

Proactive Pseudo-Intervention: Pre-informed Contrastive Learning For Interpretable Vision Models

Dong Wang¹, Yuewei Yang², Liqun Chen³ Zhe Gan⁴ Ricardo Henao⁵ Lawrence Carin⁵

^{1, 5}Elevance Health, ²Amazon, ³Cubist Systematic, ⁴Apple

⁵Duke University

dong.wang@elevancehealth.com, yueweiya@amazon.com, albert_lqchen@hotmail.com, zhe.gan@apple.com

Abstract

Deep neural networks excel at comprehending complex visual signals, delivering on par or even superior performance to that of human experts. However, ad-hoc visual explanations of model decisions often reveal an alarming level of reliance on exploiting non-causal visual cues that strongly correlate with the target label in training data. As such, deep neural nets suffer compromised generalization to novel inputs collected from different sources, and the reverse engineering of their decision rules offers limited interpretability. To overcome these limitations, we present a novel contrastive learning strategy called *Proactive Pseudo-Intervention* (PPI) that leverages proactive interventions to guard against image features with no causal relevance. We also devise a novel pre-informed salience mapping module to identify key image pixels to intervene and show it greatly facilitates model interpretability. To demonstrate the utility of our proposals, we benchmark it on both standard natural images and challenging medical image datasets. PPI-enhanced models consistently deliver superior performance relative to competing solutions, especially on out-of-domain predictions and data integration from heterogeneous sources. Further, saliency maps of models that are trained in our PPI framework are more succinct and meaningful.

Introduction

Deep neural networks hold great promise in applications requiring the analysis and comprehension of complex imagery. Recent advances in hardware, network architectures, and model optimization, along with the increasing availability of large-scale annotated datasets (Krizhevsky 2009; Deng 2012; Deng et al. 2009), have enabled these models to match and sometimes outperform human experts on a number of tasks, including natural image classification (Krizhevsky, Sutskever, and Hinton 2017), objection recognition (Girshick et al. 2014), disease diagnosis (Sajda 2006), and autonomous driving (Chen et al. 2015), among others.

While deep learning solutions have been positively recognized for their ability to learn *black-box* models in a purely data driven manner, their very nature makes them less credible for their inability to communicate the reasoning for making predictions in a way that is comprehensible to humans (Hooker et al. 2019; Rebuffi et al. 2020). This denies

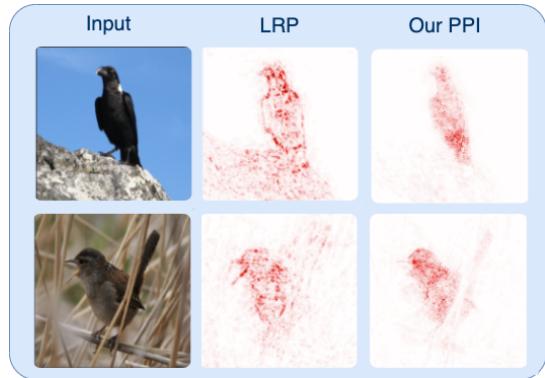


Figure 1: Interpretation for the bird-classification models using saliency maps generated by LRP (*layer-wise relevance propagation*) and our model PPI. LRP shows that naively trained deep model makes decisions based on the background cues (habitat, e.g., rocks, bulrushes) that are spuriously correlated with the bird species, while our pre-informed PPI mostly focuses on the bird anatomy, that generalizes beyond the natural habitat.

consequential applications where the reliability and trustworthiness of a prediction are of primary concern and require expert audit, e.g., in healthcare (Sajda 2006). To stimulate widespread use of deep learning models, a means of interpreting predictions is necessary. However, model interpretation techniques often reveal a concerning fact, that deep learning models tend to assimilate spurious correlations that do not necessarily capture the causal relationship between the input (image) and output (label) (Wang et al. 2020a). This issue is particularly notable in small-sample-size (weak supervision) scenarios or when the sources of non-informative variation are overwhelming, thus likely to cause severe overfitting. These can lead to catastrophic failures on deployment (Fukui et al. 2019; Wang et al. 2019b).

A growing recognition of the issues associated with the lack of interpretable predictions is well documented in recent years (Adebayo et al. 2018; Hooker et al. 2019; Rebuffi et al. 2020). Such phenomenon has energized researchers to actively seek creative solutions. Among these, two streams of work, namely *saliency mapping* (Zhao et al. 2018; Si-

monyan, Vedaldi, and Zisserman 2013; Dabkowski and Gal 2017) and *causal representation learning* (CRL) (Johansson, Shalit, and Sontag 2016; Wang et al. 2020b; Arjovsky et al. 2019), stand out as some of the most promising directions. Specifically, saliency mapping encompasses techniques for *post hoc* visualizations on the input (image) space to facilitate the interpretation of model predictions. This is done by projecting the key features used in prediction back to the input space, resulting in the commonly known *saliency maps*. Importantly, these maps do not directly contribute to model learning. Alternatively, CRL solutions are built on the principles of establishing invariance from the data, and it entails teasing out sources of variation that are spuriously associated with the model output (labels). CRL models, while emphasizing the differences between causation and correlation, are not subject to the rigor of causal inference approaches, because their goal is not to obtain accurate causal effect estimates but rather to produce robust models with better generalization ability relative to their naively learned counterparts (Arjovsky et al. 2019).

In this work, we present *Proactive Pseudo-Intervention* (PPI), a solution that accounts for the needs of causal representation identification and visual verification. Our key insight is the derivation of pre-informed saliency maps which facilitate visual verification of model predictions and enable learning that is robust to (non-causal) associations. While true causation can only be established through experimental interventions, we leverage tools from contrastive representation learning to synthesize pseudo-interventions from observational data. Our procedure is motivated by the causal argument: perturbing the non-causal features will not change the target label.

To motivate, in Figure 1 we present an example to illustrate the benefits of producing causally-informed saliency maps. In this scenario, the task is to classify two bird species (A and B) in the wild. Due to the differences in their natural habitats, A-birds are mostly seen resting on rocks, while B-birds are more commonly found among bulrushes. A deep model, trained naively, will tend to associate the background characteristics with the labels, knowing these strongly correlate with the bird species (labels) in the training set. This is confirmed by the saliency maps derived from the layer-wise relevance propagation (LRP) techniques (Bach et al. 2015): the model also attends heavily on the background features, while the difference in bird anatomy is what causally determines the label. If we were provided with an image of a bird in an environment foreign to the images in the training set, the model will be unable to make a reliable prediction, thus causing robustness concerns. This generalization issue worsens with a smaller training sample size. On the other hand, saliency maps from our PPI-enhanced model successfully focus on the bird anatomy, and thus will be robust to environmental changes captured in the input images.

PPI addresses causally-informed reasoning, robust learning, and model interpretation in a unified framework. A new saliency mapping method, named *Weight Back Propagation* (WBP), is also proposed to generate more concentrated intervention mask for PPI training. The key contributions of this paper include:

- An end-to-end contrastive representation learning strategy PPI that employs proactive interventions to identify causally relevant features.
- A fast and architecture-agnostic saliency mapping module WBP that delivers better visualization and localization performance.
- Experiments demonstrating significant performance boosts from integrating PPI and WBP relative to competing solutions, especially on out-of-domain predictions, data integration with heterogeneous sources and model interpretation.

Background

Visual Explanations. Saliency mapping collectively refers to a family of techniques to understand and interpret black-box image classification models, such as deep neural networks (Adebayo et al. 2018; Hooker et al. 2019; Rebuffi et al. 2020). These methods project the model understanding of the targets, *i.e.*, labels, and their predictions back to the input space, which allows for the visual inspection of automated reasoning and for the communication of predictive visual cues to the user or human expert, aiming to shed model insights or to build trust for deep-learning-based systems.

