# Interpretability of Automatic Infectious Disease Classification Analysis with Concept Discovery

**Elena Sizikova, Joshua Vendrow, Xu Cao, Rachel Grotheer, Jamie Haddock, Lara Kassab, Alona Kryshchenko, Thomas Merkh, R. W. M. A. Madushani, Kenny Moise, Annie Ulichney, Huy V. Vo, Chuntian Wang, Megan Coffee, Kathryn Leonard, Deanna Needell**

arXiv:2209.02415v2 [cs.CV] 14 Nov 2022

## Abstract

Automatic infectious disease classification from images can facilitate needed medical diagnoses. Such an approach can identify diseases, like tuberculosis, which remain under-diagnosed due to resource constraints and also novel and emerging diseases, like monkeypox, which clinicians have little experience or acumen in diagnosing. Avoiding missed or delayed diagnoses would prevent further transmission and improve clinical outcomes. In order to understand and trust neural network predictions, analysis of learned representations is necessary. In this work, we argue that automatic discovery of concepts, i.e., human interpretable attributes, allows for a deep understanding of learned information in medical image analysis tasks, generalizing beyond the training labels or protocols. We provide an overview of existing concept discovery approaches in medical image and computer vision communities, and evaluate representative methods on tuberculosis (TB) prediction and monkeypox prediction tasks. Finally, we propose NMFx, a general NMF formulation of interpretability by concept discovery that works in a unified way in unsupervised, weakly supervised, and supervised scenarios[1].

**Keywords:** Explainability, Non-Negative Matrix Factorization (NMF), Neural Networks

---

## 1. Introduction

Visualization of learned features is crucial for a better understanding of patterns learned by neural networks. Concept discovery approaches are a type of explainable artificial intelligence (XAI) technique that identifies high-level and human-understandable explanations for the predictive behavior of a model. In this work, we analyze infectious disease classification neural networks using concept discovery approaches.

While neural networks achieve impressive results, they mostly remain black-box models. Visualization and interpretation of learned representations remains an important challenge in their analysis (Lou et al., 2012), and is one of the main hurdles to wider adoption of data mining and AI technology for many applications (Kim and Canny, 2017; Bussmann et al., 2020; Islam et al., 2021). Specifically for healthcare (Singh et al., 2020), explainability is crucial for analysis of success and failure cases that is typically done by medical professionals. There have been instances where AI tools have been found to determine classifications based on meta or extraneous data, which would not be reliable when repeated in medical contexts. Other tools are simply black boxes. The contexts in which and the patients for which these tools may fail matter. Biological plausibility is a key factor in determining whether an intervention will be adopted in medicine. It is important that a diagnosis is based not

on a correlation that may change but on a substantive basis. AI tools may be biased and may not perform as well in certain populations, perhaps related to race, gender, age, or other characteristics, depending on the composition of training datasets. To the clinician, what matters is the patient in front of them, not the average performance. Medical ethics requires physicians to "do no harm". Medical legal liability requires clinicians to be responsible for the decisions they make with the use of AI tools. Many AI tools are developed and deployed without Randomized Controlled Trials to demonstrate replicable results across a broad array of populations and sub-populations. Other interventions in medicine require substantial validation in a real world setting. What works in the lab or in a specific population may not work in the real world clinical situations where clinicians work. The impact can be devastating in medicine if clinical diagnosis is affected; clinicians may even learn to rely on tools and be deskilled, leading to worsened outcomes than would result without the tool. In medicine, not all that seems good is in fact good. What may seem like a step forward can inadvertently be a step back.

Many XAI techniques require a particular network architecture, access to network weights, and back-propagation to generate an interpretation heat map (Chattopadhay et al., 2018; Zhou et al., 2016). Additionally, it is crucial to link generated explanations to existing human knowledge. Therefore, for studying infectious disease classification models, we seek XAI methods that are model-agnostic, fast, provide a global overview of model behavior, and generate human interpretable outputs. Concept extraction (CE) approaches are a class of XAI techniques that seek to explain the decision-making process of a neural network in human-understandable terms, or so-called concepts. While gaining popularity

