# Deep Hiearchical Multi-Label Classification Applied to Chest X-Ray Abnormality Taxonomies

Haomin Chen[a,*], Shun Miao[b], Daguang Xu[c], Gregory D. Hager[a], Adam P. Harrison[b,*]

*[a]Johns Hopkins University, Baltimore, MD, US*
*[b]PAII Inc., Bethesda, MD, US*
*[c]NVIDIA AI-Infra, Bethesda, MD, US*

**Abstract**

Chest X-rays (CXRs) are a crucial and extraordinarily common diagnostic tool, leading to heavy research for computer-aided diagnosis (CAD) solutions. However, both high classification accuracy *and* meaningful model predictions that respect and incorporate clinical taxonomies are crucial for CAD usability. To this end, we present a deep hierarchical multi-label classification (HMLC) approach for CXR CAD. Different than other hierarchical systems, we show that first training the network to model conditional probability directly and then refining it with unconditional probabilities is key in boosting performance. In addition, we also formulate a numerically stable cross-entropy loss function for unconditional probabilities that provides concrete performance improvements. Finally, we demonstrate that HMLC can be an effective means to manage missing or incomplete labels. To the best of our knowledge, we are the first to apply HMLC to medical imaging CAD. We extensively evaluate our approach on detecting abnormality labels from the CXR arm of the Prostate, Lung, Colorectal and Ovarian (PLCO) dataset, which comprises over $198,000$ manually annotated CXRs. When using complete labels, we report a mean area under the curve (AUC) of $0.887$, the highest yet reported for this dataset. These results are supported by ancillary experiments on the PadChest dataset, where we also report significant improvements, 1.2% and 4.1% in AUC and average precision, respectively over strong "flat" classifiers. Finally, we

---

*Corresponding authors. Work performed at NVIDIA
*Email addresses:* hchen135@jhu.edu (Haomin Chen ), adam.p.harrison@gmail.com (Adam P. Harrison )

demonstrate that our HMLC approach can much better handle incompletely labelled data. These performance improvements, combined with the inherent usefulness of taxonomic predictions, indicate that our approach represents a useful step forward for CXR CAD.

## 1. Introduction

Chest X-rays (CXRs) account for a large proportion of ordered image studies, *e.g.*, in the US it accounted for almost half of ordered studies in 2006 (Mettler et al., 2009). Commensurate with this importance, CXR *computer-aided diagnosis (CAD)* has received considerable research attention, both prior to the popularity of deep learning (Jaeger et al., 2013), and afterwards (Wang et al., 2017; Yao et al., 2017; Gündel et al., 2019a; Irvin et al., 2019; Bustos et al., 2019). These efforts have met success and typically approach the problem as a standard multi-label classification scenario, which attempts to make a set of individual binary predictions for each disease pattern under consideration. Yet, pushing raw performance further will likely require models that depart from standard multi-label classifiers. For instance, despite their importance to clinical understanding and interpretation (Stevens et al., 2007; Humphreys and Lindberg, 1993; Stearns et al., 2001), taxonomies of disease patterns are not typically incorporated into CXR CAD systems, or for other medical CAD domains for that matter. This observation motivates our work, which uses *hierarchical multi-label classification (HMLC)* to both push raw area under the curve (AUC) performance further and also to provide more meaningful predictions that leverage clinical taxonomies.

Organizing diagnoses or observations into ontologies and/or taxonomies is crucial within radiology, *e.g.*, RadLex (Langlotz, 2006), with CXR interpretation being no exception (Folio, 2012; Demner-Fushman et al., 2015; Dimitrovski et al., 2011). This importance should also be reflected within CAD systems. For instance, when uncertain about fine-level predictions, *e.g.*, *nodules* vs. *masses*, a CAD system should still be able to provide meaningful parent-level predictions, *e.g.*, *pulmonary nodules and*

*masses*. This parent prediction may be all the clinician is interested in anyway. Another important benefit is that observations are conditioned upon their parent being true, allowing fine-level predictors to focus solely on discriminating between siblings rather than on having to discriminate across all possible conditions. This can help improve classification performance (Bi and Kwok, 2015).

Elegantly addressing the problem of incompletely labelled data is another benefit of incorporating taxonomy. To see this, note that many CXR datasets are collected using natural language processing (NLP) approaches applied to hospital picture archiving and communication systems (PACSs) (Wang et al., 2017; Irvin et al., 2019). This is a trend that will surely increase given that PACSs remain the most viable source of large-scale medical data (Kohli et al., 2017; Harvey and Glocker, 2019). In such cases, it may not always be possible to extract fine-grained labels with confidence. For instance, imaging conditions may have only allowed a radiologist to report "opacity", instead of a more specific observation of "infiltration" vs. "atelectasis". Added to this inherent uncertainty is the fact that NLP approaches for CXR label extraction themselves can suffer from considerable levels of error and uncertainty (Irvin et al., 2019; Erdi et al., 2019). As a result, it is likely that CAD systems will increasingly be faced with incompletely labelled data, where data instances may be missing fine-grained labels, but still retain labels higher up in the clinical taxonomy. An HMLC approach can naturally handle such incompletely labelled data.

For these reasons, we present a deep HMLC approach for CXR CAD. We extensively evaluate our HMLC approach on the CXR arm of the *Prostate, Lung, Colorectal and Ovarian (PLCO)* dataset (Gohagan et al., 2000) with supporting experiments on the PadChest dataset (Bustos et al., 2019). Experiments demonstrate that our HMLC approach can push raw performance higher compared to both leading "flat" classification baselines and other HMLC alternatives. We also demonstrate that our HMLC approach can robustly handle extremely large proportions of incompletely labelled data with much less performance loss than alternatives. To the best of our knowledge, we are the first to outline an HMLC CAD system for medical imaging and the first to characterize performance when faced with incompletely labelled data.

3

*1.1. Related Work*

**CXR Classification:** Because more than one abnormality can be observed on a CXR at the same time, a CAD CXR system must operate in a multi-label setting. This is in contrast to *multi-class* approaches, which typically attempt to make a single *n*-ary prediction per image. Truly large-scale CXR classification started with the CXR14 dataset and the corresponding model (Wang et al., 2017), with many subsequents improvements both in modeling and in dataset collection (Irvin et al., 2019; Bustos et al., 2019; Johnson et al., 2019). These improvements include incorporating ensembling (Islam et al., 2017), attention mechanisms (Guan et al., 2018; Wang and Xia, 2018; Liu et al., 2019), and localizations (Yan et al., 2018; Li et al., 2018; Liu et al., 2019; Gündel et al., 2019a; Cai et al., 2018). Similar to (Gündel et al., 2019a), we also train and test on the PLCO dataset. However, (Gündel et al., 2019a) boosted their performance by incorporating the CXR14 dataset (Wang et al., 2017) and a multi-task framework that also predicted the rough locations and the lung and heart segmentations. While the contributions of these cues, when available, are important to characterize and incorporate, our HMLC approach can achieve higher AUCs[1] without extra data or auxiliary cues.

A commonality between these prior approaches is that they typically treat each label as an independent prediction, which is commonly referred to as binary relevance (BR) learning within the multi-label classification field (Zhang and Zhou, 2014). However, prior work has well articulated the limitations of BR learning (Dembczyński et al., 2012). A notable exception to this trend is (Yao et al., 2017), which modeled correlations between labels using a recurrent neural network. In contrast, our HMLC system takes a different approach by incorporating top-down knowledge to model the conditional dependence of children labels upon their parents. In this way, we make predictions conditionally independent rather than globally independent, allowing the model to focus on discriminating between siblings rather than across all disease patterns.

**Hierarchical Classification:** Given its potential to improve performance, incorporating taxonomy through hierarchical classification has been well-studied. Prior to the emergence of deep learning, seminal approaches used hierarchical and multi-label gen-

---

[1]With the caveat of using different data splits, since there is no official split.

eralizations of classic algorithms (McCallum et al., 1998; Cesa-bianchi et al., 2005; Cai, 2007; Vens et al., 2008). With the advent of deep learning, a more recent focus has been on adapting deep networks, typically convolutional neural networks (CNNs), for hierarchical classification (Redmon and Farhadi, 2017; Roy et al., 2020; Yan et al., 2015; Guo et al., 2018; Kowsari et al., 2017). Interestingly, (Cesa-bianchi et al., 2005) use an approach similar to popular approaches seen in more recent deep hierarchical *multi-class* classification of natural images (Redmon and Farhadi, 2017; Roy et al., 2020; Yan et al., 2015), *i.e.*, train classifiers to predict conditional probabilities at each node. Our approach is similar to these more recent deep approaches, except that we focus on *multi-label* classification and we also formulate a numerically stable unconditional probability fine-tuning step.