In this study, we focus on *post hoc* saliency mapping strategies, where saliency maps are constructed given an arbitrary prediction model, as opposed to relying on customized model architectures for interpretable predictions (Fukui et al. 2019; Wang et al. 2019b), or to train a separate module to explicitly produce model explanations (Fukui et al. 2019; Goyal et al. 2019; Chang et al. 2018; Fong and Vedaldi 2017; Shrikumar, Greenside, and Kundaje 2017). Popular solutions under this category include: activation mapping (Zhou et al. 2016; Selvaraju et al. 2017), input sensitivity analysis (Shrikumar, Greenside, and Kundaje 2017), and relevance propagation (Bach et al. 2015). Activation mapping based methods fail at visualizing fine-grained evidence, which is particularly important in explaining medical classification models (Du et al. 2018; Selvaraju et al. 2017; Wagner et al. 2019). Input sensitivity analysis based methods produce fine-grained saliency maps. However, these maps are generally less concentrated (Dabkowski and Gal 2017; Fong and Vedaldi 2017) and less interpretable. Relevance propagation based methods, like LRP and its variants, use complex rules to prioritize positive or large relevance, making the saliency maps visually appealing to human. However, our experiments demonstrate that LRP and its variants highlight spuriously correlated features (boarderlines and backgrounds). By contrast, our WBP backpropagates the weights through layers to compute the contributions of each input pixel, which is truly faithful to the model, and WBP tends to highlight the target objects themselves rather than the background. At the same time, the simplicity and efficiency makes WBP easily work with other advanced learning strategies for both model diagnosis and improvements during training.

Our work is in a similar spirit to (Fong and Vedaldi 2017; Dabkowski and Gal 2017; Chang et al. 2018; Wagner et al. 2019), where meaningful perturbations have been applied

to the image during model training, to improve prediction and facilitate interpretation. Pioneering works have relied on user supplied “ground-truth” explainable masks to perturb (Ross, Hughes, and Doshi-Velez 2017; Li et al. 2018; Rieger et al. 2020), however such manual annotations are costly and hence rarely available in practice. Alternatively, perturbations can be computed by solving an optimization for each image. Such strategies are costly in practice and also do not effectively block spurious features. Very recently, exploratory effort has been made to leverage tools from counterfactual reasoning (Goyal et al. 2019) and causal analysis (O’Shaughnessy et al. 2020) to derive visual explanations, but do not lend insights back to model training. Our work represents a fast, principled solution that overcomes the above limitations. It automatically derives explainable masks faithful to the model and data, without explicit supervision from user-generated explanations.

Contrastive Learning. There has been growing interest in exploiting contrastive learning (CL) techniques for representations learning (Oord, Li, and Vinyals 2018; Chen et al. 2020; He et al. 2020; Khosla et al. 2020; Tian, Krishnan, and Isola 2019). Originally devised for density estimation (Gutmann and Hyvärinen 2010), CL exploits the idea of *learning by comparison* to capture the subtle features of data, *i.e.*, positive examples, by contrasting them with negative examples drawn from a carefully crafted noise distribution. These techniques aim to avoid representation collapse, or to promote representation consistency, for downstream tasks. Recent developments, both empirical and theoretical, have connected CL to information-theoretic foundations (Tian, Krishnan, and Isola 2019; Grill et al. 2020), thus establishing them as a suite of *de facto* solutions for unsupervised representation learning (Chen et al. 2020; He et al. 2020).

The basic form of CL is essentially a binary classification task specified to discriminate positive and negative examples. In such a scenario, the binary classifier is known as the critic function. Maximizing the discriminative power wrt the critic and the representation sharpens the feature encoder. Critical to the success of CL is the choice of appropriate noise distribution, where the challenging negatives, *i.e.*, those negatives that are more similar to positive examples, are often considered more effective contrasts. In its more generalized form, CL can naturally repurpose the predictor and loss functions without introducing a new critic (Tian, Krishnan, and Isola 2019). Notably, current CL methods are not immune to spurious associations, a point we wish to improve in this work.

Causality and Interventions. From a causality perspective, humans learn via actively interacting with the environment. We intervene and observe changes in the outcome to infer causal dependencies. Machines instead learn from static observations that are unable to inform the structural dependencies for causal decisions. As such, perturbations to the external factors, *e.g.*, surroundings, lighting, viewing angles, may drastically alter machine predictions, while human recognition is less susceptible to such nuisance variations. Such difference is best explained with the *do*-notation (Pearl 2009).

Unfortunately, carrying out real interventional studies, *i.e.*, randomized control trials, to intentionally block non-

causal associations, is oftentimes not a feasible option for practical considerations, *e.g.*, due to cost and ethics. This work instead advocates the application of synthetic interventions to uncover the underlying causal features from observational data. Specifically, we proactively edit x and its corresponding label y in a data-driven fashion to encourage the model to learn potential causal associations. Our proposal is in line with the growing appreciation for the significance of establishing causality in machine learning models (Schölkopf 2019). Via promoting invariance (Arjovsky et al. 2019), such causally inspired solutions demonstrate superior robustness to superficial features that do not generalize (Wang et al. 2019a). In particular, (Suter et al. 2019; Zhang, Zhang, and Li 2020) showed the importance and effectiveness of accounting for interventional perspectives. Our work brings these causal views to construct a simple solution that explicitly optimizes visual interpretation and model robustness.

Proactive Pseudo-Intervention

Below, we describe the construction of *Proactive Pseudo-Intervention* (PPI), a causally-informed contrastive learning scheme that seeks to simultaneously improve the accuracy, robustness, generalization and interpretability of deep-learning-based computer vision models.

The PPI learning strategy, schematically summarized in Figure 2, consists of three main components: (*i*) a saliency mapping module that highlights causally relevant features; (*ii*) an intervention module that synthesizes contrastive samples; and (*iii*) the prediction module, which is standard in recent vision models, *e.g.*, VGG (Simonyan and Zisserman 2014), ResNet (He et al. 2016), and Inception Net (Szegedy et al. 2016). Motivated by the discussions from our introduction, PPI establishes a feedback loop between the saliency map module and the prediction module, which is interfaced by the synthesized contrastive examples in the intervention module. Under this configuration, the prediction module is encouraged to modify its predictions only when provided with causally-relevant synthetic interventions. Note that components (*i*) and (*ii*) do not involve any additional parameters or neural network modules, which makes our strategy readily applicable to the training of virtually any computer vision task without major customization. Details of these building blocks are given below.

Synthetic causal interventions for contrasts

Key to our formulation is the design of a synthetic intervention strategy that generates contrastive examples to reinforce causal relevance during model training. Given a causal saliency map $s_m(x)$ for an input x wrt label $y = m$, where $m = 1, \dots, M$, and M is the number of classes, the synthetic intervention consists of removing (replacing with zero) the causal information from x contained in $s_m(x)$, and then using it as the contrastive learning signal.

