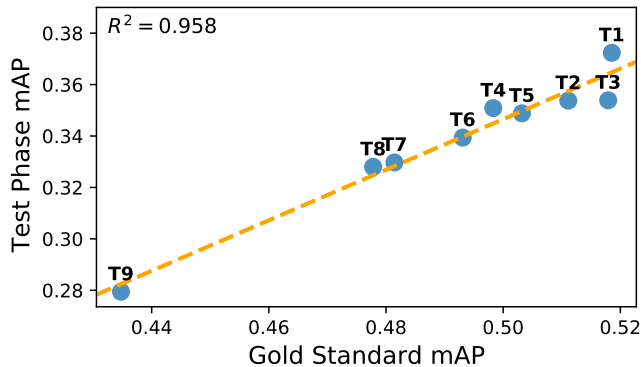
Figure 4: Comparison of performance on CXR-LT Test Phase data (Section 2.1.1) and gold standard test data (Section 2.1.2).

# 4 Discussion

## 4.1 Themes of successful solutions

As outlined in Table 1, several salient patterns emerge among top-performing challenge solutions. Compared to the vast majority of research efforts in natural image recognition, top-performing solutions used relatively high image resolution (as high as $1024 \times 1024$, used by two teams). Also, despite the recent popularity of Vision Transformers (ViTs) (Khan et al., 2022), all top-performing solutions used a convolutional neural network (CNN) as the image encoder. The most popular choice was ConvNeXt (Liu et al., 2022), followed by the EfficientNet (Tan and Le, 2021) and ResNet (He et al., 2016) architecture families. Many (7 out of 9) solutions involved ensemble learning for improved generalization, however it was not strictly necessary for high performance, evidenced by the fact that the 1st- and 4th-placed teams utilized a single well-trained model. Owing to the challenging long-tailed nature of this problem, the top 5 performing solutions all used some form of loss reweighting in order to adequately model rare classes. Additionally, all top-performing solutions utilized pretraining or transfer learning of some kind. While many used standard ImageNet-pretrained models (some leveraging the larger ImageNet21k), several teams performed additional domain-specific pretraining on publicly available, external CXR datasets, which was allowed in this competition; such a multi-stage "generalist", then domain-specific pretraining scheme has proven successful in prior works such as REMEDIS (Azizi et al., 2023). In addition, three teams used a multimodal approach, interpreting and representing the label information as text.

In summary, many successful solutions leveraged the following:
- Relatively high image resolution ($>300 \times 300$)
- Modern CNNs like ConvNeXt and EfficientNetV2
- Strong data augmentation and ensemble learning (though not strictly necessary)
- Loss reweighting to amplify tail classes
- Domain-specific pretraining on CXR data
- Multimodal learning via text-based label representations.

## 4.2 Limitations and future work

In addition to the common themes of successful solutions outlined above, it should be emphasized that the unique and often novel aspects of each team's solution also contributed to their success. For example, Kim (2023) leveraged a cross-view Transformer to aggregate information across radiographic views; Park et al. (2023) proposed a novel robust asymmetric loss (RAL), with additional experiments demonstrating improved performance on other long-tailed medical image classification tasks; Verma (2023) leveraged a vision-language foundation model, RoentGen (Chambon et al., 2022), to synthesize "tail" class exam-

ples to combat the long-tailed problem; Hong et al. (2023) took a vision-language approach leveraging Transformers pretrained on clinical text in order to learn rich representations of the multi-label disease information. One promising observation is that many teams reached very similar overall performance – measured by Test Phase mAP – with dramatically different methods. This suggests that the solutions from our participants may represent *orthogonal* contributions that, when combined, prove greater than the sum of their parts. Future work might unify the insights learned from top-performing solutions into a single long-tailed, multi-label medical image classification framework.

Regarding the data contributions of this work, we acknowledge that the CXR-LT dataset bears the same pitfalls as many other publicly available CXR benchmarks with automatically text-mined labels, namely label noise (Abdalla and Fine, 2023). We attempted to rectify this by manually annotating radiology reports to obtain a gold standard test set for additional evaluation. While this improved the label quality, this approach is limited in that it can only, at best, confirm the opinion of the individual radiologist writing the report. Even for highly trained experts, diagnosis from CXR is difficult and complex, leading to high inter-reader variability (Hopstaken et al., 2004; Sakurada et al., 2012). Ideally, a true "gold standard" dataset would consist of consensus labeling from multiple radiologists' interpretations. However, this is of course prohibitively expensive and time-consuming (Zhou et al., 2021) given the volume of labeled data required to train deep learning models and is the primary motivation for automatic disease labeling in the first place.