in the image analysis community, CE approaches are not widely used for analyzing medical image classification models. In this work, we seek to explain the behavior of neural network (NN) infectious disease classification models using post-hoc visual explanations generated by CE methodology proposed in (Collins et al., 2018; Oramas Mogrovejo et al., 2019; Posada-Moreno et al., 2022) (we provide an extensive discussion of related work in the Appendix). As a result we can visualize information encoded in features and analyze whether encoded topics consistently capture image regions that are semantically related. In comparison to existing techniques, the presented method is flexible to work in unsupervised, semi-supervised or weakly supervised fashion, and provided labels do not need to correspond to the labels that the underlying network was trained for. The extracted concepts provide a useful visualization tool to medical image professionals, and match the intuition of where doctors would search an image for the presence of disease. In summary, our contributions are as follows: (1) To understand infectious disease classification behavior, we propose to extract and analyze concepts generated by CE approaches on image datasets for tuberculosis (TB) and monkeypox, both infectious diseases that are prone to underdiagnosis and require prompt identification for treatment and to prevent others being infected. (2) We compare concepts generated by the feature factorization methods proposed by (Oramas Mogrovejo et al., 2019), (Posada-Moreno et al., 2022) and (Collins et al., 2018), extending the latter to accept weak supervision (NMFx framework). We study medical image datasets (TBX11K (Liu et al., 2020), monkeypox (Moise et al., 2022)) as well as PASCAL-VOC 2010 (Everingham et al., 2010), a public image analysis dataset.

2

## 2. Methodology

The first method we consider is proposed in (Collins et al., 2018). The general idea is to extract feature representations of images using a convolutional neural network and factorize the resulting reshaped matrix using NMF with $K$ topics. After a reshaping procedure, the weight matrix results in a set of $K$ heat maps that visually explain information encoded in the features. The process is shown in Figure 1.

Let $X \in \mathbb{R}^{n_1 \times n_2}$ denote the nonnegative data matrix of $n_2$ data points in $\mathbb{R}^{n_1}$. Lee and Seung (Lee and Seung, 1999) propose to decompose X into a topic matrix $A$ and a weight matrix $S$ using the following, Frobenius-norm optimization objective:

$$\min_{A,S} \|X - AS\|_F^2. \tag{1}$$

Here, $A \in \mathbb{R}_{\geq 0}^{n_1 \times k}$ denotes the topic matrix with $k$ topics and $S \in \mathbb{R}_{\geq 0}^{k \times n_2}$ denotes the representative weight matrix.

**NMF with Image Label Supervision** When information about the data points' labels is available, we can encode it into $Y \in \mathbb{R}^{l \times n_2}$, a binary label matrix where columns correspond to data points in $X$ and rows represent their class membership. Lee et al. (Lee et al., 2009) then propose the classical semi-supervised nonnegative factorization (SSNMF) method, whose objective is given by:

$$\min_{A,B,S} \|X - AS\|_F^2 + \lambda \|Y - BS\|_F^2, \tag{2}$$

where $\lambda \in \mathbb{R}$ is a regularization parameter and $B \in \mathbb{R}^{l \times k}$ is the trained classification matrix. The first term in objective (2) denotes the reconstruction error of the factorization and the second term denotes the classification error. This model simultaneously learns a topic model, defined by the matrix $A$, and a classification model, defined by the matrix $B$. The matrix $S$ provides a representative weight matrix which both fits the topic model and predicts labels. We refer to the resulting generalized concept extraction technique as NMFx.

## 3. Implementation Details

### 3.1. NMF.

Given $n$ images, we obtain their network feature representation $X'$ with dimension $(n, p, d_1, d_2)$, where $p$ is the number of feature maps, and $d_1, d_2$ are the width and height of each feature map, respectively. For example, using the features after the rectified linear unit (ReLU) and the last convolutional layer of VGG-16 (Simonyan and Zisserman, 2014), we have $p = 512$ and $d_1 = 14$ and $d_2 = 14$. $X'$ is then flattened and transposed into matrix $X$ of dimension $(p, n \times d_1 \times d_2)$, before being passed to the NMF or SSNMF optimization objective, and obtaining the factorization $X \approx AS$ with $K$ topics. Note that in this setting, a data point is a $p$-dimensional vector representing a location in an image. The resulting nonnegative weight matrix $S$ is of dimensionality $(K, n \times d_1 \times d_2)$ and, if a train-test split is used, we obtain $S_{\text{test}}$ of dimension $(K, n_{\text{test}} \times d_1 \times d_2)$ using nonnegative least squares, as described above. We reshape $S$ into a heat map tensor of dimension $(n, K, d_1, d_2)$ and up-sample to image resolution $(n, K, w, h)$, where $w$ and $h$ are the width and height of the input images, respectively. The reshaping procedure for $S_{\text{test}}$ is analogous to that of $S$.