For now, let us assume the causal salience map $s_m(x)$ is known; the procedure to obtain the saliency map will be addressed in the next section. For notational clarity, we use subscript i to denote entities associated with the i -th training

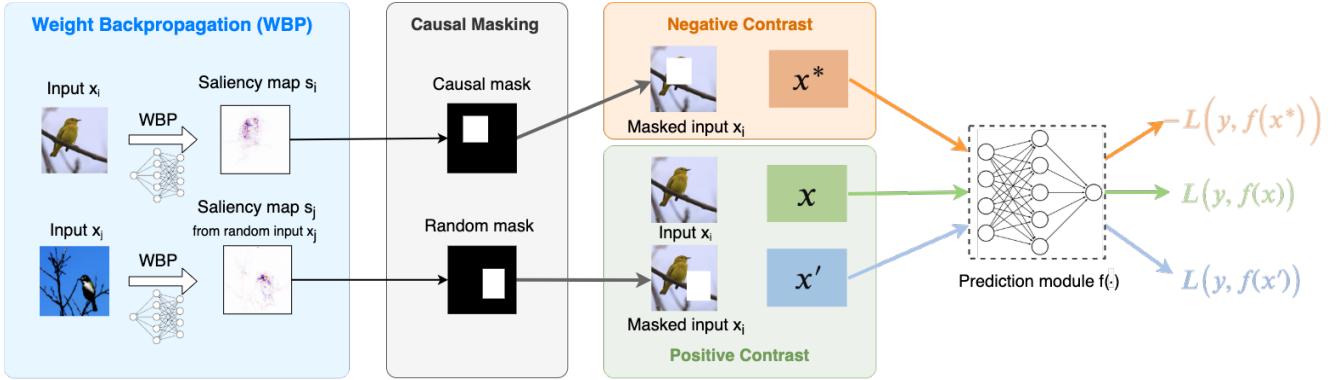


Figure 2: Illustration of the proposed PPI learning strategy. Input images are intervened by removing the saliency map based masks, which alters the input label (e.g., negative control). For positive contrast, we use the original input as well as an input masked with a random saliency map. We use WBP for the generation of saliency maps.

sample, and omit the dependency on learnable parameters. To remove causal information from \mathbf{x}_i and obtain a negative contrast \mathbf{x}_i^* , we apply the following *soft-masking*

$$\mathbf{x}_i^* = \mathbf{x}_i - T(s_m(\mathbf{x}_i)) \odot \mathbf{x}_i, \quad (1)$$

where $T(\cdot)$ is a differentiable masking function and \odot denotes element-wise (Hadamard) multiplication. Specifically, we use the thresholded sigmoid for masking:

$$T(s_m(\mathbf{x}_i)) = \frac{1}{1 + \exp(-\omega(s_m(\mathbf{x}_i) - \sigma))}, \quad (2)$$

where σ and $\omega > 0$ are the threshold and scaling parameters, respectively. We set the scaling ω so that $T(s)$ will result in a sharp transition from 0 to 1 near σ . Using equation 1 we define the contrastive loss as

$$L_{con}(\theta) = \sum_i \ell(\mathbf{x}_i^*, \neg y; f_\theta), \quad (3)$$

where f_θ is the prediction module, $\ell(\mathbf{x}, y; f_\theta)$ is the loss function we wish to optimize, e.g. cross entropy, and \neg is used to denote that the original class label has been flipped. In the binary case, $\neg y = 1 - y$, and in the multi-class case it can be interpreted accordingly, e.g., using a one vs. others cross-entropy loss. In practice, we set $\ell(\mathbf{x}, y; f_\theta) = -\ell(\mathbf{x}, y; f_\theta)$. We will show in the experiments that this simple and intuitive causal masking strategy works well in practice (see Tables 2 and 4, and Figure 5). Alternatively, we also consider a *hard-masking* approach in which a minimal bounding box covering the thresholded saliency map is removed. See the Appendix for details.

Note that we are making the implicit assumption that the saliency map is uniquely determined by the prediction module f_θ . While optimizing equation 3 explicitly attempts to improve the fit of the prediction module f_θ , it also implicitly informs the causal saliency mapping. This is sensible because if a prediction is made using non-causal features, which implies the associated saliency map $s_m(\mathbf{x})$ is also non-causal, then we should expect that after applying $s_m(\mathbf{x})$ to \mathbf{x} using equation 1, we can still expect to make the correct prediction, i.e., the true label, for both positive (the original) and negative (the intervened) samples.

Saliency map regularization. Note that naively optimizing equation 3 can lead to degenerate solutions for which any saliency map that satisfies the causal sufficiency, i.e., encompassing all causal features, is a valid causal saliency map. For example, a trivial solution where the saliency map covers the entire image may be considered causal. To protect against such degeneracy, we propose to regularize the L_1 -norm of the saliency map to encourage succinct (sparse) representations, i.e., $L_{reg} = \|s_m\|_1$, for $m = 1, \dots, M$.

Adversarial positive contrasts. Another concern with solely optimizing equation 3 is that models can easily overfit to the intervention, i.e., instead of learning to capture causal relevance, the model learns to predict interventional operations. For example, the model can learn to change its prediction when it detects that the input has been intervened, regardless of whether the image is missing causal features. So motivated, we introduce adversarial positive contrasts:

$$\mathbf{x}'_i = \mathbf{x}_i - T(s_m(\mathbf{x}_j)) \odot \mathbf{x}_i, \quad i \neq j, \quad (4)$$

where we intervene with a *false* saliency map, i.e., $s_m(\mathbf{x}_j)$ is the saliency map from a different input \mathbf{x}_j , while still encouraging the model to make the correct prediction via

$$L_{ad}(\theta) = \sum_i \ell(\mathbf{x}'_i, y; f_\theta), \quad (5)$$

where \mathbf{x}'_i is the adversarial positive contrast. The complete loss for the proposed model, $L = L_{cls} + L_{con} + L_{reg} + L_{ad}$, consists of the contrastive loss in equation 3, the regularization loss, L_{reg} , and the adversarial loss in equation 5.

Saliency Weight Backpropagation

PPI requires a module to generate saliency maps that inform decision-driving features in the (raw) pixel space. This module is expected to: i) generate high-quality saliency maps that faithfully reflect the model's focus, ii) efficient, as it will be used repeatedly in the PPI training framework. There is a tension between these two goals: advanced saliency map methods are usually time-consuming (Smilkov et al. 2017; Sundararajan, Taly, and Yan 2017). PPI finds a sweet point that better balances the trade-offs: the *Weight Back Propagation* (WBP). WBP is a novel computationally efficient

scheme for saliency mapping, and it applies to arbitrary neural architectures. Heuristically, WBP evaluates individual contributions from each pixel to the final class-specific prediction. Empirical examinations reveal that WBP results are more causally-relevant relative to competing solutions based on human judgment (refer to Figure 3 visualization).

To simplify our presentation, we first consider a vector input and a linear mapping. Let \mathbf{v}^l be the internal representation of the data at the l -th layer, with $l = 0$ being the input layer, *i.e.*, $\mathbf{v}^0 = \mathbf{v}$, and $l = L$ being the penultimate *logit* layer prior to the softmax transformation, *i.e.*, $\mathbb{P}(y|\mathbf{v}) = \text{softmax}(\mathbf{v}^L)$. To assign the relative importance to each hidden unit in the l -th layer, we notationally collapse all transformations after l into an operator denoted by $\tilde{\mathbf{W}}^l$, which we call the *saliency matrix*, satisfying,

$$\mathbf{v}^L = \tilde{\mathbf{W}}^l \mathbf{v}^l, \quad \forall l \in [0, \dots, L], \quad (6)$$

where \mathbf{v}^L is an M -dimensional vector corresponding to the M distinct classes in y . Though presented in a matrix form in a slight abuse of notation, *i.e.*, the instantiation of the operator $\tilde{\mathbf{W}}^l$ effectively depends on the input \mathbf{v} , thus all nonlinearities have been effectively absorbed into it. We posit that for an object associated with a given label $y = m$, its causal features are subsumed in the interactions between the m -th row of \mathbf{W}^0 and input \mathbf{v} , *i.e.*,

$$[\mathbf{s}_m(\mathbf{v})]_k = [\tilde{\mathbf{W}}^0]_{mk} [\mathbf{v}]_k, \quad (7)$$

where $[\mathbf{s}_m(\mathbf{v})]_k$ denotes the k -th element of the saliency map $\mathbf{s}_m(\mathbf{v})$ and $[\tilde{\mathbf{W}}^0]_{mk}$ is a single element of $\tilde{\mathbf{W}}^0$. A key observation for computation of $\tilde{\mathbf{W}}^l$ is that it can be done recursively. Specifically, let $g_l(\mathbf{v}^l)$ be the transformation at the l -th layer, *e.g.*, an affine transformation, convolution, activation, normalization, *etc.*, then it holds that

$$\tilde{\mathbf{W}}^{l+1} \mathbf{v}^{l+1} = \tilde{\mathbf{W}}^{l+1} g_l(\mathbf{v}^l) = \tilde{\mathbf{W}}^l \mathbf{v}^l. \quad (8)$$

This allows for recursive computation of $\tilde{\mathbf{W}}^l$ via

$$\tilde{\mathbf{W}}^l = G(\tilde{\mathbf{W}}^{l+1}, g_l), \quad \tilde{\mathbf{W}}^L = 1, \quad (9)$$

where $G(\cdot)$ is the update rule. We list the update rules for common transformations in deep networks in Table 1, with corresponding derivations detailed below.