Other deep approaches used complicated combinations of CNNs and recurrent neural networks (RNNs) (Guo et al., 2018; Kowsari et al., 2017), but for our CXR application we show that a much simpler approach that uses a shared trunk network for each of the output nodes can, on its own, provide important performance improvements over "flat" classifiers.

Within medical imaging, there is work on HMLC medical image retrieval using either nearest-neighbor or multi-layer perceptrons (Pourghassem and Ghassemian, 2008) or decision trees (Dimitrovski et al., 2011). However, hierarchical classifiers have not received much attention for medical imaging *CAD* and deep HMLC approaches have not been explored at all. Finally, we note that the process of producing a set of binary HMLC labels, given a set of pseudo-probability predictions, is a surprisingly rich topic (Bi and Kwok, 2015), but here we focus on producing said predictions.

**Incompletely Labelled Data:** As mentioned, another motivating factor for HMLC is its ability to handle incompletely or partially labelled data. Within the computer vision and text mining literature, there is a rich body of work on handling partial labels (Yu et al., 2014; Kong et al., 2014; Zhao and Guo, 2015; Elkan and Noto, 2008; Liu et al., 2003; Qi et al., 2011; Bucak et al., 2011; Yang et al., 2013). When missing labels are positive examples, this problem has also been called positive and unlabelled (PU) learning. Seminal PU works focus on multi-class learning (Elkan and Noto, 2008; Liu et al., 2003). There are also efforts for *multi-label* PU learning (Kong et al., 2014;

5

Yu et al., 2014; Zhao and Guo, 2015; Qi et al., 2011; Bucak et al., 2011; Yang et al., 2013), which attempt to exploit label dependencies and correlations to overcome missing annotations. However, many of these approaches do not scale well with large-scale data (Kong et al., 2014).

(Yu et al., 2014) and (Kong et al., 2014) provide two exceptions to this, tackling large-scale numbers of labels and data instances, respectively. In our case, we are only interested in the latter, as the number of observable CXR disease patterns remains manageable. We are able to take advantage of a hierarchical clinical taxonomy to model label dependencies, allowing us to avoid complex approaches to learn these dependencies, such as the stacking methods used by (Kong et al., 2014). In this way, our approach is similar to that of (Cesa-bianchi et al., 2005), who also use a hierarchy to handle PU data through an incremental linear classification scheme. However, our approach uses deep CNNs and we are the first to show how HMLC can help address the PU problem for CAD and the first to characterize performance of CXR classifiers under this scenario.

*1.2. Contributions*

Based on the above, the contributions of our work can be summarized as follows:

- Like other deep hierarchical *multi-class* classifiers, we train a classifier to predict conditional probabilities. However, we operate in the *multi-label* space and we also demonstrate that a second fine-tuning stage, trained using unconditional probabilities, can boost performance for CXR classification even further.

- To handle the unstable multiplication of prediction outputs seen in unconditional probabilities we introduce and formulate a numerically stable and principled loss function.

- Using our two-stage approach, we are the first to apply hierarchical multi-label classification (HMLC) to CXR CAD. Our straightforward, but effective, HMLC approach results in the highest mean AUC value yet reported for the PLCO dataset.

6

- In addition, we demonstrate how HMLC can serve as an effective means to handle incompletely labelled data. We are the first to characterize CXR classification performance under this scenario, and experiments demonstrate how HMLC can garner even greater boosts in classification performance.

Finally, we note that portions of this work were previously published as a conference proceeding (Chen et al., 2019). This work adds several contributions: (1) we significantly expand upon the literature review; (2) we include the derivation of the numerically stable unconditional probability loss within the main body and have made its derivation clearer; (3) we include additional results with the PadChest Bustos et al. (2019) dataset to further validate our approach; and (4) we add the motivation, discussion, and experiments on incompletely labelled data.

## 2. Materials and Methods

We introduce a two-stage method for CXR HMLC. We first outline the datasets and taxonomy we use in Section 2.1 and then overview the general concept of HMLC in Section 2.2. This is followed by Sections 2.3 and 2.4, which detail our two training stages that use conditional probability and a numerically stable unconditional probability formulation, respectively.

### 2.1. Datasets and Taxonomy

The first step in creating an HMLC system is to create the label taxonomy. In this work, our main results focus on the labels and data found within the CXR arm of the PLCO dataset (Gohagan et al., 2000), a large-scale lung cancer screening trial that collected 198 000 CXRs with image-based annotations of abnormalities obtained from multiple US clinical centers. While other large-scale datasets (Wang et al., 2017; Bustos et al., 2019; Irvin et al., 2019; Johnson et al., 2019) are *extraordinarily valuable*, their labels are generated by using NLP to extract mentioned disease patterns from radiological reports found in hospital PACSs. While medical NLP has made great strides in recent years, it still remains an active field of research, *e.g.*, NegBio still reports limitations with uncertainty detection, double-negation, and missed positive findings

7

for certain CXR terms (Peng et al., 2018). However, irrespective of the NLP's level of accuracy, there are more inherent limitations to using text-mined labels. Namely, examining a text report is no substitute for visually examining the actual radiological scan, as the text of an individual report is not a complete description of the CXR study in question. Thus, terms may not be mentioned, *e.g.*, "no change", even though they are indeed visually apparent. Additionally, a radiologist will consider lab tests, prior radiological studies, and the patient's records when writing up a report. Thus, mentioned terms, and their meaning, may well be influenced by factors that are not visually apparent. Compounding this, text which is unambiguous given the patient's records and radiological studies may be highly ambiguous when only considering text alone, *e.g.*, whether a pneumothorax is untreated or not (Oakden-Rayner, 2019). Indeed, the authors of the PadChest dataset bring up some of these caveats themselves, which are relevant even for the 27% of their radiological reports that are text-mined by hand, which presumably have no NLP errors Bustos et al. (2019). An independent study of CXR14 (Wang et al., 2017) concludes that its labels have low positive predictive value and argues that visual inspection is necessary to create radiological datasets (Oakden-Rayner, 2019). Consequently, PLCO is unique in that it is the only large-scale CXR dataset with labels generated via *visual observation* from radiologists. Although the PLCO data is older than alternatives (Wang et al., 2017; Bustos et al., 2019; Irvin et al., 2019; Johnson et al., 2019), it has greater label reliability.

Radiologists in the PLCO trial labelled 15 disease patterns, which we call "leaf labels" in our taxonomy. Because of low prevalence, we merged "left hilar abnormality" and "right hilar abnormality" into "hilar abnormality", resulting in 14 labels. From the leaf nodes, we constructed the label taxonomy shown in Figure 1. The hierarchical structure follows the PLCO trial's division of "suspicious for cancer" disease patterns vs. not, and is further partitioned using common groupings (Folio, 2012), totalling 19 leaf and non-leaf labels. While care was taken in constructing the taxonomy and we aimed for clinical usefulness, we make no specific claim as such. We instead use the taxonomy to explore the benefits of HMLC, stressing that our approach is general enough to incorporate any appropriate taxonomy. Figure 2 visually depicts examples from our chosen CXR taxonomy.
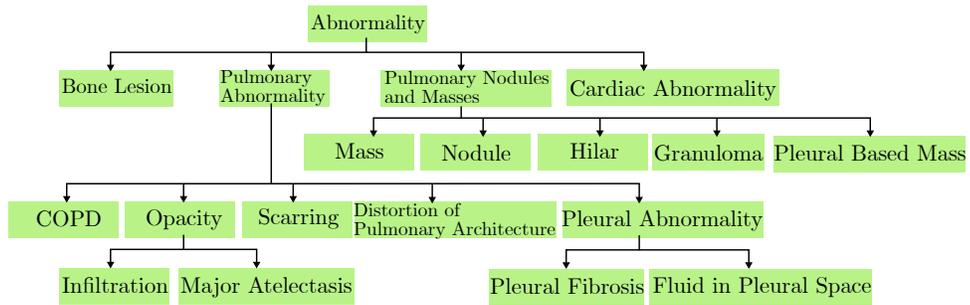
8

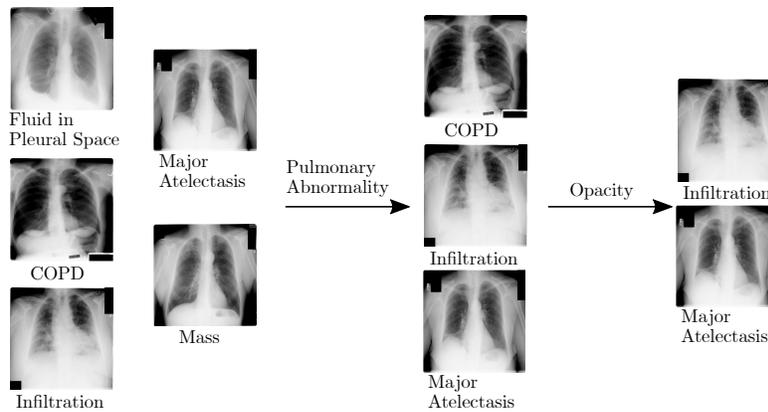Figure 1: Constructed label hierarchy from the PLCO dataset.



