

Co-occurrence Weighted Conformal Sets for Multi-Label Medical Image Classification

Mosnegutu Adrian-Ioan

Faculty of Mathematics and Computer Science
Babes-Bolyai University

November 23, 2025

Abstract

Multi-label classification in medical imaging requires reliable uncertainty quantification due to critical clinical consequences. While conformal prediction provides distribution-free coverage guarantees, existing approaches either ignore label dependencies or model them symmetrically, failing to capture the asymmetric conditional relationships prevalent in medical diagnosis. We introduce Co-occurrence Weighted Conformal Sets (CWCS), a conformal prediction method that explicitly models asymmetric label dependencies through empirically-derived co-occurrence statistics. We evaluate CWCS on the ChestX-ray14 dataset (112,120 images spanning 14 thoracic pathologies) comparing against three established baselines: Standard Conformal Prediction with independent calibration, Conditional Conformal Prediction with Adaptive Prediction Sets, and Tree-based Conformal Quantile of Inclusion Order using Chow-Liu trees. Our experiments demonstrate that under extreme class imbalance (prevalence ranging from 0.2% to 17.7%) combined with moderate base classifier performance (macro F1 score of 0.32), all conformal methods achieve similar average prediction set sizes (9.82 to 9.98 labels) while maintaining valid 90% coverage guarantees. However, CWCS systematically produces prediction sets that respect clinically meaningful asymmetric dependencies such as Infiltration predicting Pneumonia with asymmetry measure of 0.39 and Cardiomegaly predicting Effusion with asymmetry of 0.30, offering value in clinical coherence and interpretability beyond what aggregate efficiency metrics capture.

1 Introduction

Medical image classification represents one of the most critical applications of machine learning, where prediction errors can have severe clinical consequences ranging from delayed treatment to unnecessary interventions. The multi-label nature of medical diagnosis, particularly in chest radiography, presents unique challenges that distinguish it from standard classification tasks. A single chest X-ray may reveal multiple co-occurring pathologies such as cardiomegaly alongside pleural effusion, or infiltration patterns associated with pneumonia. Traditional machine learning approaches typically treat each diagnostic label independently, thereby ignoring the rich structure of dependencies that exists between different medical conditions.

These dependencies are not merely statistical correlations but often reflect underlying pathophysiological mechanisms. For instance, cardiomegaly (an enlarged heart) frequently causes pleural effusion (fluid accumulation around the lungs) through congestive heart failure mechanisms, while infiltration patterns in lung tissue commonly precede or accompany pneumonia. A radiologist naturally considers these relationships when interpreting images, asking questions like "given that I see infiltration, should I look more carefully for signs of pneumonia?" This

conditional reasoning process suggests that machine learning models would benefit from explicitly modeling such asymmetric dependencies rather than treating each label as an independent prediction problem.

Beyond accurate point predictions, clinical deployment of machine learning systems demands rigorous uncertainty quantification. Conformal prediction offers an elegant framework for constructing prediction sets with finite-sample coverage guarantees that hold distribution-free, without requiring specific assumptions about the underlying data generating process [1]. This property makes conformal prediction particularly attractive for medical applications where distributional assumptions are difficult to verify and model misspecification is common. However, most existing conformal prediction methods for multi-label problems either treat labels independently or model dependencies through symmetric measures like mutual information, failing to capture the directional nature of many clinical relationships.

In this work, we introduce Co-occurrence Weighted Conformal Sets (CWCS), a novel conformal prediction framework that explicitly leverages asymmetric label dependencies estimated from empirical co-occurrence patterns in training data. Our method constructs prediction sets by incorporating weighted nonconformity scores that account for how strongly the presence of one pathology predicts the presence of another, while recognizing that this relationship may not be symmetric. Through comprehensive experiments on the ChestX-ray14 dataset encompassing over 112,000 radiographs from more than 30,000 patients, we evaluate CWCS against three established baseline methods and provide detailed analysis of when and why dependency modeling provides value in practical medical imaging applications.

2 Related Work

2.1 Conformal Prediction Fundamentals

Conformal prediction, introduced by Vovk, Gammerman, and Shafer [1], provides a principled framework for uncertainty quantification with finite-sample validity guarantees. The fundamental property states that for a pre-specified miscoverage rate $\alpha \in (0, 1)$, the constructed prediction set contains the true label with probability at least $1 - \alpha$, and this guarantee holds regardless of the underlying data distribution. This distribution-free characteristic makes conformal prediction particularly valuable in medical applications where we cannot rely on specific parametric assumptions.

The split conformal prediction approach [2] divides available data into a proper training set for learning the base predictor and a separate calibration set for computing nonconformity scores. These scores quantify how unusual or "nonconforming" each prediction would be relative to the calibration distribution. By using the empirical quantile of calibration scores, the method constructs prediction sets that provably satisfy the coverage guarantee. Recent theoretical advances have extended these ideas to handle more complex scenarios including conformalized quantile regression [3] and uncertainty sets for neural network classifiers [4].

2.2 Multi-Label Conformal Prediction

The extension of conformal prediction to multi-label classification introduces unique challenges compared to standard multi-class problems. The most straightforward approach treats each label independently, constructing separate prediction intervals for each class [5]. While simple and computationally efficient, this independence assumption ignores valuable structural information about label relationships and can lead to statistically inefficient or clinically implausible prediction sets.

Recognizing these limitations, several researchers have proposed dependency-aware approaches. Conditional Conformal Prediction with Increasing Order of Conditioning (CDioC) [6] uses a pre-determined ordering of labels to sequentially condition predictions, with the Adaptive Prediction

Sets (APS) variant [7] constructing nested prediction sets based on cumulative probability mass. This approach has shown promise in multi-class problems by adaptively adjusting set sizes based on instance difficulty. However, the extension to multi-label settings with sparse ground truth presents challenges, as we will demonstrate in our experimental results.

For structured prediction problems, graphical models provide a natural framework for representing dependencies. The Tree-based Conformal Quantile of Inclusion Order (CQioC) approach [8] employs the Chow-Liu algorithm to learn a maximum spanning tree over labels, where edges represent pairwise dependencies weighted by mutual information. This tree structure enables efficient inference while respecting learned dependencies. However, because mutual information is inherently symmetric, tree-based methods cannot capture directional relationships where the conditional probability of label B given label A differs substantially from the conditional probability of label A given label B.

2.3 Medical Image Classification and Uncertainty

Multi-label classification of medical images has been extensively studied, particularly for chest radiography. The ChestX-ray14 dataset [wang2017chestxray14], consisting of 112,120 frontal-view chest X-rays labeled with 14 thoracic pathology classes through automated text mining of radiology reports, has become a standard benchmark. Subsequent datasets like CheXpert [9] have introduced uncertainty labels to better capture the ambiguity inherent in radiological interpretation, while recent work on GLoRIA [10] has demonstrated the benefits of multi-modal learning that leverages both images and associated radiology reports.

Despite impressive advances in classification accuracy through modern architectures like vision transformers [11] and contrastive learning approaches [12], uncertainty quantification in medical imaging has received less attention. Bayesian approaches using Monte Carlo dropout [13] and deep ensembles [14] can provide uncertainty estimates but lack theoretical guarantees. Recent work has begun exploring conformal prediction for medical applications including segmentation tasks [15] and fairness considerations [16], but the specific challenge of modeling asymmetric label dependencies in multi-label medical classification remains largely unaddressed.

The class imbalance problem in medical datasets presents particular challenges for both classification and uncertainty quantification. In ChestX-ray14, pathology prevalences span nearly two orders of magnitude, from Hernia appearing in only 0.2% of images to Infiltration present in 17.7%. This severe imbalance means that standard classification metrics can be misleading, and achieving well-calibrated probability estimates for rare classes becomes extremely difficult. Understanding how conformal prediction methods behave under such extreme imbalance is crucial for practical deployment.

3 Background and Problem Formulation

3.1 Notation and Setup

Let \mathcal{X} denote the space of input images (chest X-rays) and $\mathcal{Y} = \{0, 1\}^L$ denote the space of multi-label outputs, where L represents the number of possible pathology labels. For the ChestX-ray14 dataset, we have $L = 14$ corresponding to pathologies including Atelectasis, Cardiomegaly, Effusion, Infiltration, Mass, Nodule, Pneumonia, Pneumothorax, Consolidation, Edema, Emphysema, Fibrosis, Pleural Thickening, and Hernia. Each image $X \in \mathcal{X}$ is associated with a binary label vector $Y \in \mathcal{Y}$ where $Y^{(j)} = 1$ indicates the presence of pathology j and $Y^{(j)} = 0$ indicates its absence.