Table 1: WBP update rules for common transformations.

Transformation	$G(\cdot)$
Activation Layer	$\tilde{\mathbf{W}}^l = h \circ \tilde{\mathbf{W}}^{l+1}$
FC Layer	$\tilde{\mathbf{W}}^l = \tilde{\mathbf{W}}^{l+1} \mathbf{W}^l$
Convolutional Layer	$\tilde{\mathbf{W}}^l = \tilde{\mathbf{W}}^{l+1} \otimes [\mathbf{W}^l]_{\text{flip}_{2,3}}^{T_{0,1}}$
BN Layer	$\tilde{\mathbf{W}}^l = \frac{\tilde{\mathbf{W}}^{l+1}}{\sigma} \gamma$
Pooling Layer	Relocate/Distribute $\tilde{\mathbf{W}}^{l+1}$

Fully-connected (FC) layer. The FC transformation is the most basic operation in deep neural networks. Below we omit the bias term as it does not directly interact with the input. Assuming $g_l(\mathbf{v}^l) = \mathbf{W}^l \mathbf{v}^l$, it is readily seen that

$$\tilde{\mathbf{W}}^{l+1} \mathbf{v}^{l+1} = \tilde{\mathbf{W}}^{l+1} g_l(\mathbf{v}^l) = (\tilde{\mathbf{W}}^{l+1} \mathbf{W}^l) \mathbf{v}^l, \quad (10)$$

so $\tilde{\mathbf{W}}^l = \tilde{\mathbf{W}}^{l+1} \mathbf{W}^l$. Graphical illustration with standard affine mapping and ReLU activation can be found in the Appendix.

Nonlinear activation layer. Considering that an activation layer simply *rescales* the saliency weight matrices, *i.e.*, $\mathbf{v}^{l+1} = g_l(\mathbf{v}^l) = h^l \circ \mathbf{v}^l$, where \circ is the composition operator, we obtain $\tilde{\mathbf{W}}^l = h \circ \tilde{\mathbf{W}}^{l+1}$. Using the ReLU activation as a concrete example, we have $h(\mathbf{v}^l) = 1\{\mathbf{v}^l \geq 0\}$.

Convolutional layer. The convolution is a generalized form of linear mapping. In practice, convolutions can be expressed as tensor products of the form $\tilde{\mathbf{W}}^l = \tilde{\mathbf{W}}^{l+1} \otimes [\mathbf{W}^l]_{\text{flip}_{2,3}}^{T_{0,1}}$, where $\mathbf{W}^l \in \mathbb{R}^{D_2 \times D_1 \times (2S+1) \times (2S+1)}$ is the convolution kernel, $T_{0,1}$ is the transpose in dimensions 0 and 1 and $\text{flip}_{2,3}$ is an exchange in dimensions 2 and 3. See the Appendix for details.

Pooling and normalization layer. Summarization and standardization are two other essential operations for the success of deep neural networks, achieved by pooling and batch normalization (BN) techniques, respectively. They too can be considered as special instantiations of linear operations. Here we summarize the two most popular operations in Table 1.

Experiments

To validate the utility of our approach, we consider both natural and medical image datasets, and compare it to existing state-of-the-art solutions. All the experiments are implemented in PyTorch. The source code will be available at https://github.com/author_name/PPI. Due to space limitation, details of the experimental setup and additional analyses are deferred to the Appendix.

Datasets. We present our findings on five representative datasets: (*i*) CIFAR-10 (Krizhevsky 2009); (*ii*) ImageNet (ILSVRC2012) (Russakovsky et al. 2015); (*iii*) CUB (Wah et al. 2011), a natural image dataset with over 12k photos for classification of 200 bird species in the wild, heavily confounded by the background characteristics; (*iv*) GA (Leuschen et al. 2013), a new medical image dataset for the prediction of *geographic atrophy* (GA) using 3D *optical coherence tomography* (OCT) image volumes, characterized by small sample size (275 subjects) and highly heterogeneous (collected from 4 different facilities); and (*v*) LiDC-IDRI (Langlotz et al. 2019), a public medical dataset of 1,085 lung lesion CT images annotated by 4 radiologists. Detailed specifications are described in the Appendix.

Baselines. We compare model trained with and without PPI framework to show gains on classification, model interpretability, and cross-domain generalization. Meanwhile, we compare our proposed WBP to the following set of popular saliency mapping schemes: (*i*) Gradient: standard gradient-based salience mapping; (*ii*) GradCAM (Selvaraju et al. 2017): gradient-weighted class activation mapping; (*iii*) LRP (Bach et al. 2015): layer-wise relevance propagation and its variants. We do no consider more advanced saliency mapping schemes, like perturbation based methods, because it is too time consuming to be used for training purposes.

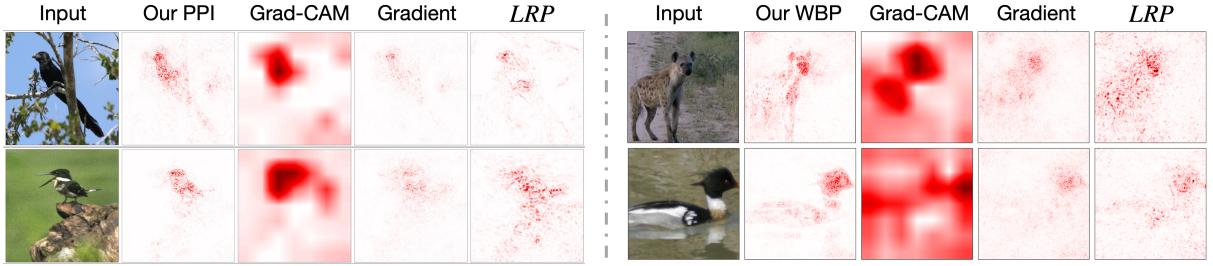


Figure 3: Visualization of the inferred saliency maps. Left: CUB dataset (PPI is based on WBP). Right: ImageNet dataset.

Hyperparameters. The final loss of the proposed model is a weighted summation of four losses: $L = L_{cls} + w_1 L_{con} + w_2 L_{reg} + w_3 L_{ad}$. The weights are simply balanced to match the magnitude of L_{cls} , i.e., $w_3 = 1$, $w_2 = 0.1$, and $w_1 = 1$ (CUB, Cifar-10, and GA) and $= 10$ (LIDC). See Appendix Sec B for details about the masking parameters σ and ω .