Figure 2: Example PLCO CXRs drawn from three levels of our taxonomy. On the left, at the higest level of taxonomy, *i.e.*, "Abnormality", disease patterns may manifest as a variety of visual features within the lung parenchyma, lung pleura, or the surrounding organs/tissues. As one progresses down the taxonomy, *i.e.*, to "Opacity", the discriminating task is narrowed into identifying the "cloudy" patterns seen in both "Infiltration" and "Major Atelectasis."
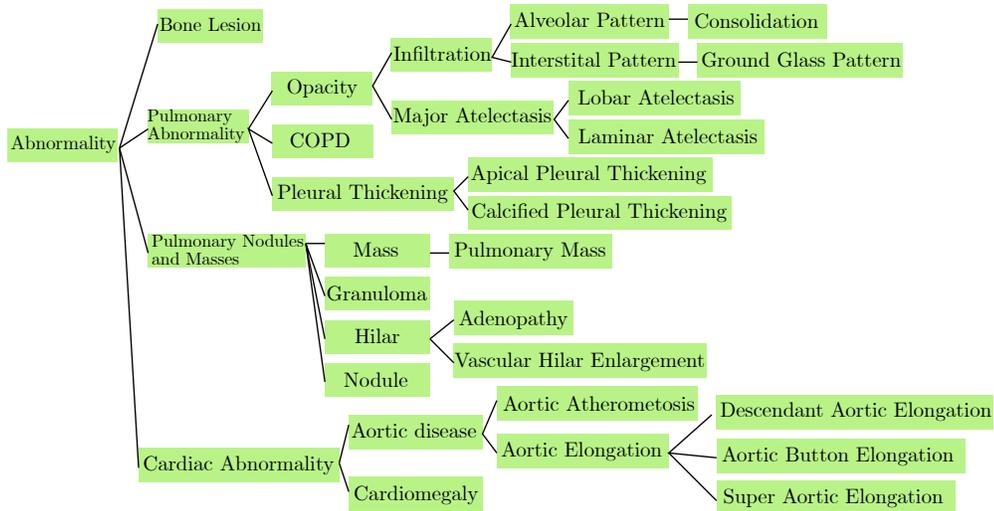
Figure 3: Constructed label hierarchy from the PadChest dataset.

As supporting validation to our main PLCO experiments, we also validate on the PadChest dataset Bustos et al. (2019), which contains $160,845$ CXRs whose labels are drawn from either manual or automatic extraction from radiological *text reports*. We focus on labels categorized as "radiological findings", which are more likely to correspond to actual disease patterns found on the CXRs Bustos et al. (2019). Any CXR with a solitary "Unchanged" label is removed, resulting in $121,242$ samples. Uniquely, PadChest offers a complete hierarchical structure for all labels. We remove labels with less than 100 manually labelled samples and only retain labels that align with our PLCO taxonomy. This both ensures we have enough statistical power for evaluation and that we are retaining PLCO-like terms that we can confidently treat as clinically significant. As a result, total 30 out of 191 labels are selected, and our supplementary includes more details of the included and excluded labels. The resulting taxonomy is shown in Figure 3. Unlike PLCO, certain parent labels can be positive with no positive children labels, *e.g.*, "Aortic Elongation".

## 2.2. Hierarchical Multi-Label Classification

With a taxonomy established, a hierarchical approach to classification must be established. Because this is a multi-label setting, all or none of the labels in Figure 1 can

be positive. The only restriction is that if a child is positive, its parent must be too. Siblings are not mutually exclusive. For PLCO, we assume that each image is associated with a set of ground-truth leaf labels and their antecedents, *i.e.*, there are no incomplete paths. However, for PadChest a ground-truth path may terminate before a leaf node. A training set, may have missing labels.

We use a DenseNet-121 (Huang et al., 2017) model as a backbone. If we use $k$ to denote the total number of leaf and non-leaf labels, we connect $k$ fully connected layers to the backbone's last feature layer to extract $k$ scalar outputs. Each output is assumed to represent the conditional probability (or its logit) given its parent is true. Thus, once the model is successfully trained, unconditional probabilities can be calculated from the output using the chain rule, *e.g.*, from the PLCO taxonomy the unconditional probability of *scarring* can be calculated as

$$P(\text{Scar.}) = P(\text{Abn.})P(\text{Pulm.}|\text{Abn.})P(\text{Scar.}|\text{Pulm.}), \qquad (1)$$

where we use abbreviations for the sake of typesetting. In this way, the predicted unconditional probability of a parent label is guaranteed to be greater than or equal to its children labels. We refer to the conditional probability in a label hierarchy as hierarchical label conditional probability (HLCP), and the unconditional probability calculated following the chain rule as hierarchical label unconditional probability (HLUP). The network outputs can be trained either conditionally or unconditionally, which we outline in the next two sections.

### 2.3. Training with Conditional Probability

Similar to prior work (Redmon and Farhadi, 2017; Roy et al., 2020; Yan et al., 2015), in the first stage of the proposed training scheme, each classifier is only trained on data conditioned upon its parent label being positive. Thus, training directly models the conditional probability. The shared part of the classifiers, *i.e.*, feature layers from the backbone network, is trained jointly by all the tasks. Specifically, for each image the losses are only calculated on labels whose parent label is also positive. For example, and once again using the PLCO taxonomy, when an image with positive *Scarring* and no other positive labels is fed into training, only the losses of *Abnormality* and the
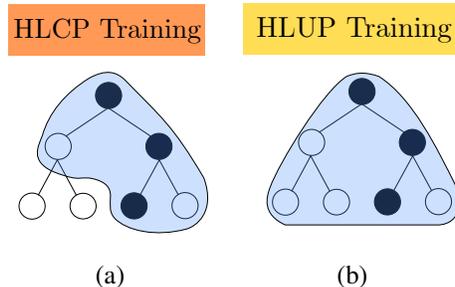
Figure 4: The HLCP and HLUP losses are depicted in (a) and (b), respectively, where black and white points are positive and negative labels, respectively. Blue areas indicate the activation area in the loss functions.

children labels of *Pulmonary Abnormality* and *Abnormality* are calculated and used for training.

Figure 4 (a) illustrates this training regimen, which we denote HLCP training. In this work, we use cross entropy (CE) loss to train the conditional probabilities, which can be written as

$$L_{HLCP} = \sum_{m \in M} CE\left(z_m, \hat{z}_m\right) * 1_{\{z_{a(m)}=1\}}, \tag{2}$$

where $M$ denotes the set of all disease patterns, and $m$ and $a(m)$ denote a disease pattern and its ancestor, respectively. Here $CE(\cdot, \cdot)$ denotes the cross entropy loss, and $z_m \in \{0, 1\}$ denotes the ground truth label of $m$, with $\hat{z}_m$ corresponding to the network's sigmoid output.

Training with conditional probability is a very effective initialization step, as it concentrates the modeling power solely on discriminating siblings under the same parent label, rather than having to discriminate across all labels, which eases convergence and reduces confounding factors. It also alleviates the problem of low label prevalence because fewer negative samples are used for each label.

### 2.4. Fine Tuning with Unconditional Probability

In the second stage, we finetune the model using an HLUP CE loss. This stage aims at improving the accuracy of unconditional probability predictions, which is what is actually used during inference and is thus critical to classification performance. Another

12

important advantage is that the final linear layer sees more negative samples. Predicted unconditional probabilities for label $m$, denoted $\hat{p}_m$, are calculated using the chain rule:

$$\hat{p}_m = \prod_{m' \in A(m)} \hat{z}_{m'}, \tag{3}$$

where $A(m)$ is the union of label $m$ and its antecedents. When training using unconditional probabilities, the loss is calculated on every classifier output for every data instance. Thus, the HLUP CE loss for each image is simply

$$L_{HLUP} = \sum_{m \in M} CE\left(z_m, \hat{p}_m\right). \tag{4}$$

Figure 4(b) visually depicts this loss.