Whenever the SSNMF is used as the objective for NMF, we also create a binary class tensor $Y'$ of dimensionality $(n, K, d_1, d_2)$. For each image $i$ of label $l$, the $d_1 \times d_2$ sub-tensor corresponding to this image and label are set to 1, otherwise it is set to 0. Subsequently, $Y'$ is reshaped and transposed to matrix Y of dimensionality $(K, n \times d_1 \times d_2)$ and used in the optimization.
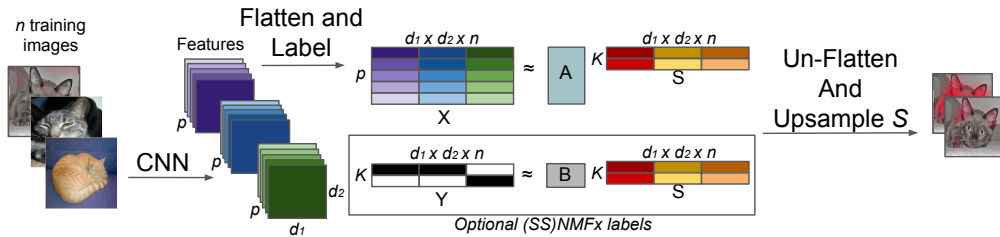
3

Figure 1: *Visual explanation of NMFx. Optional image labels can be used during the modified optimization step, (SS)NMFx. Example above is shown for $K = 2$ and $n = 3$.*

## 4. Datasets

We evaluate concept discovery using NMFx and other techniques on TBX11K (Liu et al., 2020), a public tuberculosis chest X-ray dataset. We rely on the official data split in (Liu et al., 2020) that includes multiple smaller, public sets (Jaeger et al., 2014). We also evaluate our approach on the task of monkeypox classification. For this task, we consider a subset of 356 monkeypox and 345 non-monkepox images from publicly available images in the medical literature and from social media and journalistic sources (Moise et al., 2022). These images include confirmed diagnoses. Monkeypox images were all from clade II and only included if documented as confirmed by PCR (polymerase chain reaction) testing. Monkeypox images were included from all stages of lesion evolution. Comparison images were selected by medical doctors as appearing similar to monkeypox and having similar clinical syndromes (such as herpes, syphilis, varicella, hand foot and mouth disease, and molluscum). Some images were also identified by google image searches for monkeypox images to identify similar appearing images. These images were cropped to include only the skin lesions and eliminate identifying information (such as jewelry, tattoos, facial features, and clothing) which could affect results and to ensure images were closely matched. The monkeypox and non-monkeypox image datasets were compared to show similar proportions by age (adult or child), skin type, sex and gender, and body part affected.

## 5. Results

We now analyze the proposed technique and its variations.

**Visual Topics Found in Tuberculosis Analysis** In Figure 2, we visualize topics found using NMFx in the VGG-16 network trained for a tuberculosis classification task. The colors were arbitrarily chosen. In particular, topic 1 (yellow color) corresponds to the areas of interest corresponding to predicting tuberculosis. Tuberculosis disproportionately affects the upper lung fields, which are highlighted in yellow, unlike other infections, which makes this area very important in identifying tuberculosis. Topic 2 (red) highlights areas that are less important, but can be involved, in tuberculosis diagnosis, such as the lower lung fields (potentially has dependent pleural effusions) or central sections (potentially has pericardial effusion or hilar lymphadenopathy). The Topic 3 (green) corresponds to areas outside the lung fields which are not expected to be helpful in TB diagnosis. We also find that visual topics consistently highlight similar anatomic areas across a variety of input example anatomies.
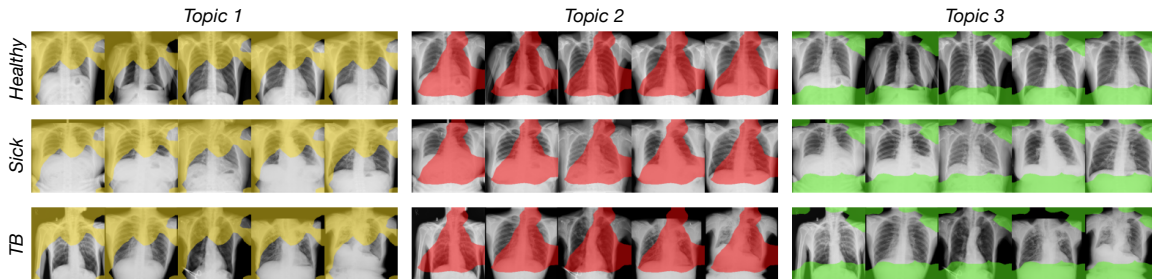
Figure 2: *Visual explanations using NMFx for a VGG-16 Tuberculosis classification task.* Clinically, Topic 1 corresponds to areas used most in diagnosis, Topic 2 corresponds to areas sometimes involved in diagnosis, and Topic 3 corresponds to unrelated areas. Colors are chosen randomly.