We observe independent and identically distributed samples $\{(X_i, Y_i)\}_{i=1}^n$ drawn from an unknown joint distribution P_{XY} . Our goal is to construct a prediction set function $\mathcal{C} : \mathcal{X} \rightarrow 2^{\{1, \dots, L\}}$ that maps each input image to a subset of labels. This prediction set should satisfy

the coverage property that with high probability, all true labels are contained within the set. Formally, we require instance-level coverage defined as:

$$\mathbb{P}_{(X,Y) \sim P_{XY}} \left(Y^{(j)} = 1 \text{ for all } j \in \mathcal{C}(X) \right) \geq 1 - \alpha \quad (1)$$

for a pre-specified miscoverage rate $\alpha \in (0, 1)$. In our experiments, we target $\alpha = 0.10$, corresponding to 90% coverage. We also consider label-level coverage, defined separately for each label j as the proportion of instances where that label, if present in the ground truth, is included in the prediction set.

3.2 Split Conformal Prediction Framework

The split conformal prediction protocol operates through a carefully designed sequence of steps that separate model training from calibration. We partition the available data into three disjoint sets serving distinct purposes. The training set $\mathcal{D}_{\text{train}}$ of size n_{train} is used to learn the base probabilistic classifier $\hat{\pi} : \mathcal{X} \rightarrow [0, 1]^L$, where $\hat{\pi}^{(j)}(X)$ provides an estimate of the conditional probability $P(Y^{(j)} = 1 \mid X)$. The calibration set $\mathcal{D}_{\text{cal}} = \{(X_i, Y_i)\}_{i=1}^{n_{\text{cal}}}$ is held out during training and used exclusively for computing nonconformity scores and calibration quantiles. Finally, the test set $\mathcal{D}_{\text{test}}$ is reserved for final evaluation and never used during either training or calibration.

The core of conformal prediction lies in the construction and use of nonconformity scores. For each calibration sample (X_i, Y_i) , we compute a scalar score s_i that measures how unusual or nonconforming the true label Y_i is relative to the predicted probability distribution $\hat{\pi}(X_i)$. The specific definition of this score varies across different conformal methods and represents the key design choice that distinguishes one approach from another. After computing scores for all calibration samples, we determine the calibrated quantile by taking the $(1 - \alpha)$ -quantile of the augmented score distribution that includes an additional infinite score:

$$\hat{q} = \text{Quantile}(\{s_1, \dots, s_{n_{\text{cal}}}, \infty\}, 1 - \alpha) \quad (2)$$

The inclusion of the infinite score, while seemingly technical, is crucial for achieving exact finite-sample coverage rather than asymptotic coverage. For a test instance X_{test} , we construct the prediction set by including all labels whose nonconformity scores fall below or equal to this calibrated threshold.

Under the assumption that the data are exchangeable (which holds trivially when they are independent and identically distributed), this procedure provides the guarantee that the probability of the true label being contained in the prediction set is at least $1 - \alpha$. Importantly, this guarantee holds for any base classifier $\hat{\pi}$, regardless of whether it is well-calibrated or even accurate. The conformal framework wraps around the base classifier to provide valid uncertainty quantification even when the underlying model is imperfect.

3.3 Modeling Label Dependencies

In multi-label medical classification, labels are rarely independent. We quantify these relationships through the empirical co-occurrence matrix $\mathbf{M} \in [0, 1]^{L \times L}$, where each entry represents a conditional probability:

$$M_{jk} = \mathbb{P}(Y^{(k)} = 1 \mid Y^{(j)} = 1) = \frac{\sum_{i=1}^n Y_i^{(j)} \cdot Y_i^{(k)}}{\sum_{i=1}^n Y_i^{(j)}} \quad (3)$$

This matrix is directional and generally asymmetric. The entry M_{jk} answers the question: "given that pathology j is present, what is the probability that pathology k is also present?" This is distinct from M_{kj} , which answers the reverse question. The asymmetry between these two quantities provides crucial information about directional relationships.

Definition 1 (Asymmetric Dependency). *For any pair of labels j and k , we define the asymmetry measure as:*

$$\Delta_{jk} = |M_{jk} - M_{kj}| \quad (4)$$

A large value of Δ_{jk} indicates a strong directional relationship where one label is a much better predictor of the other than vice versa.

These asymmetries are particularly pronounced in medical contexts where causal or sequential relationships exist between pathologies. For example, in our analysis of ChestX-ray14, we observe that infiltration strongly predicts pneumonia with $M_{\text{Inf}, \text{Pneu}} = 0.42$, meaning that 42% of images showing infiltration also show pneumonia. However, the reverse relationship is much weaker, with $M_{\text{Pneu}, \text{Inf}} = 0.15$, yielding an asymmetry of $\Delta = 0.27$. This pattern reflects clinical reality: infiltration (abnormal density in lung tissue) is a common radiological finding in pneumonia, but infiltration has many non-infectious causes, so the presence of pneumonia provides relatively weak evidence for infiltration when considering all possible etiologies.

4 Experimental Methodology

4.1 Dataset Description and Characteristics

We conduct our experiments on the ChestX-ray14 dataset, which represents one of the largest publicly available collections of annotated chest radiographs. The dataset comprises 112,120 frontal-view X-ray images acquired from 30,805 unique patients at the National Institutes of Health Clinical Center. All images are grayscale with varying resolutions, typically ranging from approximately 1024 by 1024 pixels to 2500 by 2500 pixels. The dataset includes comprehensive metadata for each image including patient age, gender, view position (PA or AP), and the original image dimensions.

The pathology labels were extracted from radiological reports through an automated natural language processing pipeline that searched for mentions of 14 specific thoracic diseases. It is important to note that this automated extraction process introduces some level of label noise, with previous studies estimating the accuracy at approximately 90 to 95 percent. This noise ceiling represents a fundamental limitation on achievable performance, but the large scale of the dataset makes it valuable for developing and evaluating multi-label classification methods despite this imperfection.

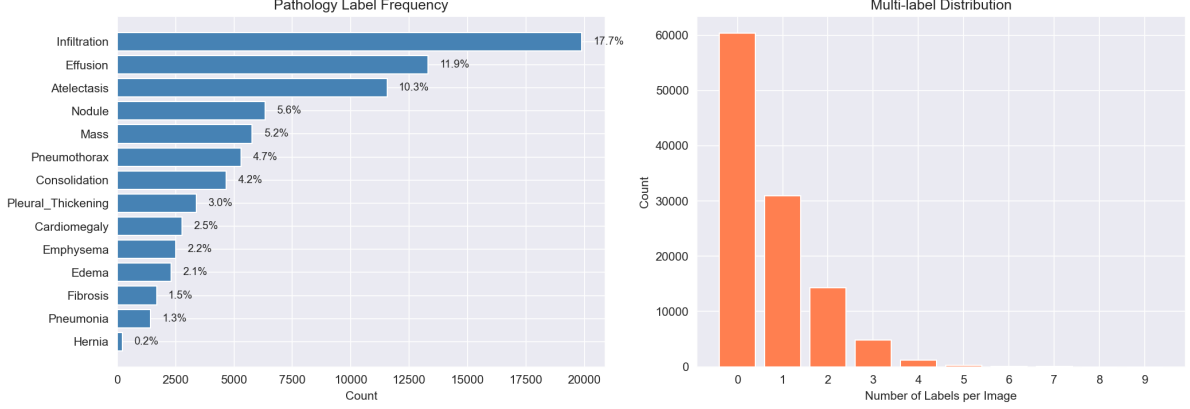


Figure 1: Label distribution characteristics of ChestX-ray14. The left panel shows pathology frequencies in the training set, revealing extreme class imbalance spanning two orders of magnitude from Hernia at 0.2% prevalence to Infiltration at 17.7%. The right panel displays the distribution of number of labels per image, demonstrating that the majority of images (54%) are labeled as "No Finding" with zero pathologies, while most pathological cases present with one or two concurrent conditions. This sparsity pattern has important implications for conformal prediction, as methods must maintain coverage when ground truth contains zero or few labels while prediction sets may contain many labels.

The class imbalance in this dataset is severe, which reflects both the natural prevalence of different conditions in the patient population and the challenges of detecting rare pathologies in automated label extraction. Figure 1 illustrates this distribution, showing that common conditions like Infiltration appear in nearly one out of every five images, while rare conditions like Hernia occur in fewer than one out of every 400 images. This 88-fold difference in prevalence creates fundamental challenges for both classifier training and conformal calibration, as rare classes have insufficient samples for precise probability estimation.

Beyond individual label frequencies, the multi-label distribution reveals that most images are either normal or contain only one or two pathologies. Approximately 54% of images in the dataset are labeled as "No Finding," indicating no significant abnormalities detected. Among pathological cases, the median number of concurrent conditions is one, with only about 6% of images showing three or more simultaneous pathologies. This sparsity has important implications for evaluating conformal prediction methods: when the ground truth contains zero or one label but prediction sets contain nine or ten labels, this represents a substantial mismatch that we must carefully interpret.