A naive way to calculate (4) would be a direct calculation. However, such an approach introduces instability during optimization, as the training would have to minimize the product of network outputs. In addition, the product of probability values within $[0, 1]$ can cause arithmetic underflow. For this reason, we derive a numerically stable formulation below.

Denoting the network's output logits as $\hat{y}_{(.)}$, the predicted unconditional probability of label $m$ can be written as:

$$\hat{p}_m = \prod_{m'} \frac{1}{1 + \exp(-y_{m'})}, \tag{5}$$

where we use $m'$ to denote $m' \in A(m)$ for notational simplicity.

The HLUP CE loss is calculated as:

$$L_{HLUP} = - z_m \log(\hat{p}_m) - (1 - z_m) \log(1 - \hat{p}_m), \tag{6}$$

$$= - z_m \log\left(\prod_{m'} \frac{1}{1 + \exp(-y_{m'})}\right)$$

$$- (1 - z_m) \log\left(1 - \left(\prod_{m'} \frac{1}{1 + \exp(-y_{m'})}\right)\right), \tag{7}$$

where $z_m$ is the ground truth label of $m$.

The formulation in (7) closely resembles several cross-entropy loss terms combined together. To see this, we can break up the second term in (7) to produce the following

13

formulation:

$$L_{HLUP} = -z_m \log\left(\prod_{m'} \frac{1}{1 + \exp(-y_{m'})}\right)$$

$$- (1 - z_m) \log\left(\prod_{m'}\left(1 - \frac{1}{1 + \exp(-y_{m'})}\right)\right) + \gamma, \tag{8}$$

where $\gamma$ is a scalar quantity that must be formulated. The log terms above can then be decomposed as

$$L_{HLUP} = \sum_{m'}\left(-z_m \log\left(\frac{1}{1 + \exp(-y_{m'})}\right)\right.$$

$$\left. -(1 - z_m) \log\left(1 - \frac{1}{1 + \exp(-y_{m'})}\right)\right) + \gamma, \tag{9}$$

$$= \sum_{m'} \ell_{m'} + \gamma, \tag{10}$$

where $\ell_m$ are individual cross entropy terms, using $z_m$ and $y_{m'}$ as the ground truth and logit input, respectively. Note that (10) allows us to take advantage of numerically stable CE implementations to calculate $\sum_{m'} \ell_{m'}$. However to satisfy (10), we will need $\gamma$ to satisfy:

$$\gamma = (1 - z_m) \log\left(\prod_{m'}\left(1 - \frac{1}{1 + \exp(-y_{m'})}\right)\right)$$

$$- (1 - z_m) \log\left(1 - \left(\prod_{m'} \frac{1}{1 + \exp(-y_{m'})}\right)\right), \tag{11}$$

$$= (1 - z_m) \log\left(\frac{\prod_{m'} \exp(-y_{m'})}{\prod_{m'}(1 + \exp(-y_{m'}))}\right)$$

$$- (1 - z_m) \log\left(\frac{\prod_{m'}(1 + \exp(-y_{m'})) - 1}{\prod_{m'}(1 + \exp(-y_{m'}))}\right), \tag{12}$$

$$= (1 - z_m) \log\left(\frac{\exp(\sum_{m'} -y_{m'})}{\prod_{m'}(1 + \exp(-y_{m'})) - 1}\right), \tag{13}$$

$$= (1 - z_m)\left(\sum_{m'} -y_{m'} - \log\left(\prod_{m'}(1 + \exp(-y_{m'})) - 1\right)\right). \tag{14}$$

If the product within the log-term of (14) is expanded, with 1 subtracted, it will result in

$$\gamma = (1 - z_m)\left(\sum_{m'} -y_{m'} - \log\left(\sum_{S \in \mathcal{P}(A(m))\setminus\{\emptyset\}} \exp\left(\sum_{j \in S} -y_j\right)\right)\right), \tag{15}$$

14

where $S$ enumerates all possible subsets of the powerset of $A(m)$, excluding the empty set. For example if there were two logits, $y_1$ and $y_2$, the summation inside the log would be:

$$\exp(-y_1) + \exp(-y_2) + \exp(-y_1 - y_2). \tag{16}$$

The expression in (15) can be written as

$$\gamma = (1 - z_m)\left(\sum_{m'} -y_{m'} - LSE\left(\left\{\sum_{j \in S} -y_j \quad \forall S \in \mathcal{P}(A(m)) \setminus \{\emptyset\}\right\}\right)\right), \tag{17}$$

where $LSE$ is the LogSumExp function. Numerically stable implementations of the LogSumExp, and its gradient, are well known. By substituting (17) into (10), a numerically stable version of the HLUP CE loss can be calculated.

Enumerating the powerset produces an obvious combinatorial explosion. However, for smaller-scale hierarchies, like that in Figure 1, it remains tractable. For larger hierarchies, an $O(|A(m)|)$ solution involves simply interpreting the LogSumExp as a smooth approximation to the maximum function, which we provide here for completeness:

$$\gamma \approx (1 - z_m)\left(\sum_{m'} -y_{m'} - \max\left(\left\{\sum_{j \in S} -y_j \quad \forall S \in \mathcal{P}(A(m)) \setminus \{\emptyset\}\right\}\right)\right), \tag{18}$$

$$= \begin{cases} (1 - z_m)\left(\sum_{m'} -y_{m'} - \sum_{j:y_j<0} -y_j\right), & \text{if } \exists\, y_{m'} < 0 \\ (1 - z_m)\left(\sum_{m'} -y_{m'} - \max(\{-y_{m'}\})\right), & \text{otherwise} \end{cases}. \tag{19}$$

## 3. Experimental

We perform two types of experiments to validate our HMLC approach. The first uses the standard completely labelled setup, helping to reveal how our use of taxonomic classification can help produce better raw classification performance than typical "flat" classifiers. The second uses incompletely labelled data under controlled scenarios to show how our HMLC approach can naturally handle such data, achieving even higher boosts in relative performance.

### 3.1. Complete Labels

**Experimental Setup** We test our HMLC approach on both the PLCO Gohagan et al. (2000) and PadChest Bustos et al. (2019) datasets, using the taxonomies of Figure 1

and Figure 3, respectively. Our emphasis is on PLCO due to its more reliable labels, but evaluations on PadChest provide important experimental support, especially given its larger taxonomy. Following accepted practices in large-scale CXR classification Wang et al. (2017); Irvin et al. (2019); Bustos et al. (2019), we split the data into single training, validation, and test sets, corresponding to 70%, 10%, and 20% of the data, respectively. Data is split at the patient level, and care was taken to balance the prevalence of each disease pattern as much as possible. As mentioned above, our HMLC approach uses a trunk network, with a final fully-connected layer outputting logit values for each of the nodes of our chosen taxonomy. Our chosen network is DenseNet-121 (Huang et al., 2017), implemented using TensorFlow. We first train with the HLCP CE loss of (2) fine-tuning from a model pretrained from ImageNet (Deng et al., 2009). We refer to this model simply as *HLCP*. To produce our final model, we then finetune the HLCP model using the HLUP CE loss of (4). We denote this final model as *HLUP-finetune*.

**Comparisons** In addition to comparing against HLCP, we also compare against three other baseline models, all using the same trunk network fine-tuned from ImageNet pretrained weights. The first, denoted *BR-leaf*, is trained using CE loss on the 14 fine-grained labels. This measures performance using a standard multi-label BR approach. The second, denoted *BR-all* is very similar, but trains a CE loss on all labels independently, including non-leaf ones. In this way, *BR-all* measures performance when one wishes to naively output non-leaf abnormality nodes, without considering label taxonomy. Finally, we also test against a model trained using the HLUP CE loss directly from ImageNet weights, rather than finetuning from the HLCP model. As such, this baseline, denoted *HLUP*, helps reveal the impact of using a two-stage approach vs. simply training an HLUP classifier in one step. For all tested models, extensive hyper-parameter searches were performed on the NVIDIA cluster to optimize mean validation fine-grained AUCs.

For comparisons to external models, we also compare to a recent DenseNet121 BR approach (Gündel et al., 2019a) trained on the PLCO data. But, we stress that direct comparisons of numbers are impossible, as (Gündel et al., 2019a) used different data splits and only evaluated on 12 fine-grained labels. In the interest of fairness we compare against both (a) their best reported numbers when only training a classifier on

16

CXR disease patterns and (b) their best reported numbers overall, in which the authors incorporated segmentation and localization cues. For (a), we use numbers reported on an earlier work (Gündel et al., 2019b), which were higher. Unfortunately, both sets of their reported numbers are based on training data that also included the ChestXRay14 dataset (Wang and Xia, 2018), providing an additional confounding factor that hampers any direct comparison.