**Visual Topics Found in Monkeypox Analysis**  In Figure 3, we look at the topics found by the EfficientNet-B3 network trained for the monkeypox classification task. We find that the topics are centered on the lesions and regions corresponding to monkeypox lesions (second row). On the other hand, in examples with non-monkeypox skin conditions (first row), the topics identify some visually similar lesions but are more scattered and occupy a smaller surface area.

### 5.1. Comparison to Other Automatic Concept Extraction Techniques

We compare the above approach to two other, representative automatic concept extraction methods: Oramas (Oramas Mogrovejo et al., 2019) and ECLAD (Posada-Moreno et al., 2022) in Figure 4. Please see Appendix for additional implementation details details. For this experiment, we consider approximately 30% of the the TBX11K Liu et al. (2020) dataset, due to computational limitations of running Collins et al. (2018); Posada-Moreno et al. (2022). In contrast to ECLAD and the method of Oramas, the NMFx method identifies larger and more consistently positioned regions in the input X-ray images.

## 6. Conclusion

In this work, we explore three automatic concept discovery approaches as explainability tools for neural networks trained to classify infectious diseases. We find that generated concepts highlight areas where disease lesions are generally found (as in tuberculosis) or noted in specific, available images (monkeypox). The highlighted areas allow clinicians to see if the tool is evaluating the same areas of interest they have themselves identified. This helps clinicians gauge plausibility and reliability. With novel and emerging diseases, like monkeypox as well as others we expect to increasingly see in the future, clinicians have little specific past clinical experience to rely on, but would identify areas of interest on radiologic imaging or dermatologic exams and it will be reassuring that the tool is focusing on the same area. This will help clinicians understand what may lead to false negative or positive results, especially if they need to override the tool's determinations. This will allow tools to be deployed more effectively in early outbreaks with an extra layer of oversight, especially as outbreaks may result in different clinical findings, whether radiologic or dermatologic, from the training
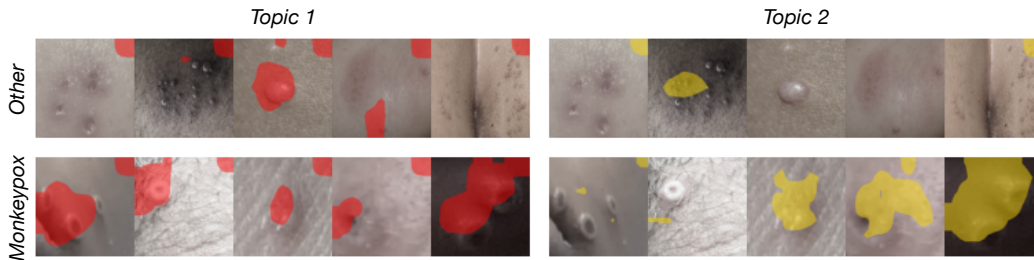
Figure 3: *Visual explanations using NMFx for a monkeypox classification task.*
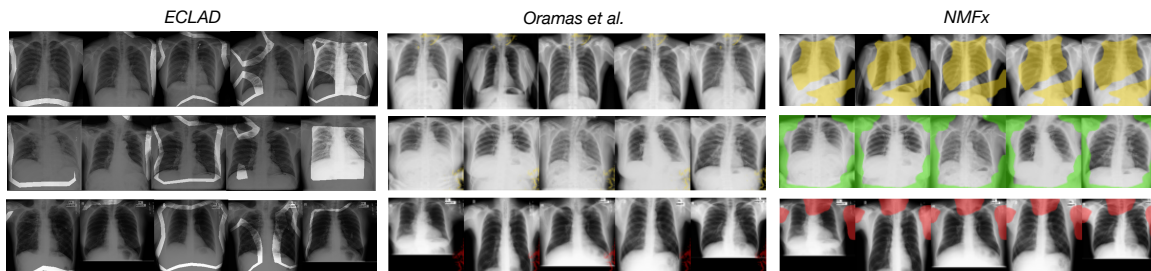


Figure 4: *Comparison of visual concepts identified by different techniques.*

datasets, as the outbreak expands. We often see as outbreaks spread beyond initial groups (whether due to geographic spread or due to different risk factors, such as we have seen with monkeypox) that clinical findings may change. Clinicians may also need to differentiate a new pathogen from other diseases as an outbreak spreads, given diseases prevalence varies by geography and between populations, and so monitoring how the tool works with different diseases will help. A lightweight tool that could provide an additional layer of oversight for new tools would help harness pattern recognition more effectively in evolving outbreaks. Inspired by the work of (Collins et al., 2018), we introduce NMFx, a lightweight and general framework for analyzing activations based on nonnegative matrix decomposition (NMF) using no concept supervision or using image labels. We show that jointly factoring images in the feature space of a classification neural network allows extracting information about localization properties of the network without any image-level category annotations.