Natural Image Datasets

Classification Gains. In this experiment, we investigate how the different pairings of PPI and saliency mapping schemes (*i.e.*, GradCAM, LRP, WBP) affect performance. In Table 2, the first row represents VGG11 model trained with only classification loss, and the following rows represent VGG11 trained with PPI with different saliency mapping schemes. We see consistent performance gains in accuracy via incorporating PPI training on both CUB and CIFAR-10 datasets. The gains are mostly significant when using our WBP for saliency mapping (improving the accuracy from 0.662 to 0.696 on CUB, and from 0.881 to 0.901 on CIFAR-10).

Table 2: Performance improvements achieved by training with PPI on CUB, CIFAR-10, and GA dataset. We report means and standard deviations (SDs) from 5-fold cross-validation for GA prediction.

Models	CUB (Acc)	Cifar-10 (Acc)	GA (AUC \pm SD)
Classification	0.662	0.881	0.877 ± 0.040
+PPI _{Gradient}	0.673	0.885	0.890 ± 0.035
+PPI _{LRP}	0.680	0.891	0.895 ± 0.037
+PPI _{GradCAM}	0.683	0.895	0.908 ± 0.036
+PPI_{WBP}	0.696	0.901	0.925 ± 0.023

Model Interpretability. In this task, we want to qualitatively and quantitatively compare the causal relevance of saliency maps generated by our proposed model and its competitors. In Figure 3(right), we show the saliency maps produced by different approaches for a VGG11 model trained on CUB. Visually, gradient-based solutions (Grad and Grad-CAM) tend to yield overly dispersed maps, indicating a lack of specificity. LRP gives more appealing saliency maps. However, these maps also heavily attend to the spurious background cues that presumably help with predictions. When trained with PPI (based on WBP), the saliency maps focus on the causal related pixels, i.e., special parts of birds

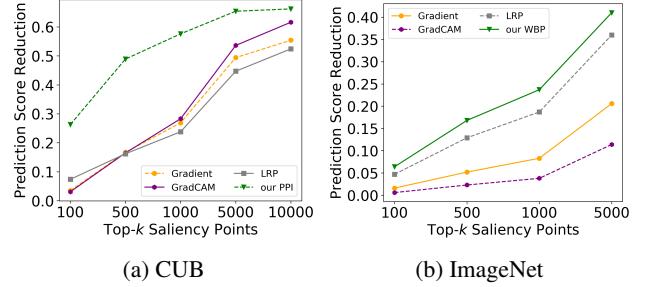


Figure 4: Quantitative evaluations of causal relevance of competing saliency maps (higher is better).

body.

To quantitatively evaluate the causal relevance of competing saliency maps, we adopt the evaluation scheme proposed in (Hooker et al. 2019), consisting of masking out the contributing saliency pixels and then calculating the reduction in prediction score. A larger reduction is considered better for accurately capturing the pixels that ‘cause’ the prediction. Results are summarized in Figure 4a, where we progressively remove the top- k saliency points, with $k = 100, 500, 1000, 5000, 10000$ ($10000 \approx 6.6\%$ of all pixels), from the CUB test input images. Our PPI consistently outperforms its counterparts, with its lead being most substantial in the low- k regime. Notably, for large k , PPI removes nearly all predictive signal. This implies PPI specifically targets the causal features. Quantitative evaluation with additional metrics are provided in the Appendix.

To test the performance of WBP itself (without being trained with PPI), we compare WBP with different approaches for a VGG11 model trained on ImageNet from PyTorch model zoo. Figure 3(left) shows that saliency maps generated by WBP more concentrate on objects themselves. Also, thanks to the fine resolution of WBP, the model pays more attention to the patterns on the fur to identify the leopard (row 1). This is more visually consistent with human judgement. Figure 4b demonstrates WBP identifies more causal pixels on ImageNet validation images.

OCT-GA: Geographic Atrophy Classification

Next we show how the proposed PPI handles the challenges of small training data and heterogeneity in medical image datasets. In this experiment (with our new dataset, that we will make public), each OCT volume image consists of 100 scans of a 512×1000 sized image (Boyer et al. 2017). We use a multi-view CNN model (Su et al. 2015) to process

such 3D OCT inputs, and use it as our baseline solution (see the Appendix). We investigate how the different saliency mapping schemes (*i.e.*, Grad, GradCAM, LRP, WBP) work with PPI. For WBP, we also tested the bounding box variant, denoted as WBP (box) (see the Appendix). In Table 2, we see consistent performance gains in AUC score via incorporating PPI training (from 0.877 to 0.925, can be improve to 0.937 by PPI with WBP(box)), accompanied by the reductions in model variation evaluated by the standard deviations of AUC from the five-fold cross-validation. The gains are most significant when using our WBP for saliency mapping. The model trained by PPI with WBP is significantly different with other baseline models based on the DeLong test for receiving operating characteristic comparisons (DeLong, DeLong, and Clarke-Pearson 1988). We further compare the saliency maps generated by these different combinations. We see that without the additional supervision from PPI, competing solutions like Grad, GradCAM and LRP sometimes yield non-sensible saliency maps (attending to image corners). Overall, PPI encourages more concentrated and less noisy saliency maps. Also, different PPI-based saliency maps agree with each other to a larger extent. Our findings are also verified by experts (co-authors, who are ophthalmologists specializing in GA) confirming that the PPI-based saliency maps are clinically relevant by focusing on retinal layers likely to contain abnormalities or lesions. These results underscore the practical value of the proposed proactive interventions.

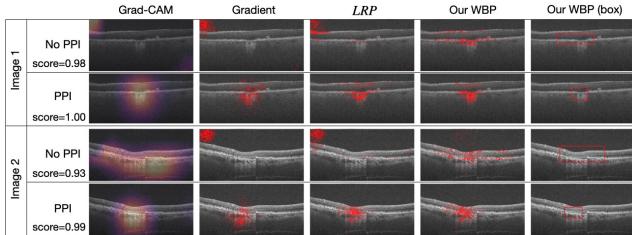


Figure 5: Saliency maps on GA dataset based on models trained with PPI and without PPI. Maps of models trained with PPI are more clinically relevant by focusing on retinal layers likely to contain abnormalities or lesions, and more concentrated.

Cross-domain generalization. Common to medical image applications is that training samples are usually integrated from a number of healthcare facilities (*i.e.*, domains), and that predictions are sometimes to be made on subjects at other facilities. Despite big efforts to standardize the image collection protocols, with different imaging systems operated by technicians with varying skills, apparent domain shifts are likely to compromise the cross-domain performance of these models. We show this phenomenon on the GA dataset in Table 3, where source samples are collected from four different hospitals in different health systems (A, B, C and D, see the Appendix for details). Each cell contains the AUC of the model trained on site X (row) and tested on site Y (column), with same-site predictions made on hold-out samples. A significant performance drop is observed for

cross-domain predictions (off-diagonals) compared to in-domain predictions (diagonals). With the application of PPI, the performance gaps between in-domain and cross-domain predictions are considerably reduced. The overall accuracy gains of PPI further justify the utility of causally-inspired modeling. Notably, site D manifests strong spurious correlation that help in-domain prediction but degrades out-of-site generalization, which is partly resolved by the proposed PPI.

Table 3: AUC results for GA prediction with or without PPI. Models are trained on one site and cross-validated on the other sites. Darker color indicates better performance.