Though many diagnostic exams are long-tailed, most publicly available medical imaging datasets only include labels for a few common findings. The CXR-LT dataset thus represents a major contribution to due its large scale (>375,000 images), multi-label nature, and long-tailed label distribution of 26 clinical findings. However, a select few large-scale, long-tailed medical imaging datasets exist, such as HyperKvasir (Borgli et al., 2020) – containing >10,000 endoscopic images labeled with 23 findings – and PadChest (Bustos et al., 2020) – containing >160,000 CXRs labeled with 174 findings.

While CXR-LT and PadChest represent meaningful contributions to long-tailed learning from CXR, it should be noted that the "true" long tail of all clinical findings is at least an order of magnitude longer than any current publicly available dataset can offer. For example, Radiology Gamuts Ontology[6] (Budovec et al., 2014) documents 4,691 unique radiological image findings. Thus, one way to enhance the CXR-LT dataset might be to include an even wider variety of automatically text-mined findings to mimic the extremely long tail of real-world CXR. However, this approach too has its own limitations. Even if we could construct a dataset with labels for up to 1,000 clinical findings, ranging from common and well-studied to exceedingly rare, and train a model on this long-tailed data, what happens when a new finding is encountered? One might argue that the only way to tackle the *true* long-tailed distribution of imaging findings is to develop a model that can adaptively generalize to previously unseen diseases. Future iterations of the CXR-LT challenge may consider this problem through the lens of zero-shot classification: can participants train a model to accurately detect a clinical finding that the model has *not been trained on*?

## 4.3    The future of long-tailed, multi-label learning

If zero-shot disease classification is the ultimate step toward clinically viable long-tailed medical image classification, then vision-language foundation models provide a very promising path forward. Several top-performing CXR-LT teams found success by encoding the label information as text, allowing for rich representation learning of the disease labels and their correlations. This allowed Kim (2023) to better handle the long-tailed, multi-label distribution via the text- and attention-based ML-Decoder classifier (Ridnik et al., 2023) and Seo et al. (2023) to exploit correlations between labels via graph learning of text-based label representations via ML-GCN (Chen et al., 2019b). While, for example, the approach of Hong et al. (2023) did not earn them 1st place in this competition, it would almost certainly prove the most useful when encountering a previously unseen disease finding, a practical scenario in real-world

---

[6] http://www.gamuts.net/about.php

clinical deployment. Since Hong et al. (2023) leverage Transformer encoders pretrained on large collections of clinical text including PubMedBERT (Gu et al., 2020) and Clinical-T5 (Lehman and Johnson, 2023), the model retains a rough *semantic* understanding of biomedical concepts via the textual representation of disease information. This in turn allows natural generalization to new findings by relating the text "query" of the potential new disease finding to relevant concepts encountered during pretraining, which is expressly not possible with standard unimodal deep learning methods.

Recent studies have shown that vision-language modeling with encoders pretrained on large collections of medical image and text data enable zero-shot disease classification, in some cases nearly reaching the performance of fully supervised approaches (Hayat et al., 2021; Tiu et al., 2022; Mishra et al., 2023; Zhang et al., 2023a). Further, such a multimodal approach can be readily combined with many of the methods employed by top-performing CXR-LT teams such as loss re-weighting, heavy data augmentation, and ensemble learning. This would allow for a computer-aided diagnosis system capable of adaptively generalizing to unseen findings, thus encapsulating the true long tail of all imaging findings.

# 5    Conclusion

In summary, we curated and publicly released a large-scale dataset of over 375,000 CXR images for long-tailed, multi-label disease classification. We then hosted an open challenge, CXR-LT, to engage with the research community on this important task. We compiled and synthesized common threads through the most successful challenge solutions, providing practical recommendations for long-tailed, multi-label medical image classification. Lastly, we identify a path forward toward tackling the "true" long tail of all imaging findings via multimodal vision-language foundation models capable of zero-shot generalization to unseen diseases, which future iterations of the CXR-LT challenge will address.

# Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

# Acknowledgments

# References

Abdalla, M., Fine, B., 2023. Hurdles to artificial intelligence deployment: Noise in schemas and "gold" labels. Radiology: Artificial Intelligence 5, e220056.