## References

Francesco Bodria, Fosca Giannotti, Riccardo Guidotti, Francesca Naretto, Dino Pedreschi, and Salvatore Rinzivillo. Benchmarking and survey of explanation methods for black box models. *arXiv preprint arXiv:2102.13076*, 2021.

Moritz Böhle, Fabian Eitel, Martin Weygandt, and Kerstin Ritter. Layer-wise relevance propagation for explaining deep neural network decisions in mri-based alzheimer's disease classification. *Frontiers in Aging Neuroscience*, 11:194, 2019.

Niklas Bussmann, Paolo Giudici, Dimitri Marinelli, and Jochen Papenbrock. Explainable ai in fintech risk management. *Frontiers in Artificial Intelligence*, 3:26, 2020.

Aditya Chattopadhay, Anirban Sarkar, Prantik Howlader, and Vineeth N Balasubramanian. Grad-cam++: Generalized gradient-based visual explanations for deep convolutional networks. In *WACV*, 2018.

Bingzhi Chen, Jinxing Li, Guangming Lu, and David Zhang. Lesion location attention guided network for multi-label thoracic disease classification in chest x-rays. *IEEE Journal of Biomedical and Health Informatics*, 2019.

Edo Collins, Radhakrishna Achanta, and Sabine Susstrunk. Deep feature factorization for concept discovery. In *ECCV*, 2018.

Arun Das and Paul Rad. Opportunities and challenges in explainable artificial intelligence (xai): A survey. *arXiv e-prints*, pages arXiv–2006, 2020.

M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The PASCAL Visual Object Classes Challenge 2010 (VOC2010) Results, 2010.

Feng-Lei Fan, Jinjun Xiong, Mengzhou Li, and Ge Wang. On interpretability of artificial neural networks: A survey. *IEEE Transactions on Radiation and Plasma Medical Sciences*, 2021.

Amirata Ghorbani, James Wexler, James Y Zou, and Been Kim. Towards automatic concept-based explanations. *NeurIPS*, 32, 2019.

Miriam Hägele, Philipp Seegerer, Sebastian Lapuschkin, Michael Bockmayr, Wojciech Samek, Frederick Klauschen, Klaus-Robert Müller, and Alexander Binder. Resolving challenges in deep learning-based analyses of histopathological images using explanation methods. *Scientific reports*, 2020.

Sara Hooker, Dumitru Erhan, Pieter-Jan Kindermans, and Been Kim. A benchmark for interpretability methods in deep neural networks. *Advances NeurIPS*, 2019.

Sheikh Rabiul Islam, William Eberle, Sheikh Khaled Ghafoor, and Mohiuddin Ahmed. Explainable artificial intelligence approaches: A survey. *arXiv preprint arXiv:2101.09429*, 2021.

Stefan Jaeger, Sema Candemir, Sameer Antani, Yì-Xiáng J Wáng, Pu-Xuan Lu, and George Thoma. Two public chest x-ray datasets for computer-aided screening of pulmonary diseases. *Quantitative Imaging in Medicine and Surgery*, 2014.

Weina Jin, Xiaoxiao Li, and Ghassan Hamarneh. Evaluating explainable ai on a multi-modal medical imaging task: Can existing algorithms fulfill clinical requirements? In *AAAI*, 2022.

Ashkan Khakzar, Shadi Albarqouni, and Nassir Navab. Learning interpretable features via adversarially robust optimization. In *MICCAI*, 2019.

Been Kim, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler, Fernanda Viegas, et al. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav). In *ICML*, 2018.

Jinkyu Kim and John Canny. Interpretable learning for self-driving cars by visualizing causal attention. In *ICCV*, 2017.

Daniel D Lee and H Sebastian Seung. Learning the parts of objects by non-negative matrix factorization. *Nature*, 1999.

Hyekyoung Lee, Jiho Yoo, and Seungjin Choi. Semi-supervised nonnegative matrix factorization. *IEEE Signal Processing Letters*, 2009.

Yun Liu, Yu-Huan Wu, Yunfeng Ban, Huifang Wang, and Ming-Ming Cheng. Rethinking computer-aided tuberculosis diagnosis. 2020.

Yin Lou, Rich Caruana, and Johannes Gehrke. Intelligible models for classification and regression. In *SIGKDD*, 2012.