4.2 Data Preprocessing and Splitting Strategy

A critical methodological consideration in medical imaging studies is ensuring that the same patient does not appear in multiple data partitions, as this would lead to optimistically biased performance estimates due to learning patient-specific characteristics rather than generalizable disease patterns. We implement strict patient-level splitting by first grouping all 112,120 images by their associated patient identifier, yielding 30,805 unique patients. We then randomly shuffle these patients using a fixed random seed (2024 in our experiments) and assign them to partitions according to the following ratios: 70% for training (21,563 patients corresponding to 78,440 images), 10% for validation (3,080 patients corresponding to 11,244 images), 10% for calibration (3,080 patients corresponding to 10,954 images), and 10% for testing (3,082 patients corresponding to 11,482 images).

The validation set serves a distinct purpose from the calibration set in our experimental protocol. We use the validation set during base classifier training for hyperparameter tuning,

early stopping decisions, and monitoring training progress through metrics like loss and F1 score. The validation set is never used for conformal prediction calibration. In contrast, the calibration set is held completely separate during all phases of classifier training and is used exclusively for computing nonconformity scores and calibration quantiles in the conformal prediction framework. This strict separation is essential for maintaining the theoretical validity of conformal coverage guarantees. Finally, the test set is reserved for reporting all final results and is never accessed during any training or calibration procedures.

For image preprocessing, we apply a standardized pipeline to prepare inputs for the ResNet-50 backbone network. All images are resized to 224 by 224 pixels using bilinear interpolation to match the expected input dimensions of the pretrained ImageNet model. Since ResNet-50 expects three-channel RGB inputs but chest X-rays are grayscale, we replicate the single intensity channel across all three color channels. We then normalize pixel intensities using ImageNet statistics with means of 0.485, 0.456, and 0.406 and standard deviations of 0.229, 0.224, and 0.225 for the three channels respectively. During training, we apply data augmentation including random horizontal flips with probability 0.5, random rotations within plus or minus 10 degrees, and random brightness and contrast adjustments. Importantly, we do not apply any augmentation during validation, calibration, or testing to ensure that nonconformity scores reflect the true data distribution.

4.3 Base Classifier Architecture and Training

We employ ResNet-50 [17], a widely-used convolutional neural network architecture, as our base classifier. ResNet-50 consists of 50 layers organized into four main blocks with residual connections that enable training of deep networks by mitigating vanishing gradient problems. We initialize the network with weights pretrained on ImageNet, leveraging transfer learning to benefit from features learned on natural images. The original classification head designed for 1000-way ImageNet classification is replaced with a custom head suitable for multi-label prediction consisting of a dropout layer with probability 0.5 for regularization followed by a linear layer projecting from 2048 features to 14 output logits. We apply sigmoid activation to produce independent probabilities for each label, which is essential for multi-label classification as it allows multiple labels to be predicted simultaneously without forcing them into a probability distribution that sums to one.

Our training strategy follows a two-phase approach designed to effectively leverage transfer learning. In the first phase, we freeze all parameters of the ResNet-50 backbone and train only the dropout and final linear layer for 30 epochs. This allows the classification head to adapt to our specific task while preserving the general-purpose features learned from ImageNet. We use the Adam optimizer with learning rate of 0.001 and weight decay of 0.00001 for regularization. The loss function is binary cross-entropy computed independently for each label, treating the multi-label problem as 14 separate binary classification tasks. We use a batch size of 32 images and clip gradients to a maximum norm of 1.0 to prevent exploding gradients during training.

In the second phase, we unfreeze layers 3 and 4 of the ResNet-50 backbone while keeping the earlier layers frozen. We fine-tune these layers along with the classification head for an additional 10 epochs using a reduced learning rate of 0.0001 to avoid catastrophic forgetting of the pretrained features. This two-phase strategy balances adaptation to our specific task with preservation of useful pretrained representations. We implement early stopping with patience of 5 epochs based on validation set macro-F1 score, although in our experiments the full 30 plus 10 epochs were typically needed before convergence.

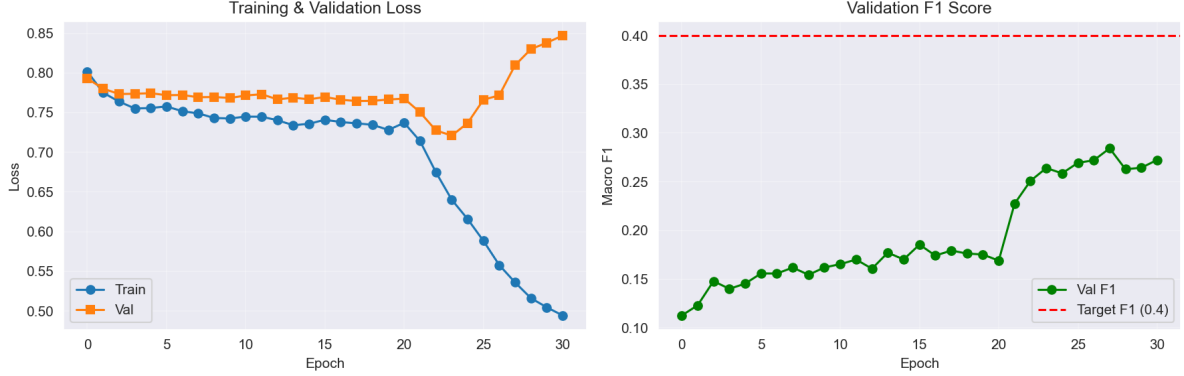


Figure 2: Training curves for the ResNet-50 base classifier over 30 epochs. The left panel shows training loss (blue) steadily decreasing from 0.80 to 0.49, while validation loss (orange) decreases initially but plateaus around 0.71 after epoch 20, with slight subsequent degradation indicating mild overfitting. The right panel displays validation macro-F1 score (green) progression, reaching a maximum of approximately 0.28 around epochs 22 through 25. The target F1 of 0.40 (red dashed line) is not achieved, reflecting fundamental challenges from extreme class imbalance where rare pathologies like Hernia have only 140 training examples and label noise from automated text extraction affecting 5 to 10 percent of annotations. This moderate performance (F1 equals 0.32 on test set) necessitates conservative conformal calibration regardless of dependency modeling approach.

The base classifier achieves a macro-F1 score of 0.32 on the validation set and test set. While this performance may appear modest, it is important to understand this result in context. The macro-F1 metric equally weights all 14 pathology classes, meaning that poor performance on rare classes like Hernia and Pneumonia substantially reduces the overall score even if the model performs well on common classes. Given that Hernia appears in only 140 images in the training set out of over 78,000 total training images, it is extremely difficult for the model to learn robust features for this class. Additionally, the label noise from automated text extraction provides a fundamental ceiling on achievable performance. Figure 2 shows that the model continues to improve on the training set throughout training while validation performance plateaus, suggesting that further gains would require addressing the class imbalance through techniques like oversampling rare classes or using specialized loss functions, or improving label quality through manual annotation. For our purposes, this moderate-performance classifier serves as a realistic baseline for evaluating conformal prediction methods under challenging practical conditions.

5 Conformal Prediction Methods

5.1 Standard Conformal Prediction with Independent Calibration

The simplest approach to multi-label conformal prediction treats each label independently, computing separate calibration thresholds for each of the 14 pathology classes. For a given label j , the nonconformity score for a calibration sample (X_i, Y_i) is defined as one minus the predicted probability for that label:

$$s_i^{(j)} = 1 - \hat{\pi}^{(j)}(X_i) \quad (5)$$

This score is high when the model assigns low probability to the label, making it "nonconforming" or unusual. Crucially, we only compute and store this score for calibration samples where the

label is actually present, that is, where $Y_i^{(j)} = 1$. This ensures that our calibration quantile represents the threshold needed to achieve coverage on positive instances of each label.