Azizi, S., Culp, L., Freyberg, J., Mustafa, B., Baur, S., Kornblith, S., Chen, T., Tomasev, N., Mitrović, J., Strachan, P., et al., 2023. Robust and data-efficient generalization of self-supervised machine learning for diagnostic imaging. Nature Biomedical Engineering , 1–24.
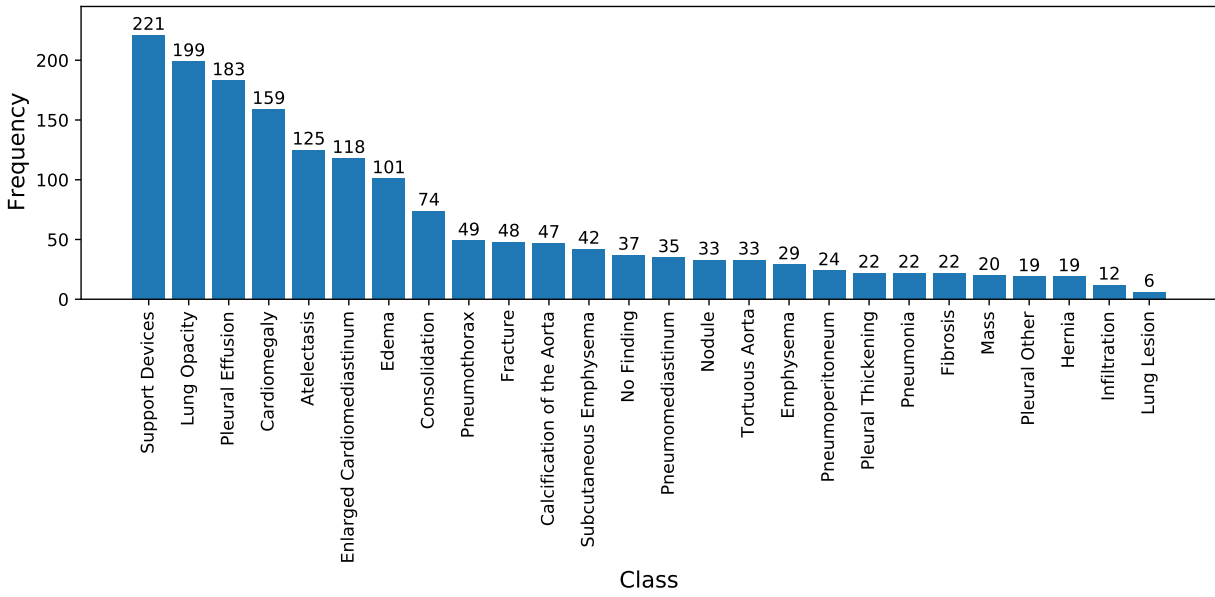
Azizi, S., Mustafa, B., Ryan, F., Beaver, Z., Freyberg, J., Deaton, J., Loh, A., Karthikesalingam, A., Kornblith, S., Chen, T., et al., 2021. Big self-supervised models advance medical image classification, in: Proceedings of the IEEE/CVF international conference on computer vision, pp. 3478–3488.

Bochkovskiy, A., Wang, C.Y., Liao, H.Y.M., 2020. Yolov4: Optimal speed and accuracy of object detection. arXiv preprint arXiv:2004.10934 .

Borgli, H., Thambawita, V., Smedsrud, P.H., Hicks, S., Jha, D., Eskeland, S.L., Randel, K.R., Pogorelov, K., Lux, M., Nguyen, D.T.D., et al., 2020. Hyperkvasir, a comprehensive multi-class image and video dataset for gastrointestinal endoscopy. Scientific data 7, 283.

Budovec, J.J., Lam, C.A., Kahn Jr, C.E., 2014. Informatics in radiology: radiology gamuts ontology: differential diagnosis for the semantic web. Radiographics 34, 254–264.

Bustos, A., Pertusa, A., Salinas, J.M., De La Iglesia-Vaya, M., 2020. Padchest: A large chest x-ray image dataset with multi-label annotated reports. Medical image analysis 66, 101797.

Chambon, P., Bluethgen, C., Delbrouck, J.B., Van der Sluijs, R., Połacin, M., Chaves, J.M.Z., Abraham, T.M., Purohit, S., Langlotz, C.P., Chaudhari, A., 2022. Roentgen: vision-language foundation model for chest x-ray generation. arXiv preprint arXiv:2211.12737 .

Chen, B., Li, J., Lu, G., Yu, H., Zhang, D., 2020. Label co-occurrence learning with graph convolutional networks for multi-label chest x-ray image classification. IEEE journal of biomedical and health informatics 24, 2292–2302.

Chen, H., Miao, S., Xu, D., Hager, G.D., Harrison, A.P., 2019a. Deep hierarchical multi-label classification of chest x-ray images, in: International conference on medical imaging with deep learning, PMLR. pp. 109–120.