Finally, we also run experiments to compare our numerically stable implementation of HLUP CE loss in (8) to: (a) the naive approach of directly optimizing (3); and (b) to a recent rescaling approximation, originally introduced for the multiplication of independent, rather than conditional probabilities, seen in multi-instance learning (Li et al., 2018). This latter approach re-scales each individual probability multiplicand (term) in (3) to guarantee that the product is greater than or equal to 1e-7. Similar to the naive approach, the product is then optimized directly using CE loss. For the PLCO dataset, based on a maximum depth of four for the taxonomy, we implement this approach by re-scaling each multiplicand in (3) to [0.02, 1].

**Evaluation Metrics** We evaluate our approach using AUC and average precision (AP), calculated across both leaf and non-leaf labels, when applicable. Additionally, we also evaluate using conditional AUC and AP metrics, which are metrics that reflect the complicated evaluation space of multi-label classification. In short, because more than one label can be positive, multi-label classification performance has exponentially more facets for evaluation than single-label or even multi-class settings. Conditional metrics are one such facet, that focus on model performance conditioned on certain non-leaf labels being positive. Here, we restrict our focus to CXRs exhibiting one or more disease patterns, *i.e.*, *abnormality* being positive. As such, this sheds light on model performance when it may be critical to discriminate what combination of disease patterns are present, which is crucial for proper CXR interpretation (Folio, 2012).

### 3.2. Incomplete Labels

**Experimental Setup** We also use the PLCO dataset (Gohagan et al., 2000) to characterize the benefits of our HMLC approach when faced with incomplete labels. However, after publication of our original work (Chen et al., 2019), the PLCO organizers

altered their data release policies and only released a subset of the original dataset, containing 88 737 labeled CXRs from 24 997 patients[2]. For this reason, we perform our incomplete labels experiments on this smaller dataset, splitting and preparing the data in an identical manner as described in Section 3.1.

To simulate a scenario where learning algorithms may be faced with incomplete labels, we removed known labels from the training set using the following controlled scheme:

1. We choose a base deletion probability, $\beta \in [0, 1]$.

2. For data instances with positive labels for "Pleural Abnormality", "Opacity", and "Pulmonary Nodules and Masses", we delete all their children labels with a probability of $\beta$. For example, if we delete the children labels of a positive "Pleural Abnormality" instance, then it is no longer known whether the "Pleural Abnormality" label corresponds to "Pleural Fibrosis", or "Fluid in Pleural Space", or both.

3. We perform the same steps for data instances with positive labels for "Pulmonary Abnormality" and "Abnormality", except with probabilities of $0.3\beta$ and $0.3^2\beta$, respectively. For example, if the children of a positive instance of "Abnormality" were deleted, then it is only known there are one or more disease patterns present, but not which one(s).

4. A higher-level deletion overrides any decision(s) at finer levels.

5. Because of their extremely low prevalence, we ignore the "Major Atelectasis" and "Distortion in Pulmonary Architecture" labels in training and evaluation.

Note that this scheme makes it more likely to have a missing fine-grained label over a higher-level label, which we posit follows most scenarios producing incomplete labels. When labels are deleted, we treat them as unknown and do not execute any training loss on them. We test our HMLC algorithm and baselines on the following $\beta$ values: $\{0, .1, .2, .3, .4, .5, .6, .7\}$, which ranges from no incompleteness to roughly 70% of fine-grained labels being deleted. To allow for stable comparisons across $\beta$ values, we also

---

[2]The first author no longer had access to the original dataset for the incomplete label experiments as he had finished his internship at NVIDIA.

ensure that if a label was deleted at a certain value of $\beta$, it will also be deleted at all higher values of $\beta$. To ease reproducibility, we publicly release our data splits[3]. All other implementation details are also identical to that of Section 3.1.

**Evaluation Metrics and Comparisons** We measure AUC values and compare our chosen model of HLUP finetune against BR-leaf and BR-all.

## 4. Results and Discussion

We focus in turn on experiments with complete and incomplete labels, which can be found in Section 4.1 and Section 4.2, respectively.

### 4.1. Complete Labels

Our complete labels experiments first focus on the benefits of our HLUP-finetune approach compared to alternative "flat" and HMLC strategies. Then, we discuss results specifically focusing on our numerically stable HLUP CE loss.

### 4.1.1. HLUP-finetune Performance

Table 1 outlines the PLCO results of our HLUP-finetune approach vs. competitors. As the table demonstrates, the standard baseline BR-leaf model produces high AUC scores, in line with prior work (Gündel et al., 2019b); however, it does not provide high-level predictions based on a taxonomy. Naively executing BR training on the entire taxonomy, *i.e.*, the BR-all model, does not improve performance. This indicates that if not properly incorporated, the label taxonomy does not benefit performance.

In contrast, the HLCP model is indeed able to match BR-leaf's performance on the fine-grained labels, despite also being able to provide high-level predictions. HLUP-finetune goes further by exceeding BR-leaf's fine-grained performance, demonstrating that our two-stage training process can produce tangible improvements. This is underscored when comparing HLUP-finetune with HLUP, which highlights that without the two-stage training, HLUP training cannot reach the same performance. If we limit ourselves to models incorporating the entire taxonomy, our final HLUP-finetune

---

[3]https://github.com/hchen135/Hierarchical-Multi-Label-Classification-X-Rays

Table 1: *PLCO* AUC and AP values across tested models. Mean values across leaf and non-leaf disease patterns are shown, as well as for leaf labels conditioned on one or more abnormalities being present.

|  | Leaf labels | | Non-leaf labels | | Leaf labels conditioned on abnormality | |
|---|---|---|---|---|---|---|
|  | AUC | AP | AUC | AP | AUC | AP |
| (Gündel et al., 2019b) | 0.865 | N/A | N/A | N/A | N/A | N/A |
| (Gündel et al., 2019a) | 0.883 | N/A | N/A | N/A | N/A | N/A |
| BR-leaf | 0.871 | 0.234 | N/A | N/A | 0.806 | 0.334 |
| BR-all | 0.867 | 0.221 | 0.852 | 0.440 | 0.808 | 0.323 |
| HLUP | 0.872 | 0.214 | 0.856 | 0.436 | 0.799 | 0.288 |
| HLCP | 0.879 | 0.229 | 0.857 | 0.440 | 0.822 | 0.329 |
| HLUP-finetune | **0.887** | **0.250** | **0.866** | **0.460** | **0.832** | **0.342** |

model outperforms BR-all by 2% and 2.9% in leaf-label mean AUC and AP values, respectively. Because HLUP-finetune shares the same labels as BR-all, the performance boosts of the former over the latter demonstrate that the additional output nodes seen in HMLC are not responsible for performance increases. Instead, it is indeed the explicit incorporation of taxonomic structure that leads to improved performance.

Figure 5 provides more details on these improvements, demonstrating that AUC values are higher for HLUP-finetune compared to the baseline method for all fine-grained and high-level disease patterns. Interested readers can find these AUC values in our supplementary materials. Although not graphed here for clarity reasons, HLUP-finetune also outperformed the HLCP method for all disease patterns. Of note is that statistically significant differences also respect the disease hierarchy, and if a child disease pattern demonstrates statistically significant improvement, so does its parent.