Kenny Moise, Cecilia Crews, Elena Sizikova, and Megan Coffee. Development of a diverse and representative database for the training of artificial intelligence tools for monkeypox classification. *In Preparation*, 2022.

José Antonio Oramas Mogrovejo, Kaili Wang, and Tinne Tuytelaars. Visual explanation by interpretation: Improving visual feedback capabilities of deep neural networks. In *ICLR*, 2019.

Andres Felipe Posada-Moreno, Nikita Surya, and Sebastian Trimpe. Eclad: Extracting concepts with local aggregated descriptors, 2022.

Sivaramakrishnan Rajaraman, Sema Candemir, George Thoma, and Sameer Antani. Visualizing and explaining deep learning predictions for pneumonia detection in pediatric chest radiographs. In *Medical Imaging 2019: Computer-Aided Diagnosis*, 2019.

Wojciech Samek, Alexander Binder, Grégoire Montavon, Sebastian Lapuschkin, and Klaus-Robert Müller. Evaluating the visualization of what a deep neural network has learned. *TNNLS*, 2016.

Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Gradcam: Visual explanations from deep networks via gradient-based localization. In *ICCV*, 2017.

Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv*, 2014.

Amitojdeep Singh, Sourya Sengupta, and Vasudevan Lakshminarayanan. Explainable deep learning models in medical image analysis. *Journal of Imaging*, 2020.

Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin Riedmiller. Striving for simplicity: The all convolutional net. *arXiv preprint arXiv:1412.6806*, 2014.

Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *ICML*, 2019.

Bas HM van der Velden, Markus HA Janse, Max AA Ragusi, Claudette E Loo, and Kenneth GA Gilhuijs. Volumetric breast density estimation on mri using explainable deep learning regression. *Scientific Reports*, 2020.

Bas HM van der Velden, Hugo J Kuijf, Kenneth GA Gilhuijs, and Max A Viergever. Explainable artificial intelligence (xai) in deep learning-based medical image analysis. *Medical Image Analysis*, 2022.

Linda Wang, Zhong Qiu Lin, and Alexander Wong. Covid-net: A tailored deep convolutional neural network design for detection of covid-19 cases from chest x-ray images. *Scientific Reports*, 2020.

Mengjiao Yang and Been Kim. Bim: Towards quantitative evaluation of interpretability methods with ground truth. *arXiv preprint arXiv:1907.09701*, 2019.

Jason Yosinski, Jeff Clune, Anh Nguyen, Thomas Fuchs, and Hod Lipson. Understanding neural networks through deep visualization. *arXiv preprint arXiv:1506.06579*, 2015.

Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *ECCV*, 2014.

Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features for discriminative localization. In *CVPR*, 2016.

## Appendix

### 6.1. Related Work

**Visual Explanation.** There have been a number of recent methods (Das and Rad, 2020) whose goal is to visualize the learned representations of convolutional neural networks (CNNs) and explain their properties. Early techniques such as Zeiler and Fergus (Zeiler and Fergus, 2014) proposed deconvolution networks to identify parts of the input image that activate each neuron unit. Guided backpropagation (Springenberg et al., 2014) evaluates the effect of each neuron with respect to the input. (Yosinski et al., 2015) introduces an optimization technique to synthesize images that highly activate specific neurons. (Zhou et al., 2016) shows that CNN layers act as unsupervised object detectors, and introduce class activation maps (CAM), a technique that uses the global average pooling (GAP) layer to generate importance heat maps for each class. This technique is limited to CNNs with a GAP layer and requires training of linear classifiers on top of the original networks. To address these limitations, Grad-CAM (Selvaraju et al., 2017) and a number of other approaches (Selvaraju et al., 2017; Chattopadhay et al., 2018) have been developed, in particular, to extend CAM to other architectures.

**XAI Classification and Evaluation.** As discussed in (Posada-Moreno et al., 2022), XAI techniques can be grouped into two classes: local and global explanation techniques. *Local methods*, also known as feature attribution techniques, seek to analyze and explain model behavior on a single data point (e.g., image). On the other hand, *global methods* seek to analyze model behavior on a group of images and extract a set of representative concepts that would explain information learned by the models. A *concept* can be any high level explanation, e.g., an object part or a super-pixel. Concepts can be provided by the user (Kim et al., 2018) or extracted automatically (Ghorbani et al., 2019). Automatic methods first identify a set of concepts and then use optimization to determine their importance. For example, (Oramas Mogrovejo et al., 2019) solves a Lasso problem to identify which which CNN activations are important for specific classes, later relying on guided back propagation (Springenberg et al., 2014) to propagate activations and generate visual explanations. (Ghorbani et al., 2019; Posada-Moreno et al., 2022) use TCAV (Kim et al., 2018) to identify top concepts.