Chen, X., Liang, C., Huang, D., Real, E., Wang, K., Liu, Y., Pham, H., Dong, X., Luong, T., Hsieh, C.J., et al., 2023. Symbolic discovery of optimization algorithms. arXiv preprint arXiv:2302.06675 .

Chen, Z.M., Wei, X.S., Wang, P., Guo, Y., 2019b. Multi-label image recognition with graph convolutional networks, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 5177–5186.

Davis, J., Goadrich, M., 2006. The relationship between precision-recall and roc curves, in: Proceedings of the 23rd international conference on Machine learning, pp. 233–240.

Fernández, A., García, S., Galar, M., Prati, R.C., Krawczyk, B., Herrera, F., 2018. Learning from imbalanced data sets. volume 10. Springer.

Goldberger, A.L., Amaral, L.A.N., Glass, L., Hausdorff, J.M., Ivanov, P.C., Mark, R.G., Mietus, J.E., Moody, G.B., Peng, C.K., Stanley, H.E., 2000. PhysioBank, PhysioToolkit, and PhysioNet: Components of a new research resource for complex physiologic signals. Circulation 101, e215–e220.

Gu, Y., Tinn, R., Cheng, H., Lucas, M., Usuyama, N., Liu, X., Naumann, T., Gao, J., Poon, H., 2020. Domain-specific language model pretraining for biomedical natural language processing. arXiv:arXiv:2007.15779.

Hayat, N., Lashen, H., Shamout, F.E., 2021. Multi-label generalized zero shot learning for the classification of disease in chest radiographs, in: Machine learning for healthcare conference, PMLR. pp. 461–477.

He, K., Zhang, X., Ren, S., Sun, J., 2016. Deep residual learning for image recognition, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 770–778.

Holste, G., Wang, S., Jiang, Z., Shen, T.C., Shih, G., Summers, R.M., Peng, Y., Wang, Z., 2022. Long-tailed classification of thorax diseases on chest x-ray: A new benchmark study, in: MICCAI Workshop on Data Augmentation, Labelling, and Imperfections, Springer. pp. 22–32.

Hong, F., Dai, T., Yao, J., Zhang, Y., Wang, Y., 2023. Bag of tricks for long-tailed multi-label classification on chest x-rays. arXiv preprint arXiv:2308.08853 .

Hopstaken, R., Witbraad, T., Van Engelshoven, J., Dinant, G., 2004. Inter-observer variation in the interpretation of chest radiographs for pneumonia in community-acquired lower respiratory tract infections. Clinical radiology 59, 743–752.

Huang, G., Liu, Z., Van Der Maaten, L., Weinberger, K.Q., 2017. Densely connected convolutional networks, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 4700–4708.

Irvin, J., Rajpurkar, P., Ko, M., Yu, Y., Ciurea-Ilcus, S., Chute, C., Marklund, H., Haghgoo, B., Ball, R., Shpanskaya, K., et al., 2019. Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison, in: Proceedings of the AAAI conference on artificial intelligence, pp. 590–597.

Jeong, J., Jeoun, B., Park, Y., Han, B., 2023. An optimized ensemble framework for multi-label classification on long-tailed chest x-ray data, in: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) Workshops, pp. 2739–2746.

Johnson, A.E., Pollard, T.J., Berkowitz, S.J., Greenbaum, N.R., Lungren, M.P., Deng, C.y., Mark, R.G., Horng, S., 2019a. Mimic-cxr, a de-identified publicly available database of chest radiographs with free-text reports. Scientific data 6, 317.

Johnson, A.E., Pollard, T.J., Greenbaum, N.R., Lungren, M.P., Deng, C.y., Peng, Y., Lu, Z., Mark, R.G., Berkowitz, S.J., Horng, S., 2019b. Mimic-cxr-jpg, a large publicly available database of labeled chest radiographs. arXiv preprint arXiv:1901.07042 .

Ju, L., Wang, X., Wang, L., Liu, T., Zhao, X., Drummond, T., Mahapatra, D., Ge, Z., 2021. Relational subsets knowledge distillation for long-tailed retinal diseases recognition, in: Medical Image Computing and Computer Assisted Intervention–MICCAI 2021: 24th International Conference, Strasbourg, France, September 27–October 1, 2021, Proceedings, Part VIII 24, Springer. pp. 3–12.

Ju, L., Wu, Y., Wang, L., Yu, Z., Zhao, X., Wang, X., Bonnington, P., Ge, Z., 2022. Flexible sampling for long-tailed skin lesion classification, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer. pp. 462–471.