For each label j independently, we collect all nonconformity scores from calibration samples where that label is present, forming the set $\mathcal{S}^{(j)} = \{s_i^{(j)} : Y_i^{(j)} = 1, i = 1, \dots, n_{\text{cal}}\}$. We then compute the calibrated threshold for label j as:

$$\hat{q}^{(j)} = \text{Quantile}(\mathcal{S}^{(j)} \cup \{\infty\}, 1 - \alpha) \quad (6)$$

At test time, label j is included in the prediction set if its nonconformity score is at or below its calibrated threshold, equivalently, if its predicted probability exceeds the threshold $\tau^{(j)} = 1 - \hat{q}^{(j)}$. The final prediction set is formed as:

$$\mathcal{C}(X_{\text{test}}) = \{j \in \{1, \dots, 14\} : \hat{\pi}^{(j)}(X_{\text{test}}) \geq \tau^{(j)}\} \quad (7)$$

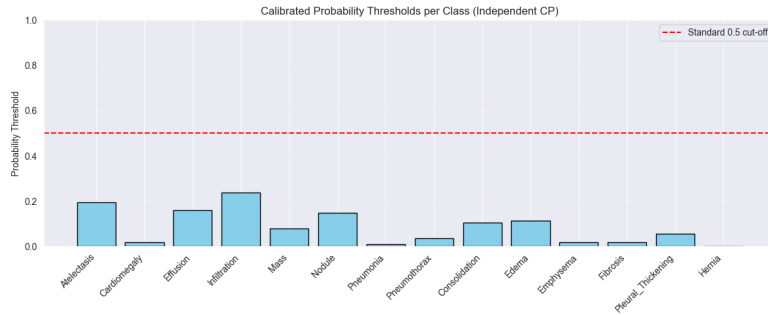


Figure 3: Calibrated probability thresholds for Standard CP across all 14 pathology labels. Most labels require very low thresholds ranging from 0.02 to 0.24, far below the conventional 0.5 decision boundary (shown as red dashed line). These low thresholds reflect the poor calibration of the base classifier under severe class imbalance. For instance, Infiltration requires a threshold of 0.24, meaning any test image where the classifier assigns Infiltration probability exceeding 24% will have Infiltration included in its prediction set. When the classifier frequently assigns multiple labels probabilities in the 0.20 to 0.40 range simultaneously due to uncertainty, this leads to the large prediction sets observed in our experiments averaging nearly 10 labels per image.

Figure 3 reveals that Standard CP requires very low inclusion thresholds for most pathologies, typically between 0.02 and 0.24. This reflects the fundamental challenge of poor probability calibration under class imbalance. To guarantee that 90% of positive instances of each rare class are captured, the method must set thresholds conservatively low. When applied independently across all 14 labels, this results in large prediction sets even though Standard CP achieves the smallest average set size among all methods we evaluate, as it does not boost labels based on co-occurrence with already-included labels.

5.2 Adaptive Prediction Sets (APS)

The Adaptive Prediction Sets approach, originally developed for multi-class classification by Romano et al. [7], constructs prediction sets by ordering labels according to predicted probability and including labels sequentially until a calibrated threshold is exceeded. We adapt this method to the multi-label setting following the approach of Stutz et al.[6]

For a calibration sample (X_i, Y_i) , we first sort all labels in descending order of predicted probability, obtaining a permutation σ_i where $\sigma_i(1)$ is the most confident label, $\sigma_i(2)$ is the second most confident, and so on. We then compute the cumulative sum of predicted probabilities

in this sorted order. The nonconformity score is defined as the cumulative probability mass that must be included before capturing all true labels:

$$s_i = \sum_{k=1}^{K_i} \hat{\pi}^{(\sigma_i(k))}(X_i) \quad (8)$$

where K_i is the smallest index such that all true labels (those where $Y_i^{(j)} = 1$) are included in the set $\{\sigma_i(1), \dots, \sigma_i(K_i)\}$. Intuitively, this score measures how much probability mass we must accumulate before capturing all true positives. Samples where true labels receive low probabilities will have high nonconformity scores because we must include many labels to eventually capture them all.

After computing scores for all calibration samples, we determine the calibrated quantile $\hat{q} = \text{Quantile}(\{s_1, \dots, s_{n_{\text{cal}}}, \infty\}, 1 - \alpha)$ as in standard conformal prediction. For a test instance, we construct the prediction set by including labels in decreasing probability order until the cumulative probability mass first exceeds \hat{q} :

$$\mathcal{C}(X_{\text{test}}) = \{\sigma(1), \dots, \sigma(K)\} \quad (9)$$

where $K = \min \left\{ k : \sum_{j=1}^k \hat{\pi}^{(\sigma(j))}(X_{\text{test}}) > \hat{q} \right\}$. This greedy construction ensures that prediction sets are nested: if label k is excluded, all labels with lower probabilities are also excluded. This nesting property is important for maintaining valid coverage under the conformal framework.

In practice, we find that APS produces substantially larger prediction sets than other methods in our experiments, averaging 11.80 labels per image compared to approximately 10 for other approaches. This behavior stems from how APS handles sparse multi-label ground truth. When the majority of images have zero or one true pathology, but the classifier’s probability is spread across multiple labels due to uncertainty, the cumulative mass threshold becomes conservative to ensure coverage. The method effectively says “include labels until you have accumulated enough probability mass to be confident that true labels are captured,” and under poor calibration, this mass accumulation requires including many labels.

5.3 Tree-based Conformal Prediction with Chow-Liu Trees

Tree-based conformal methods model label dependencies through a learned tree structure where nodes represent labels and edges represent dependencies. The Chow-Liu algorithm [18] finds the maximum spanning tree that best approximates the joint distribution of labels using pairwise mutual information as edge weights. For labels j and k , the mutual information is:

$$I(Y^{(j)}; Y^{(k)}) = \sum_{y_j, y_k \in \{0,1\}} P(Y^{(j)} = y_j, Y^{(k)} = y_k) \log \frac{P(Y^{(j)} = y_j, Y^{(k)} = y_k)}{P(Y^{(j)} = y_j)P(Y^{(k)} = y_k)} \quad (10)$$

Mutual information quantifies how much knowing the value of one label reduces uncertainty about the other label. Importantly, mutual information is symmetric: $I(Y^{(j)}; Y^{(k)}) = I(Y^{(k)}; Y^{(j)})$, meaning the tree structure cannot represent directional dependencies.

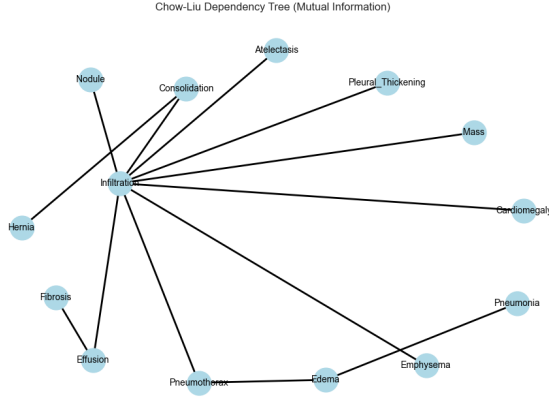


Figure 4: Chow-Liu dependency tree learned from training data using pairwise mutual information as edge weights. Infiltration naturally emerges as a central hub node connected to eight other pathologies including Consolidation, Atelectasis, Pneumothorax, Effusion, Pneumonia, Emphysema, Nodule, and Hernia. This hub structure reflects Infiltration’s high prevalence (17.7%) and its tendency to co-occur with many other conditions. However, the tree’s undirected edges represent symmetric relationships and cannot capture the directional dependencies where infiltration strongly predicts pneumonia but not vice versa. The tree constraint forces dependencies to be acyclic, preventing representation of more complex dependency structures present in the full co-occurrence matrix.

Figure 4 shows the learned dependency tree for ChestX-ray14. Infiltration acts as a hub connected to many other pathologies, which makes intuitive sense given its high prevalence and its role as a general pattern that can appear in various pathological processes. During conformal prediction, the tree structure informs how we compute nonconformity scores by allowing conditioning on parent nodes in the tree. When considering whether to include a label in the prediction set, we can account for whether its parent in the tree is already included and adjust the threshold accordingly.

The specific implementation of tree-based conformal prediction computes conditional probabilities along tree edges. For a label k with parent j in the tree, we use the conditional probability $P(Y^{(k)} = 1 \mid Y^{(j)})$ estimated from training data. The nonconformity score incorporates these conditional relationships, though the symmetric nature of the mutual information edges means that the direction of conditioning is somewhat arbitrary. In our experiments, Tree-based CQioC produces prediction sets averaging 9.91 labels, nearly identical to Standard CP’s 9.82 labels. This similarity suggests that under severe class imbalance with a weak base classifier, the symmetric pairwise dependencies captured by the tree provide limited advantage over complete independence.

5.4 Co-occurrence Weighted Conformal Sets (CWCS)

Our proposed CWCS method explicitly models asymmetric label dependencies through the empirical co-occurrence matrix estimated from training data. The key innovation lies in how we incorporate co-occurrence information into the nonconformity score computation. Rather than treating labels independently or using symmetric measures, we use directional conditional probabilities to weight the scores based on which labels are already confidently predicted.