	With PPI	A	B	C	D	Mean	STD
A	1.000	0.906	0.877	0.865	0.912	0.061	
B	0.851	0.975	0.863	0.910	0.900	0.056	
C	0.954	0.875	0.904	0.931	0.916	0.034	
D	0.824	0.846	0.853	0.904	0.857	0.034	
No PPI	A	B	C	D	Mean	STD	
A	1.000	0.854	0.832	0.827	0.878	0.082	
B	0.810	0.874	0.850	0.906	0.860	0.040	
C	0.860	0.779	0.873	0.862	0.843	0.043	
D	0.748	0.792	0.836	0.961	0.834	0.092	

LIDC-IDRI: Lung Lesions Classification

To further examine the practical advantages of the proposed PPI in real-world applications, we benchmark its utility on LIDC-IDRI; a public lung CT scan dataset (Armato III et al. 2011). We followed the preprocessing steps outlined in (Kohl et al. 2018) to prepare the data, and adopted the experimental setup from (Selvan and Dam 2020) to predict lesions. We use Inception_v3 (Szegedy et al. 2016) as our base model for both standard classification and PPI-enhanced training with various saliency mapping schemes. See the Appendix for details.

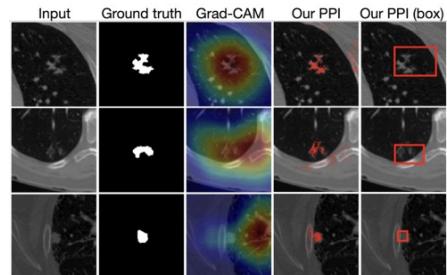


Figure 6: Saliency maps on LIDC-IDRI. Saliency maps of PPI+WBP are mostly consistent with the ground truths.

Lesion classification. We first compare PPI to other specialized SOTA network architectures. Table 4 summarizes AUC scores of Tensor Net-X (Efthymiou, Hidary, and Leichenauer 2019), DenseNet (Huang et al. 2017), LoTeNet (Selvan and Dam 2020), Inception_v3 (Szegedy et al. 2016), as well as our Inception_v3 trained with and without PPI_{WBP} . The proposed $\text{PPI}_{WBP}^{(box)}$ leads the performance chart by a considerable margin, improving Inception_v3 from 0.92 to 0.94.

Table 4: LIDC-IDRI classification AUC results.

Models	AUC
Tensor Net-X (Efthymiou, Hidary, and Leichenauer 2019)	0.823
DenseNet (Huang et al. 2017)	0.829
LoTeNet (Selvan and Dam 2020)	0.874
Inception_v3 (Szegedy et al. 2016)	0.921
+PPI _{GradCAM}	0.933
+PPI _{Gradient}	0.930
+PPI _{LRP}	0.931
+PPI _{WBP}	0.935
+PPI _{WBP(box)}	0.941

Weakly-supervised image segmentation. In Figure 6, we compare saliency maps generated by GradCAM, WBP, WBP (box) to the ground truth lesion masks from expert annotations. Note that we have only supplied patch-label labels during training, not the pixel-level expert segmentation masks, which constitute a challenging task of weakly-supervised image segmentation. In line with the observations from the GA experiment, our PPI-training enhanced WBP saliency maps are mostly consistent with the expert segmentations. Together with Table 4, Figure 6 confirms that the proposed PPI+WBP improves both the classification performance and model interpretability.

Conclusions

We have presented *Proactive Pseudo-Intervention* (PPI), a novel interpretable computer vision framework that organically integrates saliency mapping, causal reasoning, synthetic intervention and contrastive learning. PPI couples saliency mapping with contrastive training by creating artificially intervened negative samples absent of causal features. To communicate model insights and facilitate pre-informed reasoning, we derived an architecture-agnostic saliency mapping scheme called *Weight Back Propagation* (WBP), which faithfully captures the causally-relevant pixels/features for model prediction. Visual inspections of the saliency maps show that WBP, is more robust to spurious features compared to competing approaches. Empirical results on natural and medical datasets verify the combination of PPI and WBP consistently delivers performance boosts across a wide range of tasks relative to competing solutions, and the gains are most significant where the application is complicated by small sample size, data heterogeneity, or confounded with spurious correlations.

References

- Adebayo, J.; Gilmer, J.; Muelly, M.; Goodfellow, I.; Hardt, M.; and Kim, B. 2018. Sanity checks for saliency maps. In *Advances in Neural Information Processing Systems*, 9505–9515.
- Arjovsky, M.; Bottou, L.; Gulrajani, I.; and Lopez-Paz, D. 2019. Invariant risk minimization. *arXiv preprint arXiv:1907.02893*.
- Armato III, S. G.; McLennan, G.; Bidaut, L.; McNitt-Gray, M. F.; Meyer, C. R.; Reeves, A. P.; Zhao, B.; Aberle, D. R.; Henschke, C. I.; Hoffman, E. A.; et al. 2011. The lung image database consortium (LIDC) and image database resource initiative (IDRI): a completed reference database of lung nodules on CT scans. *Medical physics*, 38(2): 915–931.
- Bach, S.; Binder, A.; Montavon, G.; Klauschen, F.; Müller, K.-R.; and Samek, W. 2015. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PloS one*, 10(7): e0130140.
- Boyer, D. S.; Schmidt-Erfurth, U.; van Lookeren Campagne, M.; Henry, E. C.; and Brittain, C. 2017. The pathophysiology of geographic atrophy secondary to age-related macular degeneration and the complement pathway as a therapeutic target. *Retina (Philadelphia, Pa.)*, 37(5): 819.
- Chang, C.-H.; Creager, E.; Goldenberg, A.; and Duvenaud, D. 2018. Explaining Image Classifiers by Counterfactual Generation. In *International Conference on Learning Representations*.
- Chattopadhyay, A.; Sarkar, A.; Howlader, P.; and Balasubramanian, V. N. 2018. Grad-cam++: Generalized gradient-based visual explanations for deep convolutional networks. In *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*, 839–847. IEEE.
- Chen, C.; Seff, A.; Kornhauser, A.; and Xiao, J. 2015. Deep-driving: Learning affordance for direct perception in autonomous driving. In *Proceedings of the IEEE International Conference on Computer Vision*, 2722–2730.
- Chen, T.; Kornblith, S.; Norouzi, M.; and Hinton, G. 2020. A simple framework for contrastive learning of visual representations. *arXiv preprint arXiv:2002.05709*.
- Dabkowski, P.; and Gal, Y. 2017. Real time image saliency for black box classifiers. In *Advances in Neural Information Processing Systems*, 6967–6976.
- DeLong, E. R.; DeLong, D. M.; and Clarke-Pearson, D. L. 1988. Comparing the areas under two or more correlated receiver operating characteristic curves: a nonparametric approach. *Biometrics*, 837–845.
- Deng, J.; Dong, W.; Socher, R.; Li, L.-J.; Li, K.; and Fei-Fei, L. 2009. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, 248–255. Ieee.
- Deng, L. 2012. The mnist database of handwritten digit images for machine learning research [best of the web]. *IEEE Signal Processing Magazine*, 29(6): 141–142.
- Dhurandhar, A.; Chen, P.-Y.; Luss, R.; Tu, C.-C.; Ting, P.; Shanmugam, K.; and Das, P. 2018. Explanations based on the missing: Towards contrastive explanations with pertinent negatives. In *Advances in Neural Information Processing Systems*, 592–603.
- Du, M.; Liu, N.; Song, Q.; and Hu, X. 2018. Towards explanation of dnn-based prediction with guided feature inversion. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 1358–1367.
- Efthymiou, S.; Hidary, J.; and Leichenauer, S. 2019. TensorNetwork for machine learning. *arXiv preprint arXiv:1906.06329*.