Of particular note, when considering AUCs conditioned on one or more abnormalities being present (last column of Table 1), the gap between all HMLC approaches and "flat" classifiers increases even more. As can be seen in such settings, HLUP-finetune still exhibits increased performance over the baseline models and also the next-best hierarchical model. Importantly, if we compare the conditional AUCs between BR-all
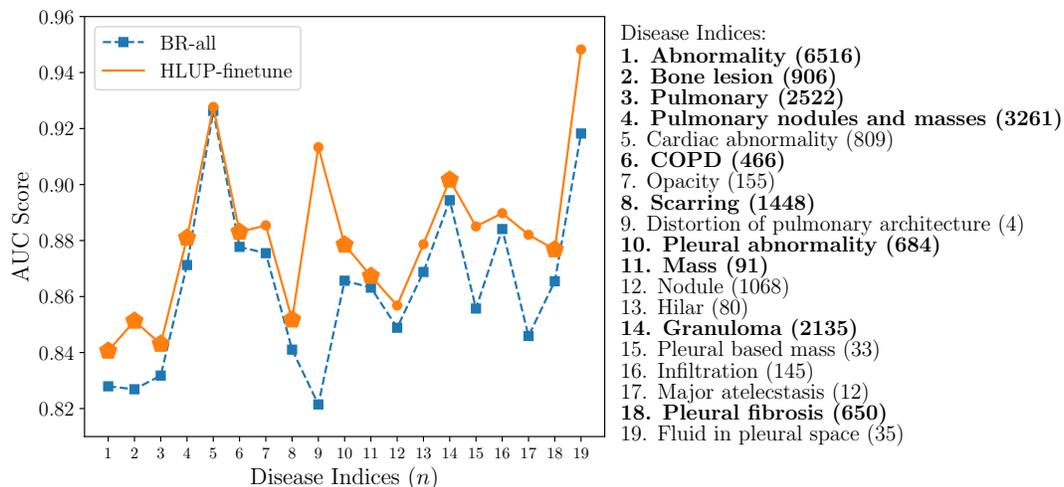
Figure 5: Comparison of AUC scores for all fine-grained and high-level (non-leaf) disease patterns for the BR-all and HLUP-finetune models. The dashed line separates the fine-grained from the high-level (non-leaf) disease patterns. Boldface labels and larger graph markers denote disease patterns exhibiting statistically significant improvement ($p < 0.05$) using the StAR software implementation (Vergara et al., 2008) of the non-parametric test of (DeLong et al., 1988).

and HLUP-finetune, we see a 2.4% increase. This indicates that HMLC is particularly effective at differentiating the exact combination of abnormalities present within an image. This may reduce the amount of spurious and distracting predictions upon deployment, but more investigation is required to quantify this.

We also note that HLUP-finetune managed to outperform (Gündel et al., 2019a)'s AUC numbers, despite the latter incorporating almost twice the amount of data and also including additional localization and segmentation tasks. However, we again note that (Gündel et al., 2019a) used a different data split and only 12 fine-grained labels, so such comparisons can only be taken so far.

Experiments on PadChest further support these results, with trends mirroring that of the PLCO experiments. As can be seen in Table 2, HLUP-finetune outperforms both the BR baselines and HMLC alternatives. Moreover, just like the PLCO experiments, when evaluating AUC and AP conditioned on one or more abnormalities being present, the performance gaps between HLUP-finetune and alternatives further increase. The

Table 2: *PadChest* AUC and AP values across tested models. Mean values across leaf and non-leaf disease patterns are shown, as well as for leaf labels conditioned on one or more abnormalities being present.

| | Leaf labels | | Non-leaf labels | | Leaf labels conditioned on abnormality | |
|---|---|---|---|---|---|---|
| | AUC | AP | AUC | AP | AUC | AP |
| BR-leaf | 0.825 | 0.104 | N/A | N/A | 0.743 | 0.212 |
| BR-all | 0.825 | 0.110 | 0.820 | 0.221 | 0.739 | 0.204 |
| HLUP | 0.831 | 0.114 | 0.828 | 0.220 | 0.752 | 0.211 |
| HLCP | 0.831 | 0.135 | 0.833 | 0.240 | 0.765 | 0.244 |
| HLUP-finetune | **0.837** | **0.145** | **0.840** | **0.253** | **0.778** | **0.261** |

relative performance improvements demonstrate that our HMLC approach generalizes well to a different CXR dataset outside of PLCO, even though PadChest uses a different taxonomy and was collected with very different patient populations at a much later date.

The PLCO and PadChest performance boosts are in line with prior work that reported improved classification performance when exploiting taxonomy, *e.g.*, for text classification (McCallum et al., 1998; Dumais and Chen, 2000), but here we use HMLC in a more modern deep-learning setting and for an imaging-based CAD application. In particular, given that taxonomy and ontology are crucial within medicine, the use of hierarchy is natural. Because the algorithmic approach we take remains very simple, our HMLC approach may be an effective method for many other medical classification tasks outside of CXRs.

The discussion of the performance boosts garnered by HMLC are very important, but it should also be noted that HMLC provides inherent benefits outside of raw classification performance. By ensuring that clinical taxonomy is respected, *i.e.*, a parent label's pseudo-probability will always be greater than or equal to any of its children's, HMLC provides a more interpretable and understandable set of predictions that better match the top-down structure of medical ontology.

In addition to exploring the benefits of the conceptual approach of HMLC to CXR

Table 3: Comparison of AUCs produced using different HLUP CE loss implementations for PLCO.

| HLUP (naive) | HLUP (rescale) | HLUP (ours) | HLUP-finetune (naive) | HLUP-finetune (rescale) | HLUP-finetune (ours) |
|---|---|---|---|---|---|
| 0.864 | 0.853 | 0.872 | 0.886 | 0.867 | 0.887 |

classification, our work also demonstrates that a two-stage HLUP finetuning approach can provide performance boosts over the more common one-stage HLCP training seen in many prior deep-learning works (Redmon and Farhadi, 2017; Roy et al., 2020; Yan et al., 2015). As such, our two-stage approach may also prove useful to hierarchical classifiers seen in other domains, such as computer vision or text classification.

### 4.1.2. Numerically Stable HLUP

Table 3 demonstrates that our numerically stable HLUP CE loss results in much better AUCs compared to the competitor rescaling approach (Li et al., 2018) and to naive HLUP training when starting from ImageNet weights. However, there were no performance improvements when compared to the naive approach when finetuning from the HLCP weights. We hypothesize that the predictions for the HLCP are already at a sufficient quality that the numerical instabilities of the naive HLUP CE loss are not severe enough to impair performance. Nonetheless, given the improvements when training from ImageNet weights, these results indicate that our HLCP CE loss does indeed provide tangible improvements in convergence stability. We expect these improvements to be greater given taxonomies of greater depth, and our formulation should also prove valuable to multi-instance setups which must optimize CE loss over the product of large numbers of probabilities, *e.g.*, the 256 multiplicands seen in (Li et al., 2018).

### 4.2. Incomplete Labels

Figure 6 shows the results of our incompletely labelled experiments. As can be seen when all labels are present, *i.e.*, $\beta = 0$, the results mirror that of Section 4.1, with HLUP-finetune outperforming the baseline models and the BR-all providing no
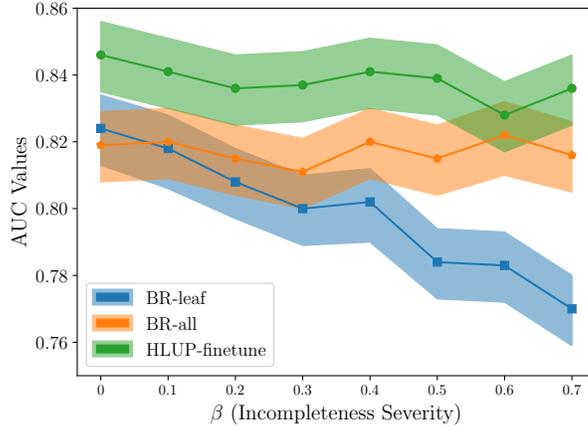
Figure 6: Mean AUC scores under different levels of label incompleteness with confidence intervals representing the 2.5th and 97.5th percentiles of 5000 resampling with replacement bootstrap rounds (Dekking et al., 2005).

improvements over BR-leaf. As the incompleteness severity increases, BR-leaf's performance drastically drops, while BR-all and HLUP-finetune are much better able to manage label incompleteness. At the highest $\beta$ level, the performance gap between HLUP-finetune and BR-leaf almost reaches 7%. Per-abnormality AUC values can be found in our supplementary materials.

Our results demonstrate that incorporating hierarchy can be an effective means to manage incomplete labels. Specifically, while HLUP-finetune's performance does indeed drop as the incompleteness severity increases, it does so at a drastically reduced rate compared to the standard BR-leaf classifier. Interestingly, BR-all, which trains all outputs but without incorporating a taxonomy, also manages to retain an equally graceful performance drop. However, HLUP-finetune's roughly 2% AUC performance advantage over BR-all indicates that properly incorporating the taxonomic hierarchy is necessary to boost classification performance. We suspect the anomaly at $\beta = 0.6$ is due to variability caused by the randomness of the training procedure and we reran our experiments at this $\beta$ value which confirmed this. Ideally, running multiple training runs at each $\beta$ value would allow us to produce confidence bars that take into account

24

effects from random weight initialization and sampling, but time and computational resources did not allow us to perform this extremely demanding set of experiments. Finally, HLUP-finetune has the added important benefit of producing predictions that respect the taxonomy, which is something that BR-all does not do. Thus, these results indicate that when possible, incorporating a HMLC approach can be an effective means to manage incompletely labelled data. As the prevalence of text-mined PACS medical imaging data increases, we expect the need for approaches to gracefully handle missing labels to increase, and our HMLC approach may provide a useful cornerstore of future work in this direction.