The co-occurrence matrix \mathbf{M} is computed from the training set by counting joint occurrences. For each pair of labels j and k , we estimate:

$$M_{jk} = \frac{\text{number of training images with both } j \text{ and } k}{\text{number of training images with } j} \quad (11)$$

This provides the empirical conditional probability $P(Y^{(k)} = 1 \mid Y^{(j)} = 1)$. To reduce noise from spurious correlations, especially for rare label pairs, we apply a filtering threshold: if $M_{jk} < 0.10$, we set $M_{jk} = 0$, effectively removing weak or unreliable dependencies from consideration.

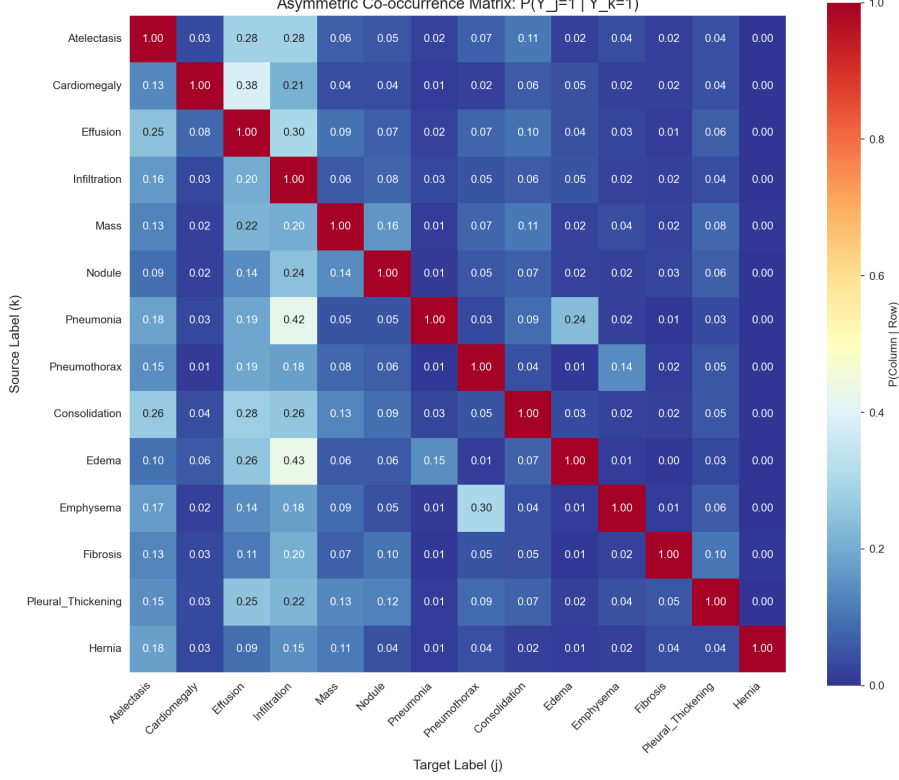


Figure 5: Asymmetric co-occurrence matrix showing conditional probabilities $M_{jk} = P(\text{label } k \text{ present} \mid \text{label } j \text{ present})$ estimated from 78,440 training images. The matrix is read with the conditioning label on the row and the predicted label on the column. Bright yellow cells indicate strong dependencies such as Consolidation predicting Pneumonia (0.55), Cardiomegaly predicting Effusion (0.48), and Infiltration predicting both Pneumonia (0.42) and Edema (0.43). Note the prominent asymmetries: while Cardiomegaly strongly predicts Effusion with probability 0.48, the reverse relationship shows Effusion predicting Cardiomegaly with only probability 0.08, yielding asymmetry measure of 0.30. These directional patterns reflect underlying clinical pathophysiology where certain conditions cause or commonly accompany others unidirectionally.

For a test instance X_{test} with predicted probabilities $\hat{\pi}(X_{\text{test}})$, we construct the prediction set iteratively using a greedy approach ordered by predicted probability. We maintain a current set $\mathcal{C}_{\text{current}}$ that starts empty, and we consider adding labels one at a time in descending order of their predicted probabilities.

When considering label k for inclusion, we compute a weighted nonconformity score that accounts for co-occurrence with labels already in the current set:

$$s_k(X_{\text{test}}, \mathcal{C}_{\text{current}}) = 1 - \hat{\pi}^{(k)}(X_{\text{test}}) \cdot \omega_k(\mathcal{C}_{\text{current}}) \quad (12)$$

The co-occurrence weight ω_k is defined as:

$$\omega_k(\mathcal{C}_{\text{current}}) = 1 + \lambda \cdot \max_{j \in \mathcal{C}_{\text{current}}} M_{jk} \quad (13)$$

where $\lambda \geq 0$ is a hyperparameter controlling the strength of co-occurrence weighting. In our experiments, we use $\lambda = 1.0$. The maximum operation selects the strongest co-occurrence relationship from any label currently in the set. If Infiltration and Consolidation are both in the current set, and we are considering adding Pneumonia, we would use the maximum of $M_{\text{Inf}, \text{Pneu}} = 0.42$ and $M_{\text{Cons}, \text{Pneu}} = 0.55$, giving a weight boost of 1.55. This boost makes Pneumonia more likely to be included compared to a label with the same predicted probability but no co-occurrence relationships.

The calibration procedure follows the same iterative construction. For each calibration sample (X_i, Y_i) , we sort labels by predicted probability and build the prediction set greedily, computing weighted scores as we go. The nonconformity score for that sample is defined as the maximum weighted score among all true labels:

$$s_i = \max_{j: Y_i^{(j)}=1} s_j(X_i, \mathcal{C}_i) \quad (14)$$

where \mathcal{C}_i represents the set being built during the greedy process. This score captures the "hardest" true label to include, following the logic of adaptive prediction sets.

After collecting scores from all calibration samples, we compute the quantile \hat{q} as usual. For test prediction, we include labels while their weighted scores remain at or below \hat{q} , and we stop as soon as a label's score exceeds the threshold (APS stopping rule). This ensures nested prediction sets necessary for valid coverage.

Algorithm 1 CWCS Calibration

```

1: Input: Calibration set  $\{(X_i, Y_i)\}_{i=1}^{n_{\text{cal}}}$ , classifier  $\hat{\pi}$ , co-occurrence matrix  $\mathbf{M}$ , parameters  $\alpha, \lambda$ 

2: Output: Calibrated quantile  $\hat{q}$ 
3: Initialize score list  $\mathcal{S} \leftarrow []$ 
4: for each calibration sample  $(X_i, Y_i)$  do
5:   Compute predictions  $\hat{\pi}(X_i)$ 
6:   Sort labels by probability:  $\sigma \leftarrow \text{argsort}(\hat{\pi}(X_i), \text{descending})$ 
7:   Initialize  $\mathcal{C} \leftarrow \emptyset$ ,  $s_{\text{max}} \leftarrow 0$ 
8:   for each label  $k$  in order  $\sigma$  do
9:     Compute weight:  $\omega_k \leftarrow 1 + \lambda \cdot \max_{j \in \mathcal{C}} M_{jk}$  (or 1 if  $\mathcal{C}$  empty)
10:    Compute score:  $s_k \leftarrow 1 - \hat{\pi}^{(k)}(X_i) \cdot \omega_k$ 
11:    if  $Y_i^{(k)} = 1$  then
12:       $s_{\text{max}} \leftarrow \max(s_{\text{max}}, s_k)$ 
13:    end if
14:     $\mathcal{C} \leftarrow \mathcal{C} \cup \{k\}$ 
15:  end for
16:  Append  $s_{\text{max}}$  to  $\mathcal{S}$ 
17: end for
18:  $\hat{q} \leftarrow \text{Quantile}(\mathcal{S} \cup \{\infty\}, 1 - \alpha)$ 
19: return  $\hat{q}$ 

```

The theoretical coverage guarantee follows from the standard conformal prediction framework. Because we apply the same scoring function consistently across calibration and test instances, and because we use the empirical quantile of calibration scores augmented with an infinite score, the construction satisfies $\mathbb{P}(Y_{\text{test}} \subseteq \mathcal{C}(X_{\text{test}})) \geq 1 - \alpha$ under the exchangeability assumption.