Khan, S., Naseer, M., Hayat, M., Zamir, S.W., Khan, F.S., Shah, M., 2022. Transformers in vision: A survey. ACM computing surveys (CSUR) 54, 1–41.

Kim, C., Kim, G., Yang, S., Kim, H., Lee, S., Cho, H., 2023. Chest x-ray feature pyramid sum model with diseased area data augmentation method, in: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) Workshops, pp. 2757–2766.

Kim, D., 2023. Chexfusion: Effective fusion of multi-view features using transformers for long-tailed chest x-ray classification, in: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) Workshops, pp. 2702–2710.

Lehman, E., Johnson, A., 2023. Clinical-t5: Large language models built using mimic clinical text.

Lin, T.Y., Goyal, P., Girshick, R., He, K., Dollár, P., 2017. Focal loss for dense object detection, in: Proceedings of the IEEE international conference on computer vision, pp. 2980–2988.

Liu, Z., Mao, H., Wu, C.Y., Feichtenhofer, C., Darrell, T., Xie, S., 2022. A convnet for the 2020s, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 11976–11986.

Mishra, A., Mittal, R., Jestin, C., Tingos, K., Rajpurkar, P., 2023. Improving zero-shot detection of low prevalence chest pathologies using domain pre-trained language models. arXiv preprint arXiv:2306.08000 .

Nguyen, H.Q., Lam, K., Le, L.T., Pham, H.H., Tran, D.Q., Nguyen, D.B., Le, D.D., Pham, C.M., Tong, H.T., Dinh, D.H., et al., 2022. Vindr-cxr: An open dataset of chest x-rays with radiologist's annotations. Scientific Data 9, 429.

Nguyen-Mau, T.H., Huynh, T.L., Le, T.D., Nguyen, H.D., Tran, M.T., 2023. Advanced augmentation and ensemble approaches for classifying long-tailed multi-label chest x-rays, in: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) Workshops, pp. 2729–2738.

Park, W., Park, I., Kim, S., Ryu, J., 2023. Robust asymmetric loss for multi-label long-tailed learning, in: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) Workshops, pp. 2711–2720.

Pavao, A., Guyon, I., Letournel, A.C., Tran, D.T., Baro, X., Escalante, H.J., Escalera, S., Thomas, T., Xu, Z., 2023. Codalab competitions: An open source platform to organize scientific challenges. Journal of Machine Learning Research 24, 1–6. URL: http://jmlr.org/papers/v24/21-1436.html.

Pennington, J., Socher, R., Manning, C.D., 2014. Glove: Global vectors for word representation, in: Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP), pp. 1532–1543.

Pinheiro, P.O., Collobert, R., 2015. From image-level to pixel-level labeling with convolutional networks, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 1713–1721.

Rethmeier, N., Augenstein, I., 2022. Long-tail zero and few-shot learning via contrastive pretraining on and for small data, in: Computer Sciences & Mathematics Forum, MDPI. p. 10.

Ridnik, T., Ben-Baruch, E., Zamir, N., Noy, A., Friedman, I., Protter, M., Zelnik-Manor, L., 2021a. Asymmetric loss for multi-label classification, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 82–91.

Ridnik, T., Lawen, H., Noy, A., Ben Baruch, E., Sharir, G., Friedman, I., 2021b. Tresnet: High performance gpu-dedicated architecture, in: proceedings of the IEEE/CVF winter conference on applications of computer vision, pp. 1400–1409.

Ridnik, T., Sharir, G., Ben-Cohen, A., Ben-Baruch, E., Noy, A., 2023. Ml-decoder: Scalable and versatile classification head, in: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision, pp. 32–41.

Sakurada, S., Hang, N.T., Ishizuka, N., Toyota, E., Hung, L.D., Chuc, P.T., Lien, L.T., Thuong, P.H., Bich, P.T., Keicho, N., et al., 2012. Inter-rater agreement in the assessment of abnormal chest x-ray findings for tuberculosis between two asian countries. BMC infectious diseases 12, 1–8.

Seo, H., Lee, M., Cheong, W., Yoon, H., Kim, S., Kang, M., 2023. Enhancing multi-label long-tailed classification on chest x-rays through ml-gcn augmentation, in: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) Workshops, pp. 2747–2756.

Seyyed-Kalantari, L., Liu, G., McDermott, M., Chen, I.Y., Ghassemi, M., 2020. Chexclusion: Fairness gaps in deep chest x-ray classifiers, in: BIOCOMPUTING 2021: proceedings of the Pacific symposium, World Scientific. pp. 232–243.