Finally, there exist a number of techniques for evaluating visual explanation and neural interpretations methods (Samek et al., 2016; Yang and Kim, 2019; Hooker et al., 2019; Oramas Mogrovejo et al., 2019), in part due to the ambiguity of discriminative features between the classes. Please see (Bodria et al., 2021) for a survey. Since the goal of our work is to visualize concepts learned by disease classification networks, we ask a board-certified doctor to analyze and comment on the information encoded in the concepts.

**XAI for Medical Image Analysis.** Nowadays, researchers in medical imaging are increasingly using XAI to explain the results of their AI algorithms. Van der Velden et al. summarize related XAI research in their survey (van der Velden et al., 2022), finding that most of the XAI papers in the medical analysis used post-hoc explanations as contrasted with model-based explanations, which means the explanation was provided on a deep neural network (DNN) that had already been trained instead of being incorporated in the training process (van der Velden et al., 2022). Mainstream XAI methods for medical data include CAM (Khakzar et al., 2019), Grad-

CAM (Chen et al., 2019), local interpretable model-agnostic explanations (LIME) (Rajaraman et al., 2019), layer-wise relevance propagation (LRP) (Böhle et al., 2019; Hägele et al., 2020), and Shapley additive explanations (SHAP) (van der Velden et al., 2020). The newest attribution-based methods, such as GSInquire and concept-based methods, such as TCAV, also show high performance in many new tasks (Singh et al., 2020; Wang et al., 2020; Posada-Moreno et al., 2022). Additionally, (Fan et al., 2021) reviews applications of interpretability in medicine in a comprehensive taxonomy. Compared with the rapid development at the methodological level of single modality, evaluations on XAI for multi-modal medical imaging task is challenging and still immature (Jin et al., 2022). XAI in medical imaging analysis gives insight into how machine learning and neural networks can make AI-based clinical decisions more understandable and help medical doctors diagnose patients accurately.

### 6.2. Additional Experimental Results

**Feature, Optimization, and Post-Processing Details** We evaluate the VGG-16 (Simonyan and Zisserman, 2014) and the EfficientNet-B3 Tan and Le (2019) networks as feature extractors. In VGG-16, features are extracted after the ReLU following the last convolutional layer (layer 29). In EfficientNet-B3, features are extracted after the `top_activation` layer. In each case, each objective is optimized until convergence. For all experiments, the default number of topics is the total number of classes, unless otherwise specified. All experiments were performed in Python, on a single node in a cluster with 200 GB memory and a 32GB GPU. Implementation details of each method are described below.

|          | Monkeypox | Tuberculosis |
|----------|-----------|--------------|
| Accuracy | 0.8864    | 0.9837       |

Table 1: Classification accuracy of neural networks on two considered medical classification tasks.



(a) Correctly classified monkeypox (TP).



(b) Correctly classified non-monkeypox (TN).



(c) False predictions of monkeypox (FP).



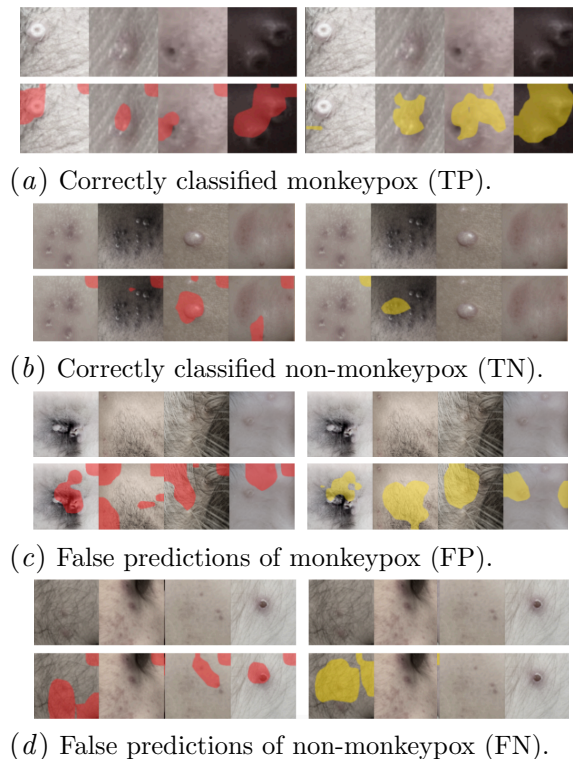(d) False predictions of non-monkeypox (FN).