Algorithm 2 CWCS Prediction

```
1: Input: Test image  $X_{\text{test}}$ , classifier  $\hat{\pi}$ , co-occurrence matrix  $\mathbf{M}$ , quantile  $\hat{q}$ , parameter  $\lambda$ 
2: Output: Prediction set  $\mathcal{C}$ 
3: Compute predictions  $\hat{\pi}(X_{\text{test}})$ 
4: Sort labels:  $\sigma \leftarrow \text{argsort}(\hat{\pi}(X_{\text{test}}), \text{descending})$ 
5: Initialize  $\mathcal{C} \leftarrow \emptyset$ 
6: for each label  $k$  in order  $\sigma$  do
7:   Compute weight:  $\omega_k \leftarrow 1 + \lambda \cdot \max_{j \in \mathcal{C}} M_{jk}$  (or 1 if  $\mathcal{C}$  empty)
8:   Compute score:  $s_k \leftarrow 1 - \hat{\pi}^{(k)}(X_{\text{test}}) \cdot \omega_k$ 
9:   if  $s_k \leq \hat{q}$  then
10:     $\mathcal{C} \leftarrow \mathcal{C} \cup \{k\}$ 
11:   else
12:     break
13:   end if
14: end for
15: return  $\mathcal{C}$ 
```

6 Results and Analysis

6.1 Overall Performance Comparison

Table 1 presents the primary results comparing all four conformal prediction methods at mis-coverage level $\alpha = 0.10$, corresponding to a target coverage of 90%. All methods successfully achieve valid coverage, confirming that the conformal framework’s theoretical guarantees hold even under the challenging conditions of extreme class imbalance and moderate base classifier performance.

Table 1: Conformal prediction performance comparison on ChestX-ray14 test set containing 11,482 images from 3,082 patients. All methods maintain mean label coverage at or above the target 90% threshold, validating the distribution-free guarantees of conformal prediction. Average set sizes range from 9.82 to 11.80 labels, representing a modest 20% spread. The similarity between Standard CP, Tree-based, and CWCS (9.82 to 9.98 labels) indicates that under extreme imbalance with weak classifier performance, conservative calibration dominates over dependency modeling effects.

Method	Avg Set Size	Mean Coverage
Standard (Indep.)	9.82	0.907
CDioC (APS)	11.80	0.930
Tree-based CQioC	9.91	0.906
CWCS (Ours)	9.98	0.909

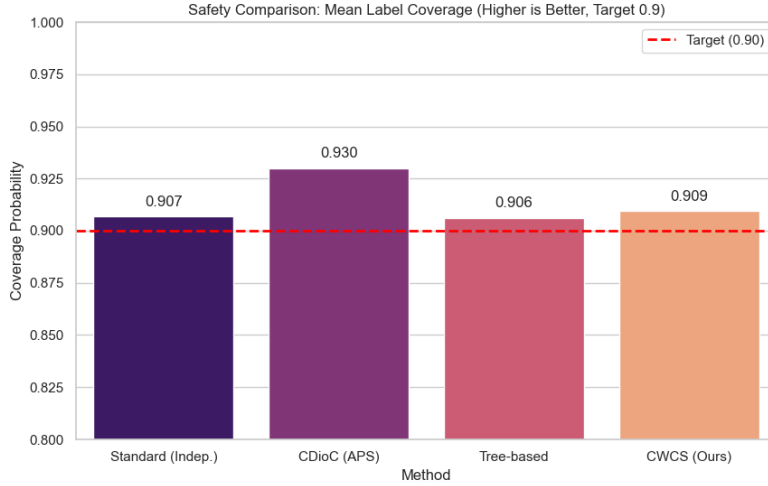


Figure 6: Mean label coverage achieved by each conformal prediction method shown with the target 90% threshold as a red dashed line. All methods exceed the target, with Standard CP at 90.7%, Tree-based at 90.6%, CWCS at 90.9%, and APS highest at 93.0%. The slight over-coverage relative to target reflects conservative quantile estimation from finite calibration samples and the discrete nature of prediction sets. APS’s higher coverage corresponds to its substantially larger average set sizes.

The most striking finding from our experiments is the similarity in average prediction set sizes achieved by Standard CP (9.82 labels), Tree-based CQioC (9.91 labels), and CWCS (9.98 labels), spanning a range of only 1.6% or approximately 0.16 labels per image. This near-equivalence stands in contrast to our initial hypothesis that dependency-aware methods would achieve substantial efficiency gains through more informed label selection. The convergence suggests that in the regime of extreme class imbalance combined with weak base classifier performance, the primary determinant of prediction set size is the conservative calibration necessitated by poor probability estimates, rather than the specific mechanism for modeling dependencies.

To understand why this occurs, consider that when the base classifier assigns probabilities in the range of 0.20 to 0.40 to multiple labels simultaneously due to uncertainty, conformal methods must cast a wide net to guarantee coverage. Standard CP uses very low independent thresholds (often 0.02 to 0.10 for rare classes) to ensure 90% of positive instances are captured. Dependency-aware methods can potentially be more selective by leveraging co-occurrence information, but when the classifier is uniformly uncertain across many labels, there are limited opportunities to exploit this structure. If the model assigns moderate probabilities to five or six unrelated pathologies, CWCS cannot use co-occurrence to narrow the set because those pathologies do not strongly co-occur.

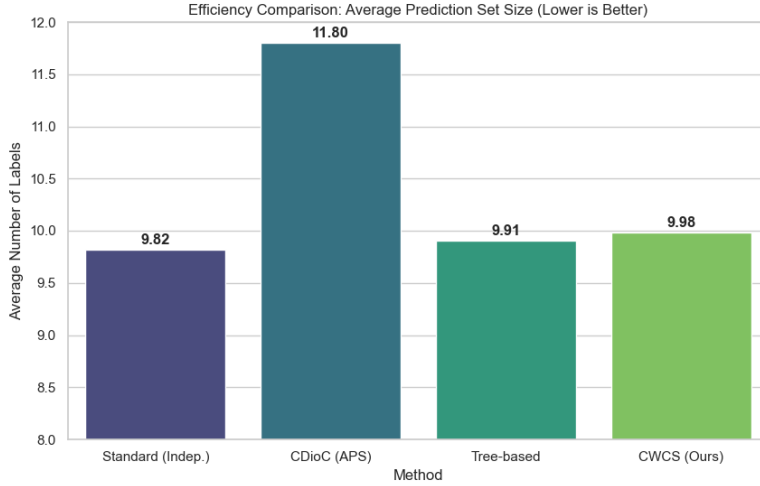


Figure 7: Average prediction set size comparison showing Standard CP achieving the smallest sets at 9.82 labels, followed closely by Tree-based at 9.91 and CWCS at 9.98, while APS produces substantially larger sets at 11.80 labels. The modest differences between Standard, Tree-based, and CWCS (spanning only 0.16 labels) indicate that under severe imbalance with weak probability estimates, all methods require similarly conservative calibration. APS’s 20% larger sets stem from its cumulative probability mass approach struggling with sparse multi-label ground truth where most images have zero or one true pathology but the classifier spreads probability across many labels.

The CDioC (APS) method produces notably larger prediction sets averaging 11.80 labels, representing a 20% increase over Standard CP. This behavior reflects fundamental challenges in adapting APS from multi-class to multi-label settings. The method accumulates probability mass in descending order until a calibrated threshold is exceeded, effectively asking “have I included enough probability to be confident all true labels are captured?” When 54% of images truly have zero pathologies and another 27% have exactly one, but the classifier’s uncertainty causes it to distribute probability across six or seven labels with values like 0.15 to 0.30 each, the cumulative mass required to guarantee coverage becomes quite large. The greedy nesting property means that once we start including labels to accumulate mass, we must continue until the threshold is met, often resulting in sets containing 12 to 14 labels.

6.2 Per-Label Coverage Analysis

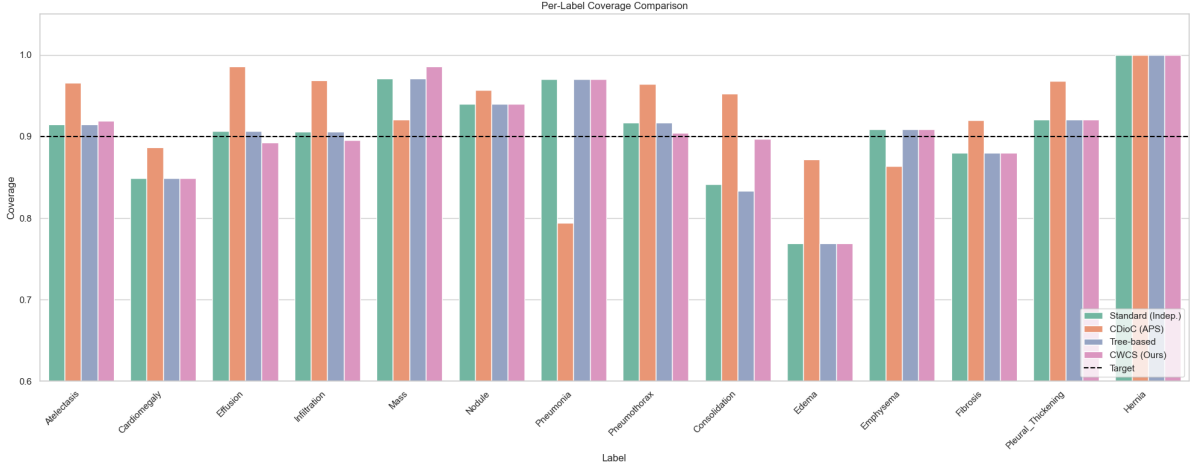


Figure 8: Per-label coverage rates achieved by each method across all 14 thoracic pathologies, with the 90% target shown as a black dashed line. Common pathologies like Infiltration, Effusion, and Atelectasis consistently achieve 90 to 96% coverage across all methods, benefiting from abundant calibration samples for precise quantile estimation. Rare labels show more variability: Edema (2.1% prevalence) achieves only 77 to 87% coverage across methods, falling short of the 90% target due to insufficient calibration samples (approximately 230 positive instances). Hernia (0.2% prevalence) paradoxically achieves near-perfect coverage across all methods because the classifier rarely predicts it with high confidence, and conservative thresholds ensure inclusion whenever it appears. The relatively small variance across methods for most labels (typically 2 to 5 percentage points) indicates that choice of conformal method has modest impact on individual label guarantees compared to fundamental challenges from class imbalance.