Tan, M., Le, Q., 2021. Efficientnetv2: Smaller models and faster training, in: International conference on machine learning, PMLR. pp. 10096–10106.

Tiu, E., Talius, E., Patel, P., Langlotz, C.P., Ng, A.Y., Rajpurkar, P., 2022. Expert-level detection of pathologies from unannotated chest x-ray images via self-supervised learning. Nature Biomedical Engineering 6, 1399–1406.

Verma, A., 2023. How can we tame the long-tail of chest x-ray datasets? arXiv preprint arXiv:2309.04293 .

Wang, G., Wang, P., Cong, J., Liu, K., Wei, B., 2023. Bb-gcn: A bi-modal bridged graph convolutional network for multi-label chest x-ray recognition. arXiv preprint arXiv:2302.11082 .

Wang, S., Lin, M., Ding, Y., Shih, G., Lu, Z., Peng, Y., 2022. Radiology text analysis system (radtext): Architecture and evaluation, in: 2022 IEEE 10th International Conference on Healthcare Informatics (ICHI), pp. 288–296. doi:10.1109/ICHI54592.2022.00050.

Wang, X., Lian, L., Miao, Z., Liu, Z., Yu, S.X., 2020. Long-tailed recognition by routing diverse distribution-aware experts. arXiv preprint arXiv:2010.01809 .

Wang, X., Peng, Y., Lu, L., Lu, Z., Bagheri, M., Summers, R.M., 2017. Chestx-ray8: Hospital-scale chest x-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 2097–2106.

Xie, Q., Luong, M.T., Hovy, E., Le, Q.V., 2020. Self-training with noisy student improves imagenet classification, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 10687–10698.

Xie, S., Girshick, R., Dollár, P., Tu, Z., He, K., 2017. Aggregated residual transformations for deep neural networks, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 1492–1500.

Yamagishi, Y., Hanaoka, S., 2023. Effect of stage training for long-tailed multi-label image classification, in: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) Workshops, pp. 2721–2728.

Yang, Z., Pan, J., Yang, Y., Shi, X., Zhou, H.Y., Zhang, Z., Bian, C., 2022. Proco: Prototype-aware contrastive learning for long-tailed medical image classification, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer. pp. 173–182.

Yun, S., Han, D., Oh, S.J., Chun, S., Choe, J., Yoo, Y., 2019. Cutmix: Regularization strategy to train strong classifiers with localizable features, in: Proceedings of the IEEE/CVF international conference on computer vision, pp. 6023–6032.

Zhang, H., Cisse, M., Dauphin, Y.N., Lopez-Paz, D., 2017. mixup: Beyond empirical risk minimization. arXiv preprint arXiv:1710.09412 .

Zhang, R., Haihong, E., Yuan, L., He, J., Zhang, H., Zhang, S., Wang, Y., Song, M., Wang, L., 2021a. Mbnm: multi-branch network based on memory features for long-tailed medical image recognition. Computer Methods and Programs in Biomedicine 212, 106448.

Zhang, X., Wu, C., Zhang, Y., Xie, W., Wang, Y., 2023a. Knowledge-enhanced visual-language pre-training on chest radiology images. Nature Communications 14, 4542.

Zhang, Y., Cheng, Y., Huang, X., Wen, F., Feng, R., Li, Y., Guo, Y., 2021b. Simple and robust loss design for multi-label learning with missing labels. arXiv preprint arXiv:2112.07368 .

Zhang, Y., Kang, B., Hooi, B., Yan, S., Feng, J., 2023b. Deep long-tailed learning: A survey. IEEE Transactions on Pattern Analysis and Machine Intelligence .

Zhou, S.K., Greenspan, H., Davatzikos, C., Duncan, J.S., Van Ginneken, B., Madabhushi, A., Prince, J.L., Rueckert, D., Summers, R.M., 2021. A review of deep learning in medical imaging: Imaging traits, technology trends, case studies with progress highlights, and future promises. Proceedings of the IEEE 109, 820–838.

Zhu, K., Wu, J., 2021. Residual attention: A simple but effective method for multi-label recognition, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 184–193.

Supplementary Figure 1: Long-tailed distribution of the CXR-LT gold standard test set.