## 5. Conclusions

We have presented a two-stage approach for deep HMLC of CXRs that combines conditional training with an unconditional probability fine-tuning step. To effect the latter, we introduce a new and numerically stable formulation for HLUP CE loss, which we expect would also prove valuable in other training scenarios involving the multiplication of probability predictions, *e.g.*, multi-instance learning. Through comprehensive evaluations, we report the highest mean AUC on the PLCO dataset yet, outperforming hierarchical and non-hierarchical alternatives. Supporting experiments on the PadChest dataset confirm these results. We also show performance improvements conditioned on one or more abnormalities being present, *i.e.*, predicting the specific combination of disease patterns, which is crucial for CXR interpretation. Experiments with incompletely labelled data also demonstrate that our two-stage HMLC approach is an effective means to handle missing labels within training data.

There are several interesting avenues of future work. For instance, while the straightforward HMLC approach we take enjoys the virtue of being easy to implement and tune, it is possible that more sophisticated approaches, *e.g.*, using hierarchical features or dedicated classifiers, may garner even further improvements. Prior work using classic, non deep-learning approaches, explored these options McCallum et al. (1998); Cesa-bianchi et al. (2005); Dumais and Chen (2000); Cai (2007); Vens et al. (2008), and their insights should be applied today. Another important topic of future work should

be on incorporating uncertainty within HMLC. This would allow a model, when appropriate, to predict high confidence for non-leaf label predictions but lower confidence for leaf label predictions, enhancing its usefulness in deployment scenarios. Future work should also consider applications outside of CXRs both within and without medical imaging, *e.g.*, genomics or proteomics. Finally, one issue for further investigation is to better understand the implications of the annotation noise described by (Gündel et al., 2019a), both for training and for evaluation. Relevant to this work, assessing label noise at higher levels of hierarchy should be an important focus going forward.

## Acknowledgements

## References

Bi, W., Kwok, J.T., 2015. Bayes-optimal hierarchical multilabel classification. IEEE Transactions on Knowledge and Data Engineering 27, 2907–2918. doi:10.1109/TKDE.2015.2441707.

Bucak, S.S., Jin, R., Jain, A.K., 2011. Multi-label learning with incomplete class assignments, in: CVPR 2011, IEEE, Colorado Springs, CO, USA. pp. 2801–2808. doi:10.1109/CVPR.2011.5995734.

Bustos, A., Pertusa, A., Salinas, J.M., de la Iglesia-Vayá, M., 2019. Padchest: A large chest x-ray image dataset with multi-label annotated reports. ArXiv abs/1901.07441.

Cai, J., Lu, L., Harrison, A.P., Shi, X., Chen, P., Yang, L., 2018. Iterative attention mining for weakly supervised thoracic disease pattern localization in chest x-rays, in: Frangi, A.F., Schnabel, J.A., Davatzikos, C., Alberola-López, C., Fichtinger, G.

(Eds.), Medical Image Computing and Computer Assisted Intervention – MICCAI 2018, Springer International Publishing, Cham. pp. 589–598.

Cai, L., 2007. Exploiting Known Taxonomies in Learning Overlapping Concepts , 6.

Cesa-bianchi, N., Gentile, C., Tironi, A., Zaniboni, L., 2005. Incremental algorithms for hierarchical classification , 233–240.

Chen, H., Miao, S., Xu, D., Hager, G.D., Harrison, A.P., 2019. Deep hierarchical multi-label classification of chest x-ray images, in: Cardoso, M.J., Feragen, A., Glocker, B., Konukoglu, E., Oguz, I., Unal, G., Vercauteren, T. (Eds.), Proceedings of The 2nd International Conference on Medical Imaging with Deep Learning, PMLR, London, United Kingdom. pp. 109–120.

Dekking, F., Kraaikamp, C., Lopuhaä, H., Meester, L., 2005. A modern introduction to probability and statistics. Understanding why and how. Springer-Verlag London.

DeLong, E.R., DeLong, D.M., Clarke-Pearson, D.L., 1988. Comparing the areas under two or more correlated receiver operating characteristic curves: A nonparametric approach. Biometrics 44, 837–845.

Dembczyński, K., Waegeman, W., Cheng, W., Hüllermeier, E., 2012. On label dependence and loss minimization in multi-label classification. Mach. Learn. 88, 5–45. doi:10.1007/s10994-012-5285-8.

Demner-Fushman, D., Shooshan, S.E., Rodriguez, L., Antani, S., Thoma, G.R., 2015. Annotation of chest radiology reports for indexing and retrieval, in: Müller, H., Jimenez del Toro, O.A., Hanbury, A., Langs, G., Foncubierta Rodriguez, A. (Eds.), Multimodal Retrieval in the Medical Domain, Springer International Publishing, Cham. pp. 99–111.

Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L., 2009. ImageNet: A Large-Scale Hierarchical Image Database, in: CVPR09.

Dimitrovski, I., Kocev, D., Loskovska, S., Deroski, S., 2011. Hierarchical annotation of medical images. Pattern Recogn. 44, 2436–2449. doi:10.1016/j.patcog.2011.03.026.

Dumais, S., Chen, H., 2000. Hierarchical classification of Web content, in: Proceedings of the 23rd annual international ACM SIGIR conference on Research and development in information retrieval - SIGIR '00, ACM Press, Athens, Greece. pp. 256–263. doi:10.1145/345508.345593.

Elkan, C., Noto, K., 2008. Learning classifiers from only positive and unlabeled data, in: Proceeding of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining - KDD 08, ACM Press, Las Vegas, Nevada, USA. p. 213. doi:10.1145/1401890.1401920.

Erdi, C., Ecem, S., Ernst, T.S., Keelin, M., Bram, v.G., 2019. Handling label noise through model confidence and uncertainty: application to chest radiograph classification. doi:10.1117/12.2514290.

Folio, L., 2012. Chest imaging: An algorithmic approach to learning. Springer.

Gohagan, J.K., Prorok, P.C., Hayes, R.B., Kramer, B.S., 2000. The prostate, lung, colorectal and ovarian (plco) cancer screening trial of the national cancer institute: History, organization, and status. Controlled Clinical Trials 21, 251S – 272S. doi:https://doi.org/10.1016/S0197-2456(00)00097-0.

Guan, Q., Huang, Y., Zhong, Z., Zheng, Z., Zheng, L., Yang, Y., 2018. Diagnose like a Radiologist: Attention Guided Convolutional Neural Network for Thorax Disease Classification. arXiv:1801.09927 [cs] ArXiv: 1801.09927.

Gündel, S., Ghesu, F.C., Grbic, S., Gibson, E., Georgescu, B., Maier, A., Comaniciu, D., 2019a. Multi-task Learning for Chest X-ray Abnormality Classification on Noisy Labels. arXiv:1905.06362 [cs] ArXiv: 1905.06362.

Gündel, S., Grbic, S., Georgescu, B., Liu, S., Maier, A., Comaniciu, D., 2019b. Learning to recognize abnormalities in chest x-rays with location-aware dense networks, in: Vera-Rodriguez, R., Fierrez, J., Morales, A. (Eds.), Progress in Pattern Recognition, Image Analysis, Computer Vision, and Applications, Springer International Publishing, Cham. pp. 757–765.

Guo, Y., Liu, Y., Bakker, E.M., Guo, Y., Lew, M.S., 2018. CNN-RNN: a large-scale hierarchical image classification framework. Multimedia Tools and Applications 77, 10251–10271. doi:10.1007/s11042-017-5443-x.

Harvey, H., Glocker, B., 2019. A standardised approach for preparing imaging data for machine learning tasks in radiology, in: Artificial Intelligence in Medical Imaging. Springer, pp. 61–72.

Huang, G., Liu, Z., v. d. Maaten, L., Weinberger, K.Q., 2017. Densely connected convolutional networks, in: 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 2261–2269. doi:10.1109/CVPR.2017.243.

Humphreys, B.L., Lindberg, D.A., 1993. The UMLS project: making the conceptual connection between users and the information they need. Bulletin of the Medical Library Association 81, 170–177.