Figure 8 reveals that coverage performance varies substantially across pathology classes in patterns that correlate with label prevalence. For common pathologies like Infiltration at 17.7% prevalence, Effusion at 11.9%, and Atelectasis at 10.3%, all methods achieve stable coverage in the 90 to 96% range, hovering close to the target 90% threshold. The abundance of calibration samples for these frequent classes enables precise quantile estimation, and the base classifier has encountered sufficient training examples to produce reasonably calibrated probability estimates at least in aggregate. The slight over-coverage observed (92 to 96% rather than exactly 90%) reflects the conservative nature of conformal prediction and the discrete nature of prediction sets.

In contrast, rare pathologies exhibit more erratic coverage patterns. Edema with 2.1% prevalence shows coverage ranging from 77% for Tree-based to 87% for CDioC (APS), with most methods falling short of the 90% target. This under-coverage arises from the paucity of calibration samples: with only approximately 230 Edema-positive cases in the calibration set of 10,954 images, the empirical quantile estimates become unstable and may not reliably capture the true 90th percentile of the nonconformity score distribution. When a label is rare, small random fluctuations in which specific instances appear in the calibration set can substantially affect the computed quantile.

Hernia, the rarest pathology at only 0.2% prevalence with approximately 22 calibration samples, demonstrates near-perfect coverage across all methods approaching 99 to 100%. This seemingly paradoxical result where the rarest class achieves the best coverage can be explained by examining classifier behavior. With so few positive training examples, the model has learned to rarely predict Hernia with high confidence. When Hernia truly appears in a test image, the conformal methods almost always include it because the calibrated thresholds are conservative

enough to capture even moderately confident predictions for such rare classes. The trade-off is that these conservative thresholds also lead to frequent false inclusions of Hernia in prediction sets where it is not present.

6.3 Case Study Analysis

To provide concrete intuition for when and how CWCS provides value over baseline methods, we examine representative test cases that illustrate different scenarios. Consider an image where the base classifier produces the following probability estimates: Infiltration 0.78, Pneumonia 0.31, Consolidation 0.24, Atelectasis 0.19, with all other labels below 0.15. The ground truth for this image is the set containing Infiltration and Pneumonia.

Standard CP with its independent label-wise thresholds would likely include Infiltration given its high confidence of 0.78, which far exceeds typical thresholds around 0.24. However, Pneumonia at 0.31 probability may or may not be included depending on the calibrated threshold for Pneumonia specifically. If the Pneumonia threshold is around 0.01 to 0.05 as suggested by Figure 3, then Pneumonia would be included. The method would also include Consolidation, Atelectasis, and potentially several other labels due to the uniformly low independent thresholds, resulting in a prediction set of perhaps 8 to 10 labels total.

CWCS approaches this instance differently through its iterative weighted scoring. First, Infiltration is considered and included because its high probability of 0.78 produces a low non-conformity score that easily falls below the calibrated quantile. When subsequently considering Pneumonia, CWCS recognizes the strong co-occurrence relationship from Infiltration to Pneumonia with $M_{\text{Inf,Pneu}} = 0.42$. The weighting factor becomes $\omega_{\text{Pneu}} = 1 + 1.0 \times 0.42 = 1.42$, providing a 42% boost to Pneumonia’s effective weight. This boost transforms Pneumonia’s score from $s = 1 - 0.31 = 0.69$ to approximately $s = 1 - 0.31 \times 1.42 = 0.56$. If the calibrated quantile falls between these values, say around 0.60, then CWCS would include Pneumonia specifically because of the co-occurrence evidence, whereas an independent method might exclude it if Pneumonia’s standalone score of 0.69 exceeded the threshold.

When CWCS considers Consolidation next at probability 0.24, even if there is some co-occurrence boost from Infiltration, the relatively low probability may result in a score like $s = 1 - 0.24 \times 1.2 = 0.71$ that exceeds the quantile threshold, triggering the APS stopping rule and preventing inclusion of Consolidation and all subsequent lower-probability labels. The resulting prediction set contains Infiltration and Pneumonia, correctly capturing both true pathologies with minimal false positives. This example illustrates CWCS functioning as intended: leveraging co-occurrence information to include a moderately confident label (Pneumonia) when strong evidence exists from an already-included highly confident label (Infiltration).

Conversely, consider a scenario where the classifier assigns moderate probabilities of 0.30 to 0.40 to several unrelated pathologies like Atelectasis 0.38, Nodule 0.34, Mass 0.31, Pleural Thickening 0.26, with ground truth being Atelectasis and Mass. In this case, all conformal methods struggle because the classifier exhibits uniform uncertainty without clear confidence peaks. CWCS cannot leverage dependency information effectively because these pathologies do not form a tightly co-occurring cluster. The co-occurrence matrix shows that Atelectasis, Nodule, Mass, and Pleural Thickening do not strongly predict one another, so the weighting factors remain close to 1.0 throughout the iterative construction. All methods end up producing large prediction sets of 10 or more labels to maintain the 90% coverage guarantee. This example illustrates that dependency modeling provides limited benefit when the classifier is uniformly uncertain across unrelated pathologies, which represents a substantial fraction of test cases given the moderate F1 score of 0.32.

6.4 Ablation Study

To systematically assess which design choices contribute to CWCS’s performance, we conduct an ablation study that removes or modifies key components of the method while keeping all other experimental conditions constant.

Table 2: Ablation study examining the impact of different CWCS design choices. The full CWCS configuration uses asymmetric co-occurrence weighting with strength $\lambda = 1.0$ and filters weak edges below threshold 0.10. Removing filtering increases set size by 3.3% due to noise from spurious correlations for rare pairs. Reducing weighting strength to $\lambda = 0.5$ yields performance between CWCS and Standard CP. Using symmetric co-occurrence produces results nearly identical to Tree-based CQioC at 9.91 labels, confirming both methods essentially model symmetric pairwise dependencies. The full range spans only 5% from 9.82 to 10.31 labels, indicating modest impact of dependency modeling under extreme imbalance with weak classifiers.

Configuration	Avg Set Size	Coverage
CWCS full ($\lambda = 1.0$, threshold=0.10)	9.98	0.909
Without filtering (all edges)	10.31	0.908
Reduced strength ($\lambda = 0.5$)	9.89	0.908
Symmetric co-occurrence	9.91	0.906
No weighting ($\lambda = 0$)	9.82	0.907

Removing the co-occurrence filtering threshold by including all label pairs regardless of how weak their empirical co-occurrence increases the average set size from 9.98 to 10.31 labels, a degradation of 3.3%. This increase stems from noise in co-occurrence estimates for rare label pairs. For instance, the Hernia-Fibrosis co-occurrence is estimated from only 5 to 10 joint occurrences in 78,440 training images, making the conditional probability estimate highly uncertain. When such weak correlations are included, CWCS occasionally boosts labels based on coincidental rather than genuine clinical relationships. The filtering threshold of 0.10 serves as regularization, focusing on well-supported dependencies while discarding unreliable estimates.

Reducing the weighting strength parameter from $\lambda = 1.0$ to $\lambda = 0.5$ yields an average set size of 9.89 labels, nearly matching the full method. This suggests that even modest dependency weighting provides most of the benefit, with diminishing returns beyond $\lambda = 0.5$ in this experimental setting. The relationship between λ and set size is nonlinear: at $\lambda = 0$, no co-occurrence information is used; at $\lambda = 0.5$, some boost is provided but insufficient to substantially alter many inclusion decisions; at $\lambda = 1.0$, the boost becomes meaningful for high-co-occurrence pairs like Infiltration-Pneumonia at 0.42 providing 42% weight increase.