Supplementary Table 1: Final test phase results of the CXR-LT 2023 competition. Presented is area under the receiver operating characteristic curve (AUROC) of each team's final model on all 26 classes evaluated on the test set. The best AUROC for a given class is highlighted in bold.

| | T1 | T2 | T3 | T4 | T5 | T6 | T7 | T8 | T9 |
|---|---|---|---|---|---|---|---|---|---|
| Atelectasis | **0.838** | 0.828 | 0.828 | 0.827 | 0.826 | 0.828 | 0.825 | 0.819 | 0.790 |
| Calcification of the Aorta | **0.920** | 0.896 | 0.891 | 0.893 | 0.892 | 0.871 | 0.875 | 0.871 | 0.857 |
| Cardiomegaly | 0.815 | 0.810 | 0.809 | 0.807 | 0.810 | 0.809 | **0.818** | 0.803 | 0.770 |
| Consolidation | **0.794** | 0.787 | 0.789 | 0.783 | 0.784 | 0.789 | 0.781 | 0.779 | 0.733 |
| Edema | **0.858** | 0.855 | 0.857 | 0.857 | 0.854 | 0.855 | 0.853 | 0.851 | 0.830 |
| Emphysema | **0.916** | 0.903 | 0.911 | 0.909 | 0.899 | 0.907 | 0.893 | 0.895 | 0.871 |
| Enlarged Cardiomediastinum | **0.619** | 0.617 | 0.615 | 0.617 | 0.615 | 0.615 | 0.61 | 0.605 | 0.557 |
| Fibrosis | **0.922** | 0.916 | 0.921 | 0.919 | 0.907 | 0.912 | 0.900 | 0.901 | 0.838 |
| Fracture | **0.852** | 0.795 | 0.799 | 0.814 | 0.788 | 0.799 | 0.756 | 0.774 | 0.722 |
| Hernia | 0.914 | **0.916** | 0.910 | 0.909 | 0.909 | 0.905 | 0.899 | 0.891 | 0.853 |
| Infiltration | 0.630 | 0.615 | 0.627 | **0.634** | 0.615 | 0.621 | 0.609 | 0.611 | 0.58 |
| Lung Lesion | **0.802** | 0.791 | 0.790 | 0.791 | 0.795 | 0.769 | 0.760 | 0.769 | 0.717 |
| Lung Opacity | **0.800** | 0.784 | 0.786 | 0.783 | 0.781 | 0.786 | 0.780 | 0.774 | 0.737 |
| Mass | **0.821** | 0.811 | 0.811 | 0.800 | 0.807 | 0.814 | 0.803 | 0.789 | 0.727 |
| Nodule | **0.846** | 0.806 | 0.809 | 0.804 | 0.796 | 0.803 | 0.769 | 0.781 | 0.737 |
| Pleural Effusion | **0.930** | 0.921 | 0.922 | 0.923 | 0.922 | 0.909 | 0.917 | 0.917 | 0.901 |
| Pleural Other | **0.910** | 0.878 | 0.876 | 0.882 | 0.890 | 0.857 | 0.886 | 0.872 | 0.751 |
| Pleural Thickening | **0.888** | 0.835 | 0.843 | 0.842 | 0.847 | 0.793 | 0.848 | 0.819 | 0.763 |
| Pneumomediastinum | 0.918 | 0.934 | 0.933 | 0.933 | **0.939** | 0.930 | 0.917 | 0.911 | 0.846 |
| Pneumonia | 0.657 | 0.650 | 0.657 | **0.658** | 0.651 | 0.654 | 0.641 | 0.645 | 0.609 |
| Pneumoperitoneum | 0.908 | 0.901 | 0.900 | 0.900 | **0.913** | 0.881 | 0.893 | 0.862 | 0.815 |
| Pneumothorax | **0.886** | 0.876 | 0.878 | 0.880 | 0.873 | 0.875 | 0.858 | 0.858 | 0.833 |
| Subcutaneous Emphysema | **0.990** | 0.978 | 0.987 | 0.986 | 0.982 | 0.982 | 0.98 | 0.975 | 0.965 |
| Support Devices | **0.956** | 0.946 | 0.952 | 0.952 | 0.947 | 0.948 | 0.938 | 0.942 | 0.934 |
| Tortuous Aorta | **0.841** | 0.834 | 0.821 | 0.825 | 0.829 | 0.817 | 0.824 | 0.814 | 0.777 |
| No Finding | **0.861** | 0.855 | 0.858 | 0.855 | 0.855 | 0.845 | 0.851 | 0.849 | 0.821 |
| Mean | **0.850** | 0.836 | 0.838 | 0.838 | 0.836 | 0.830 | 0.826 | 0.822 | 0.782 |

Supplementary Table 2: Gold standard test set results from CXR-LT 2023 participants. Presented is area under the receiver operating characteristic (AUROC) of each team's final model on all 26 classes evaluated on our human-annotated gold standard test set. The best AUROC for a given class is highlighted in bold.