- Erhan, D.; Bengio, Y.; Courville, A.; and Vincent, P. 2009. Visualizing higher-layer features of a deep network. *University of Montreal*, 1341(3): 1.
- Fong, R.; Patrick, M.; and Vedaldi, A. 2019. Understanding deep networks via extremal perturbations and smooth masks. In *Proceedings of the IEEE International Conference on Computer Vision*, 2950–2958.
- Fong, R. C.; and Vedaldi, A. 2017. Interpretable explanations of black boxes by meaningful perturbation. In *Proceedings of the IEEE International Conference on Computer Vision*, 3429–3437.
- Fukui, H.; Hirakawa, T.; Yamashita, T.; and Fujiyoshi, H. 2019. Attention branch network: Learning of attention mechanism for visual explanation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 10705–10714.
- Girshick, R.; Donahue, J.; Darrell, T.; and Malik, J. 2014. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 580–587.
- Goyal, Y.; Wu, Z.; Ernst, J.; Batra, D.; Parikh, D.; and Lee, S. 2019. Counterfactual Visual Explanations. In *ICML*.
- Grill, J.-B.; Strub, F.; Altché, F.; Tallec, C.; Richemond, P.; Buchatskaya, E.; Doersch, C.; Avila Pires, B.; Guo, Z.; Gheshlaghi Azar, M.; et al. 2020. Bootstrap Your Own Latent-A New Approach to Self-Supervised Learning. *Advances in Neural Information Processing Systems*, 33.
- Gutmann, M.; and Hyvärinen, A. 2010. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, 297–304.
- He, K.; Fan, H.; Wu, Y.; Xie, S.; and Girshick, R. 2020. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 9729–9738.
- He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 770–778.
- Hooker, S.; Erhan, D.; Kindermans, P.-J.; and Kim, B. 2019. A benchmark for interpretability methods in deep neural networks. In *Advances in Neural Information Processing Systems*, 9737–9748.
- Huang, G.; Liu, Z.; Van Der Maaten, L.; and Weinberger, K. Q. 2017. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 4700–4708.
- Johansson, F.; Shalit, U.; and Sontag, D. 2016. Learning representations for counterfactual inference. In *International conference on machine learning*, 3020–3029.
- Khosla, P.; Teterwak, P.; Wang, C.; Sarna, A.; Tian, Y.; Isola, P.; Maschinot, A.; Liu, C.; and Krishnan, D. 2020. Supervised contrastive learning. *arXiv preprint arXiv:2004.11362*.
- Kohl, S.; Romera-Paredes, B.; Meyer, C.; De Fauw, J.; Ledam, J. R.; Maier-Hein, K.; Eslami, S. A.; Rezende, D. J.; and Ronneberger, O. 2018. A probabilistic u-net for segmentation of ambiguous images. In *Advances in Neural Information Processing Systems*, 6965–6975.
- Krizhevsky, A. 2009. Learning Multiple Layers of Features from Tiny Images. *Master's thesis, University of Tront*.
- Krizhevsky, A.; Sutskever, I.; and Hinton, G. E. 2017. Imagenet classification with deep convolutional neural networks. *Communications of the ACM*, 60(6): 84–90.
- Langlotz, C. P.; Allen, B.; Erickson, B. J.; Kalpathy-Cramer, J.; Bigelow, K.; Cook, T. S.; Flanders, A. E.; Lungren, M. P.; Mendelson, D. S.; Rudie, J. D.; et al. 2019. A roadmap for foundational research on artificial intelligence in medical imaging: from the 2018 NIH/RSNA/ACR/The Academy Workshop. *Radiology*, 291(3): 781–791.
- Leuschen, J. N.; Schuman, S. G.; Winter, K. P.; McCall, M. N.; Wong, W. T.; Chew, E. Y.; Hwang, T.; Srivastava, S.; Sarin, N.; Clemons, T.; et al. 2013. Spectral-domain optical coherence tomography characteristics of intermediate age-related macular degeneration. *Ophthalmology*, 120(1): 140–150.
- Li, K.; Wu, Z.; Peng, K.-C.; Ernst, J.; and Fu, Y. 2018. Tell me where to look: Guided attention inference network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 9215–9223.
- Mahendran, A.; and Vedaldi, A. 2016. Salient deconvolutional networks. In *European Conference on Computer Vision*, 120–135. Springer.
- Montavon, G. 2019. Gradient-based vs. propagation-based explanations: an axiomatic comparison. In *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning*, 253–265. Springer.
- Montavon, G.; Binder, A.; Lapuschkin, S.; Samek, W.; and Müller, K.-R. 2019. Layer-wise relevance propagation: an overview. In *Explainable AI: interpreting, explaining and visualizing deep learning*, 193–209. Springer.
- Oord, A. v. d.; Li, Y.; and Vinyals, O. 2018. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*.
- O'Shaughnessy, M.; Canal, G.; Connor, M.; Rozell, C.; and Davenport, M. 2020. Generative causal explanations of black-box classifiers. *Advances in Neural Information Processing Systems*, 33.
- Pearl, J. 2009. *Causality*. Cambridge university press.
- Petsiuk, V.; Das, A.; and Saenko, K. 2018. Rise: Randomized input sampling for explanation of black-box models. *arXiv preprint arXiv:1806.07421*.
- Rebuffi, S.-A.; Fong, R.; Ji, X.; and Vedaldi, A. 2020. There and Back Again: Revisiting Backpropagation Saliency Methods. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 8839–8848.
- Ribeiro, M. T.; Singh, S.; and Guestrin, C. 2016. "Why should I trust you?" Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, 1135–1144.

- Rieger, L.; Singh, C.; Murdoch, W.; and Yu, B. 2020. Interpretations are useful: penalizing explanations to align neural networks with prior knowledge. In *International Conference on Machine Learning*, 8116–8126. PMLR.
- Ross, A. S.; Hughes, M. C.; and Doshi-Velez, F. 2017. Right for the right reasons: training differentiable models by constraining their explanations. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, 2662–2670.
- Russakovsky, O.; Deng, J.; Su, H.; Krause, J.; Satheesh, S.; Ma, S.; Huang, Z.; Karpathy, A.; Khosla, A.; Bernstein, M.; Berg, A. C.; and Fei-Fei, L. 2015. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3): 211–252.
- Sajda, P. 2006. Machine learning for detection and diagnosis of disease. *Annu. Rev. Biomed. Eng.*, 8: 537–565.
- Schölkopf, B. 2019. Causality for machine learning. *arXiv preprint arXiv:1911.10500*.
- Selvan, R.; and Dam, E. B. 2020. Tensor Networks for Medical Image Classification. In *Medical Imaging with Deep Learning*.
- Selvaraju, R. R.; Cogswell, M.; Das, A.; Vedantam, R.; Parikh, D.; and Batra, D. 2017. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE International Conference on Computer Vision*, 618–626.
- Seo, D.; Oh, K.; and Oh, I.-S. 2019. Regional multi-scale approach for visually pleasing explanations of deep neural networks. *IEEE Access*, 8: 8572–8582.
- Shrikumar, A.; Greenside, P.; and Kundaje, A. 2017. Learning Important Features Through Propagating Activation Differences. In *International Conference on Machine Learning*, 3145–3153.
- Simonyan, K.; Vedaldi, A.; and Zisserman, A. 2013. Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034*.
- Simonyan, K.; and Zisserman, A. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*.
- Smilkov, D.; Thorat, N.; Kim, B.; Viégas, F.; and Wattenberg, M. 2017. Smoothgrad: removing noise by adding noise. *arXiv preprint arXiv:1706.03825*.
- Su, H.; Maji, S.; Kalogerakis, E.; and Learned-Miller, E. 2015. Multi-view convolutional neural networks for 3d shape recognition. In *Proceedings of the IEEE international conference on computer vision*, 945–953.
- Sundararajan, M.; Taly, A.; and Yan, Q. 2017. Axiomatic attribution for deep networks. In *International Conference on Machine Learning*, 3319–3328. PMLR.
- Suter, R.; Miladinovic, D.; Schölkopf, B.; and Bauer, S. 2019. Robustly disentangled causal mechanisms: Validating deep representations for interventional robustness. In *International Conference on Machine Learning*, 6056–6065. PMLR.
- Szegedy, C.; Vanhoucke, V.; Ioffe, S.; Shlens, J.; and Wojna, Z. 2016. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2818–2826.
- Tian, Y.; Krishnan, D.; and Isola, P. 2019. Contrastive multiview coding. *arXiv preprint arXiv:1906.05849*.
- Wagner, J.; Kohler, J. M.; Gindele, T.; Hetzel, L.; Wiedemer, J. T.; and Behnke, S. 2019. Interpretable and fine-grained visual explanations for convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 9097–9107.
- Wah, C.; Branson, S.; Welinder, P.; Perona, P.; and Belongie, S. 2011. The Caltech-UCSD Birds-200-2011 Dataset. Technical Report CNS-TR-2011-001, California Institute of Technology.
- Wang, H.; He, Z.; Lipton, Z. C.; and Xing, E. P. 2019a. Learning robust representations by projecting superficial statistics out. In *International Conference on Learning Representations*.
- Wang, H.; Wang, Z.; Du, M.; Yang, F.; Zhang, Z.; Ding, S.; Mardziel, P.; and Hu, X. 2020a. Score-CAM: Score-weighted visual explanations for convolutional neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 24–25.
- Wang, L.; Wu, Z.; Karanam, S.; Peng, K.-C.; Singh, R. V.; Liu, B.; and Metaxas, D. N. 2019b. Sharpen focus: Learning with attention separability and consistency. In *Proceedings of the IEEE International Conference on Computer Vision*, 512–521.
- Wang, T.; Huang, J.; Zhang, H.; and Sun, Q. 2020b. Visual Commonsense Representation Learning via Causal Inference. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 378–379.
- Zeiler, M. D.; and Fergus, R. 2014. Visualizing and understanding convolutional networks. In *European conference on computer vision*, 818–833. Springer.
- Zhang, C.; Zhang, K.; and Li, Y. 2020. A Causal View on Robustness of Neural Networks. *arXiv preprint arXiv:2005.01095*.
- Zhao, Y.; Zheng, Y.; Zhao, Y.; Liu, Y.; Chen, Z.; Liu, P.; and Liu, J. 2018. Uniqueness-driven saliency analysis for automated lesion detection with applications to retinal diseases. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 109–118. Springer.
- Zhou, B.; Khosla, A.; Lapedriza, A.; Oliva, A.; and Torralba, A. 2016. Learning deep features for discriminative localization. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2921–2929.