Irvin, J., Rajpurkar, P., Ko, M., Yu, Y., Ciurea-Ilcus, S., Chute, C., Marklund, H., Haghgoo, B., Ball, R., Shpanskaya, K., Seekins, J., Mong, D., Halabi, S., Sandberg, J., Jones, R., Larson, D., Langlotz, C., Patel, B., Lungren, M., Ng, A., 2019. Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison. Proceedings of the AAAI Conference on Artificial Intelligence 33, 590–597. doi:10.1609/aaai.v33i01.3301590.

Islam, M.T., Aowal, M.A., Minhaz, A.T., Ashraf, K., 2017. Abnormality Detection and Localization in Chest X-Rays using Deep Convolutional Neural Networks. arXiv:1705.09850 [cs] ArXiv: 1705.09850.

Jaeger, S., Karargyris, A., Candemir, S., Siegelman, J., Folio, L., Antani, S., Thoma, G., 2013. Automatic screening for tuberculosis in chest radiographs: a survey. Quantitative Imaging in Medicine and Surgery 3, 89–99. doi:10.3978/j.issn.2223-4292.2013.04.03.

Johnson, A.E.W., Pollard, T.J., Berkowitz, S.J., Greenbaum, N.R., Lungren, M.P., Deng, C.y., Mark, R.G., Horng, S., 2019. MIMIC-CXR, a de-identified publicly

available database of chest radiographs with free-text reports. Scientific Data 6, 1–8. URL: `https://www.nature.com/articles/s41597-019-0322-0`, doi:`10.1038/s41597-019-0322-0`.

Kohli, M.D., Summers, R.M., Geis, J.R., 2017. Medical image data and datasets in the era of machine learning: Whitepaper from the 2016 c-mimi meeting dataset session, in: Journal of Digital Imaging.

Kong, X., Wu, Z., Li, L.J., Zhang, R., Yu, P.S., Wu, H., Fan, W., 2014. Large-Scale Multi-Label Learning with Incomplete Label Assignments. SIAM. pp. 920–928. doi:`10.1137/1.9781611973440.105`.

Kowsari, K., Brown, D.E., Heidarysafa, M., Jafari Meimandi, K., Gerber, M.S., Barnes, L.E., 2017. Hdltex: Hierarchical deep learning for text classification, in: 2017 16th IEEE International Conference on Machine Learning and Applications (ICMLA), pp. 364–371. doi:`10.1109/ICMLA.2017.0-134`.

Langlotz, C.P., 2006. Radlex: A new method for indexing online educational materials. RadioGraphics 26, 1595–1597. doi:`10.1148/rg.266065168`. pMID: 17102038.

Li, Z., Wang, C., Han, M., Xue, Y., Wei, W., Li, L.J., Fei-Fei, L., 2018. Thoracic Disease Identification and Localization with Limited Supervision, in: 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, IEEE, Salt Lake City, UT. pp. 8290–8299. doi:`10.1109/CVPR.2018.00865`.

Liu, B., Dai, Y., Li, X., Lee, W.S., Yu, P.S., 2003. Building text classifiers using positive and unlabeled examples, in: Third IEEE International Conference on Data Mining, pp. 179–186. doi:`10.1109/ICDM.2003.1250918`.

Liu, H., Wang, L., Nan, Y., Jin, F., Wang, Q., Pu, J., 2019. Sdfn: Segmentation-based deep fusion network for thoracic disease classification in chest x-ray images. Computerized Medical Imaging and Graphics 75, 66 – 73. doi:`https://doi.org/10.1016/j.compmedimag.2019.05.005`.

McCallum, A., Rosenfeld, R., Mitchell, T.M., Ng, A.Y., 1998. Improving Text Classification by Shrinkage in a Hierarchy of Classes, in: ICML.

Mettler, F.A., Bhargavan, M., Faulkner, K., Gilley, D.B., Gray, J.E., Ibbott, G.S., Lipoti, J.A., Mahesh, M., McCrohan, J.L., Stabin, M.G., Thomadsen, B.R., Yoshizumi, T.T., 2009. Radiologic and nuclear medicine studies in the united states and worldwide: Frequency, radiation dose, and comparison with other radiation sources—1950–2007. Radiology 253, 520–531. URL: `https://doi.org/10.1148/radiol.2532082010`, doi:`10.1148/radiol.2532082010`, arXiv:`https://doi.org/10.1148/radiol.2532082010`. pMID: 19789227.

Oakden-Rayner, L., 2019. Exploring large scale public medical image datasets. Academic radiology 27, 106–112.

Peng, Y., Wang, X., Lu, L., Bagheri, M., Summers, R., Lu, Z., 2018. NegBio: a high-performance tool for negation and uncertainty detection in radiology reports, in: AMIA 2018 Informatics Summit 2018.

Pourghassem, H., Ghassemian, H., 2008. Content-based medical image classification using a new hierarchical merging scheme. Computerized Medical Imaging and Graphics 32, 651 – 661.

Qi, Z., Yang, M.W., Zhang, Z., Zhang, Z., 2011. Mining partially annotated images, in: KDD, pp. 1199–1207. doi:`10.1145/2020408.2020592`.

Redmon, J., Farhadi, A., 2017. Yolo9000: Better, faster, stronger, in: 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 6517–6525. doi:`10.1109/CVPR.2017.690`.

Roy, D., Panda, P., Roy, K., 2020. Tree-CNN: A hierarchical deep convolutional neural network for incremental learning. Neural Networks 121, 148–160. URL: `https://doi.org/10.1016%2Fj.neunet.2019.09.010`, doi:`10.1016/j.neunet.2019.09.010`.

Stearns, M.Q., Price, C., Spackman, K.A., Wang, A.Y., 2001. SNOMED clinical terms: overview of the development process and project status. Proceedings of the AMIA Symposium , 662–666.

Stevens, R., Aranguren, M.E., Wolstencroft, K., Sattler, U., Drummond, N., Horridge, M., Rector, A., 2007. Using OWL to model biological knowledge. International Journal of Human-Computer Studies 65, 583 – 594. doi:https://doi.org/10.1016/j.ijhcs.2007.03.006.

Vens, C., Struyf, J., Schietgat, L., Džeroski, S., Blockeel, H., 2008. Decision trees for hierarchical multi-label classification. Machine Learning 73, 185. doi:10.1007/s10994-008-5077-3.

Vergara, I.A., Norambuena, T., Ferrada, E., Slater, A.W., Melo, F., 2008. StAR: a simple tool for the statistical comparison of ROC curves. BMC bioinformatics 9, 265–265. doi:10.1186/1471-2105-9-265.

Wang, H., Xia, Y., 2018. ChestNet: A Deep Neural Network for Classification of Thoracic Diseases on Chest Radiography. arXiv:1807.03058 [cs] ArXiv: 1807.03058.

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: 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 3462–3471. doi:10.1109/CVPR.2017.369.

Yan, C., Yao, J., Li, R., Xu, Z., Huang, J., 2018. Weakly Supervised Deep Learning for Thoracic Disease Classification and Localization on Chest X-rays, in: Proceedings of the 2018 ACM International Conference on Bioinformatics, Computational Biology, and Health Informatics, ACM, New York, NY, USA. pp. 103–110. doi:10.1145/3233547.3233573. event-place: Washington, DC, USA.

Yan, Z., Zhang, H., Piramuthu, R., Jagadeesh, V., DeCoste, D., Di, W., Yu, Y., 2015. Hd-cnn: Hierarchical deep convolutional neural networks for large scale visual recognition, in: 2015 IEEE International Conference on Computer Vision (ICCV), pp. 2740–2748. doi:10.1109/ICCV.2015.314.

Yang, S.J., Jiang, Y., Zhou, Z.H., 2013. Multi-Instance Multi-Label Learning with Weak Label , 7.

Yao, L., Poblenz, E., Dagunts, D., Covington, B., Bernard, D., Lyman, K., 2017. Learning to diagnose from scratch by exploiting dependencies among labels. CoRR abs/1710.10501. `arXiv:1710.10501`.

Yu, H.F., Jain, P., Kar, P., Dhillon, I., 2014. Large-scale multi-label learning with missing labels, in: Xing, E.P., Jebara, T. (Eds.), Proceedings of the 31st International Conference on Machine Learning, PMLR, Bejing, China. pp. 593–601.

Zhang, M.L., Zhou, Z.H., 2014. A Review on Multi-Label Learning Algorithms. IEEE Transactions on Knowledge and Data Engineering 26, 1819–1837. `doi:10.1109/TKDE.2013.39`.

Zhao, F., Guo, Y., 2015. Semi-supervised multi-label learning with incomplete labels, in: Proceedings of the 24th International Conference on Artificial Intelligence, AAAI Press. pp. 4062–4068.