| | T1 | T2 | T3 | T4 | T5 | T6 | T7 | T8 | T9 |
|---|---|---|---|---|---|---|---|---|---|
| Atelectasis | 0.702 | 0.698 | 0.706 | 0.694 | **0.713** | 0.699 | 0.693 | 0.692 | 0.687 |
| Calcification of the Aorta | **0.928** | 0.918 | 0.926 | 0.922 | 0.911 | 0.890 | 0.855 | 0.895 | 0.876 |
| Cardiomegaly | 0.798 | 0.811 | 0.812 | 0.807 | 0.817 | **0.824** | 0.821 | 0.808 | 0.764 |
| Consolidation | 0.746 | 0.760 | 0.766 | 0.762 | 0.752 | **0.766** | 0.748 | 0.728 | 0.719 |
| Edema | **0.839** | 0.837 | 0.833 | 0.822 | 0.833 | 0.831 | 0.820 | 0.828 | 0.809 |
| Emphysema | 0.863 | **0.881** | 0.876 | 0.864 | 0.876 | 0.872 | 0.877 | 0.853 | 0.854 |
| Enlarged Cardiomediastinum | 0.574 | **0.588** | 0.583 | 0.561 | 0.580 | 0.577 | 0.564 | 0.542 | 0.575 |
| Fibrosis | 0.884 | 0.858 | 0.880 | 0.868 | **0.886** | 0.867 | 0.841 | 0.852 | 0.872 |
| Fracture | **0.887** | 0.823 | 0.856 | 0.837 | 0.828 | 0.841 | 0.808 | 0.794 | 0.809 |
| Hernia | 0.938 | 0.928 | 0.929 | 0.944 | 0.925 | 0.940 | **0.945** | 0.889 | 0.909 |
| Infiltration | **0.674** | 0.571 | 0.578 | 0.598 | 0.565 | 0.571 | 0.575 | 0.584 | 0.633 |
| Lung Lesion | 0.651 | 0.704 | 0.674 | **0.712** | 0.650 | 0.699 | 0.637 | 0.654 | 0.699 |
| Lung Opacity | 0.688 | 0.690 | **0.700** | 0.676 | 0.699 | 0.686 | 0.693 | 0.698 | 0.667 |
| Mass | 0.828 | 0.861 | 0.851 | 0.842 | 0.841 | **0.884** | 0.829 | 0.783 | 0.733 |
| Nodule | **0.801** | 0.706 | 0.713 | 0.737 | 0.728 | 0.791 | 0.734 | 0.724 | 0.730 |
| Pleural Effusion | **0.880** | 0.860 | 0.865 | 0.870 | 0.864 | 0.847 | 0.853 | 0.861 | 0.851 |
| Pleural Other | **0.893** | 0.833 | 0.858 | 0.852 | 0.890 | 0.760 | 0.836 | 0.869 | 0.767 |
| Pleural Thickening | **0.826** | 0.767 | 0.759 | 0.751 | 0.777 | 0.785 | 0.784 | 0.773 | 0.72 |
| Pneumomediastinum | 0.953 | **0.982** | 0.974 | 0.973 | 0.98 | 0.972 | 0.982 | 0.967 | 0.849 |
| Pneumonia | 0.631 | **0.677** | 0.663 | 0.650 | 0.659 | 0.653 | 0.641 | 0.656 | 0.615 |
| Pneumoperitoneum | 0.882 | 0.881 | 0.903 | 0.905 | **0.907** | 0.888 | 0.884 | 0.890 | 0.824 |
| Pneumothorax | 0.942 | 0.934 | **0.944** | 0.943 | 0.926 | 0.940 | 0.905 | 0.924 | 0.906 |
| Subcutaneous Emphysema | 0.986 | 0.973 | **0.988** | 0.988 | 0.987 | 0.978 | 0.978 | 0.973 | 0.985 |
| Support Devices | 0.948 | 0.941 | 0.946 | **0.950** | 0.948 | 0.941 | 0.918 | 0.939 | 0.934 |
| Tortuous Aorta | **0.861** | 0.831 | 0.822 | 0.828 | 0.833 | 0.829 | 0.811 | 0.824 | 0.806 |
| No Finding | 0.944 | **0.957** | 0.953 | 0.943 | 0.952 | 0.95 | 0.952 | 0.944 | 0.938 |
| Mean | **0.829** | 0.818 | 0.821 | 0.819 | 0.82 | 0.819 | 0.807 | 0.805 | 0.790 |