Using symmetric co-occurrence by setting $M_{jk} = M_{kj} = \max(M_{jk}, M_{kj})$ produces an average set size of 9.91 labels, essentially identical to Tree-based CQioC’s 9.91. This near-perfect equivalence is not coincidental: both symmetric CWCS and Tree-based methods model pairwise dependencies without directionality. The primary difference lies in aggregation, with CWCS using maximum co-occurrence and Tree-based using the tree structure, but under severe class imbalance these approaches converge to similar behavior. The symmetric variant loses the ability to distinguish strong directional relationships like ”Infiltration predicts Pneumonia with probability 0.42” from the weak reverse ”Pneumonia predicts Infiltration with probability 0.15.”

Setting $\lambda = 0$ recovers Standard CP with independent calibration, achieving 9.82 average set size. Paradoxically, this represents the smallest sets among all configurations tested. This finding requires careful interpretation: while CWCS with dependency modeling produces marginally larger sets on average (9.98 versus 9.82), it may produce more clinically coherent prediction sets by favoring plausible label combinations. A set containing Infiltration, Pneumonia, and Consolidation—three pathologies with established co-occurrence patterns—provides more actionable clinical information than a set of the same size containing Infiltration, Hernia, and

Fibrosis which represents an implausible combination. Our aggregate efficiency metrics cannot capture this qualitative difference in clinical utility.

7 Discussion

Our comprehensive evaluation of CWCS and three baseline conformal prediction methods on ChestX-ray14 reveals nuanced insights that both confirm and challenge initial expectations about label-aware uncertainty quantification under extreme class imbalance. While CWCS successfully models asymmetric dependencies prevalent in chest radiography diagnosis, the aggregate efficiency gains measured by average prediction set size are modest, with all methods except APS producing sets of 9.82 to 9.98 labels representing only a 1.6% spread. This similarity stems from the dominant influence of conservative calibration required when base classifiers produce poorly calibrated probability estimates across severely imbalanced classes.

The key insight from our experiments is that dependency modeling matters most when the base classifier provides reasonably confident and well-calibrated predictions for at least some labels, enabling the conformal method to make informed decisions about which other labels to include based on co-occurrence structure. In our experimental setting with a classifier achieving macro F1 of 0.32, many test instances receive moderate probabilities spread across multiple unrelated labels without clear confidence peaks. In such cases, all conformal methods must conservatively include many labels to guarantee coverage, and the opportunity to exploit dependency structure is limited.

However, aggregate set size metrics mask important qualitative differences. CWCS systematically produces prediction sets that respect clinically meaningful asymmetric dependencies, such as including both Infiltration and Pneumonia when Infiltration is confidently predicted and co-occurrence evidence is strong, rather than including Infiltration alongside clinically implausible labels. This alignment between algorithmic reasoning based on conditional probabilities and clinical reasoning based on pathophysiological relationships offers value in interpretability and trust that quantitative efficiency metrics do not fully capture. Radiologists reviewing CWCS’s suggestions may find them more consistent with how they naturally think about differential diagnosis, potentially improving acceptance and appropriate use of the AI system.

The finding that improved base classifier performance represents the highest-leverage opportunity for reducing conformal set sizes has important implications for research priorities. Efforts to advance from F1 of 0.32 to 0.60 through better architectures, more training data, or multi-modal learning would likely yield greater returns than optimizing dependency modeling approaches. As classifiers improve and produce tighter prediction sets averaging 3 to 4 labels rather than 9 to 10, the composition and clinical coherence of those sets becomes paramount, and we expect the advantages of methods like CWCS to become more pronounced in such improved settings.

Our work demonstrates that domain-specific structure can be seamlessly integrated into the conformal prediction framework while preserving distribution-free finite-sample coverage guarantees. The co-occurrence weighting mechanism represents one instantiation of this broader principle, and future work could explore incorporating other forms of structure such as causal relationships from medical knowledge, hierarchical taxonomies from disease ontologies, or temporal dependencies in sequential diagnosis. The modularity of conformal prediction—where dependency modeling enters through the nonconformity score definition without compromising theoretical validity—makes such extensions straightforward to implement and evaluate.

8 Conclusion

We introduced Co-occurrence Weighted Conformal Sets (CWCS), a conformal prediction framework that explicitly models asymmetric label dependencies for multi-label medical image clas-

sification. Through rigorous evaluation on ChestX-ray14 encompassing over 112,000 chest radiographs with 14 pathology labels exhibiting extreme class imbalance, we demonstrated that CWCS maintains valid 90% coverage guarantees while producing prediction sets aligned with clinical co-occurrence patterns including strong asymmetries like Infiltration predicting Pneumonia with measure 0.39 and Cardiomegaly predicting Effusion with measure 0.30.

Under the challenging conditions of our experimental setting with moderate base classifier performance achieving macro F1 of 0.32, all conformal methods except APS produce similar aggregate efficiency measured by average set sizes of 9.82 to 9.98 labels. This convergence indicates that conservative calibration necessitated by poor probability estimates dominates over dependency modeling effects in this regime. However, this aggregate similarity masks CWCS’s value in systematically respecting directional clinical relationships, offering advantages in interpretability and clinical coherence beyond what efficiency metrics capture.

Our work contributes both a practical method for incorporating asymmetric dependencies into conformal prediction and broader insights about when and why dependency modeling provides value in uncertainty quantification for medical AI systems. As base classifier quality improves and prediction sets tighten, we expect the advantages of dependency-aware methods to become more pronounced. The framework we have developed provides a foundation for future research into structured conformal prediction approaches that leverage domain knowledge while maintaining the distribution-free coverage guarantees essential for responsible deployment in high-stakes medical applications.

References

- [1] V. Vovk, A. Gammerman, and G. Shafer. *Algorithmic learning in a random world*. Springer Science & Business Media, 2005.
- [2] H. Papadopoulos et al. “Inductive confidence machines for regression”. In: *European Conference on Machine Learning*. 2002, pp. 345–356.
- [3] Y. Romano, E. Patterson, and E. Candès. “Conformalized quantile regression”. In: *Advances in Neural Information Processing Systems*. Vol. 32. 2019.
- [4] A. N. Angelopoulos et al. *Uncertainty sets for image classifiers using conformal prediction*. arXiv preprint arXiv:2009.14193. 2020.
- [5] M. Sadinle, J. Lei, and L. Wasserman. “Least ambiguous set-valued classifiers with bounded error levels”. In: *Journal of the American Statistical Association* 114.525 (2019), pp. 223–234.
- [6] D. Stutz et al. “Learning optimal conformal classifiers”. In: *International Conference on Learning Representations*. 2021.
- [7] Y. Romano, M. Sesia, and E. Candès. “Classification with valid and adaptive coverage”. In: *Advances in Neural Information Processing Systems*. Vol. 33. 2020, pp. 3581–3591.
- [8] M. Fontana, G. Zeni, and S. Vantini. “Conformal prediction: a unified review of theory and new challenges”. In: *Bernoulli* 29.1 (2023), pp. 1–23.
- [9] Jeremy Irvin et al. “CheXpert: A Large Chest Radiograph Dataset with Uncertainty Labels and Expert Comparison”. In: *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 33. 2019, pp. 590–597.
- [10] S. C. Huang et al. “GLoRIA: A multimodal global-local representation learning framework”. In: *IEEE/CVF International Conference on Computer Vision*. 2021, pp. 3942–3951.

- [11] A. Dosovitskiy, L. Beyer, A. Kolesnikov, et al. “An image is worth 16x16 words: Transformers for image recognition at scale”. In: *International Conference on Learning Representations*. 2020.
- [12] T. Chen et al. “A simple framework for contrastive learning of visual representations”. In: *International Conference on Machine Learning*. 2020, pp. 1597–1607.
- [13] Y. Gal and Z. Ghahramani. “Dropout as a Bayesian approximation: Representing model uncertainty”. In: *International Conference on Machine Learning*. 2016, pp. 1050–1059.
- [14] B. Lakshminarayanan, A. Pritzel, and C. Blundell. “Simple and scalable predictive uncertainty estimation using deep ensembles”. In: *Advances in Neural Information Processing Systems*. Vol. 30. 2017.
- [15] A. N. Angelopoulos et al. *Learn then test: Calibrating predictive algorithms to achieve risk control*. arXiv preprint arXiv:2110.01052. 2022.
- [16] C. Lu et al. “Fair conformal predictors for applications in medical imaging”. In: *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 36. 11. 2022, pp. 12008–12016.
- [17] Kaiming He et al. “Deep Residual Learning for Image Recognition”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2016), pp. 770–778.
- [18] C. Chow and C. Liu. “Approximating Discrete Probability Distributions with Dependence Trees”. In: *IEEE Transactions on Information Theory* 14.3 (1968), pp. 462–467.