Figure 5: Examples of predictions in the monkeypox classification tasks with corresponding NMFx topics.

**Classification Accuracy** For reference, in Table 1, we list the test accuracies of the neural networks trained for both classification tasks.

**Ablation Study: Prediction Analysis** In Figure 5, we analyze examples that correctly and incorrectly classified and study their corresponding topics. We find that the model correctly classifies examples (both

monkeypox and non-monkeypox) which are prototypical examples of disease. In these cases, the NMFx topics also consistently identify areas corresponding to disease. On the other hand, incorrect predictions correspond to out of distribution images. In these cases, the visual topics are scattered and not consistently identifying a particular region.

**Ablation Study: Effect of Label Supervision and Number of Topics** We can use the supervised version of NMFx to specify labels that the input images belong to, and visually check whether the extracted topics can be grouped into labels. To demonstrate the ability of supervision to impact the heat maps produced by our method, we apply both unsupervised NMFx and supervised NMFx on a subset of images from seven classes (cow, cat, dog, bird, car, aeroplane, and bicycle) of the PASCAL-VOC 2010 (Everingham et al., 2010) data set. We find that our method works well for both black and white and color images, in particular, handling the larger variation in objects and their backgrounds found in PASCAL-VOC 2010. In Figure 6, we display the heat maps produced for test images resulting from running unsupervised NMFx at (a) $K = 2$ and (d) $K = 3$, as well as supervised NMFx with three supervision labeling. In, (b) we use K=2 with two classes: animals (cow, cat, dog, bird) and vehicles (car, aeroplane, bicycle). In (c), we use $K = 2$ with two classes: flying objects (bird, aeroplane) and non-flying objects (cow, cat, dog, car, bicycle). In (e), we use $K = 3$ with three classes: land animals (cow, dog, cat), flying objects (bird, aeroplane), and land vehicles (car, bicycle). In each case, adding supervision causes the heat maps to better match the provided class labels. For $K = 2$, we see that each of the two groupings of the classes in (b) and (c) shift the heat maps in the desired direction from the heat maps produced
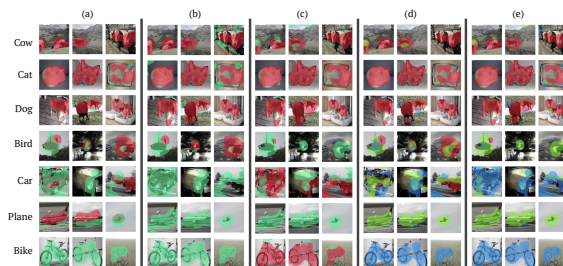


Figure 6: *Visualization of heat maps produced by both unsupervised NMFx and supervised NMFx on a subset of images from the PASCAL-VOC 2010, as we vary the provided labels to demonstrate their impact on the heat maps. In (a) and (d), we run unsupervised NMFx with $K = 2$ and $K = 3$ topics, respectively. In (b) and (c), we report results of (SS)NMFx with $K = 2$ topics and trained with class labels that separate (b) animals from vehicles and (c) flying objects from non-flying objects. In (e), we run (SS)NMFx with $K = 3$ topics and provided class labels that separate land animals, land vehicles, and flying objects.*

in (a). Here, (c) especially shifts the heat maps to group together objects that were distinct in (a), such as grouping cows with bikes (non-flying objects) and grouping birds with aeroplanes (flying objects), to respect the provided labels. Similarly, for $K = 3$, the supervision leads to bikes and cars (land vehicles) being grouped together. This suggests that with supervision, our method is able to identify class-distinct features. We have also compared NMF and SSNMF results on the medical image applications discussed earlier, however, we did not observe consistently different performance.

**Comparison to Other Techniques** The first method we compare to is the method of (Oramas Mogrovejo et al., 2019), which follows a two step approach: first, a Lasso optimization problem is set up to identify the importance of features to the class of interest. Second, once filters that are important for prediction a certain class are identified, guided back-propagation (Springenberg et al., 2014) is applied to generate a saliency heatmap. In comparison, our method generates heatmaps in a single step, and does not require backpropagation. For a fair comparison, we re-implement the method of (Oramas Mogrovejo et al., 2019) and compare to our approach using the same network architecture and layers. The second method we compare to is ECLAD (Posada-Moreno et al., 2022), a recent method for automatic explainable concept extraction. Similar to (Ghorbani et al., 2019), ECLAD automatically finds concepts, but additionally locates them within the image. The underlying algorithm rescales multiple levels of the activation map and analyzes each level at the pixel level. We use the default released implementation of ECLAD with the VGG-16 backbone.