Appendix

Derivation of Convolutional Weight Backpropagation

Let's denote the input variable as $\mathbf{I} \in \mathbb{R}^{H \times W}$, the convolutional filter weight as $\mathbf{W} \in \mathbb{R}^{(2S+1) \times (2S+1)}$, the output variable as $\mathbf{O} \in \mathbb{R}^{H \times W}$, and the weight backpropagate to O as $\hat{\mathbf{W}} \in \mathbb{R}^{H \times W}$. We omit the bias here because it does not directly interact with the input variables. We denote \otimes as the convolutional operator. We have

$$\mathbf{O} = \mathbf{I} \otimes \mathbf{W} \quad (11)$$

$$\mathbf{O}_{i,j} = \sum_{i'=-S}^S \sum_{j'=-S}^S \mathbf{I}_{i+i', j+j'} \mathbf{W}_{i'+S, j'+S} \quad (12)$$

$$\sum_i \sum_j \mathbf{O}_{i,j} \tilde{\mathbf{W}}_{i,j} = \sum_i \sum_j \sum_{i'=-S}^S \sum_{j'=-S}^S \mathbf{I}_{i+i', j+j'} \mathbf{W}_{i'+S, j'+S} \tilde{\mathbf{W}}_{i,j} \quad (13)$$

$$\sum_i \sum_j \mathbf{O}_{i,j} \tilde{\mathbf{W}}_{i,j} = \sum_i \sum_j \mathbf{I}_{i,j} \sum_{i'=-S}^S \sum_{j'=-S}^S \tilde{\mathbf{W}}_{i+i', j+j'} \mathbf{W}_{-i'+S, -j'+S} \quad (14)$$

$$\sum_i \sum_j \mathbf{O}_{i,j} \tilde{\mathbf{W}}_{i,j} = \sum_i \sum_j \mathbf{I}_{i,j} (\tilde{\mathbf{W}} \otimes [\mathbf{W}]_{flip_{i,j}})_{i,j} \quad (15)$$

Hence the weight backpropagate through a convolutional layer is $\tilde{\mathbf{W}}^l = \tilde{\mathbf{W}}^{l+1} \otimes [\mathbf{W}^l]_{flip}$. For the 3D cases, $\mathbf{I}^l \in \mathbb{R}^{D_1 \times H \times W}$, the weight back propagates to \mathbf{O}^l is $\tilde{\mathbf{W}}^{l+1} \in \mathbb{R}^{D_2 \times H \times W}$ and the convolutional weight is $\mathbf{W}^l \in \mathbb{R}^{D_2 \times D_1 \times (2S+1) \times (2S+1)}$. To match the depth of $\tilde{\mathbf{W}}^{l+1}$, the \mathbf{W}^l is transposed in the first two dimensions. So $\tilde{\mathbf{W}}^l = \tilde{\mathbf{W}}^{l+1} \otimes [\mathbf{W}^l]_{flip_{2,3}}^{T_{0,1}}$. If the convolutional layer is downsizing the input variable (*i.e.*, strides), the $\tilde{\mathbf{W}}^{l+1}_{ijk}$ is padded with zeros around the weights (left, right, up, and down) to for the input elements that the convolutional filter strides over. The number of padding zeros is equal to the number of strides minus 1.

Details on Causal Masking

In this work, we consider three types of causal masking: (*i*) the point-wise soft causal masking defined by Equation (2) in the main text, (*ii*) hard masking, and (*iii*) box masking. For the hard masking, for each image, we keep points with WBP weight larger than k times of the standard deviation of WBP weights of the whole image. We test k from 1 to 7 and achieve similar results. As $k = 7$ performs slightly better, we set k as 7 for all experiments. For the box masking, we use the center of mass for these kept points as the center to draw a box. The height and width of this box is defined as $center_{h/w} \pm 1.2std_{h/w}$. In this way at least 90% of filtered points are contained in the box. For the soft masking, we set ω to 100 and σ to 0.25. We have also experimented

with image-adaptive thresholds instead of a fixed σ for all inputs, *i.e.*, set the threshold as mean value plus k times of the standard deviation of WBP weights of the whole image. We repeat the experiments a few times and the results are consistent. The experiment comparison of these masking methods mention above is conducted on LIDC dataset.

Table 5: AUC on LIDC from different causal masking methods

Models	AUC
WBP-soft (fixed σ)	0.931
WBP-soft (adaptive σ)	0.941
WBP-hard (point)	0.935
WBP-hard (box)	0.941

Related Work

In this work, we propose a contrastive causal representation learning strategy, *i.e.*, Proactive Pseudo-Intervention (PPI), that leverages proactive interventions to identify causally-relevant image features. This approach is complemented with a novel causal salience map visualization module, *i.e.*, Weight Back Propagation (WBP), that identifies important pixels in the raw input image, which greatly facilitates interpretability of predictions.

Prior related works will be discussed in this section. Compared with alternative post-hot saliency mapping methods, WBP outperforms these methods as both a standalone causal saliency map and a trainable model for model interpretation. Compared with other trainable interpretation models, the proposed PPI+WBP improves both model performance and model interpretations.

Post-hoc Saliency Maps

We compare WBP with other post-hoc saliency mapping methods to show why WBP is able to target the causal features, and generate more succinct and reliable saliency maps.

Perturbation Based Methods These methods make perturbations to individual inputs or neurons and monitor the impact on output neurons in the network. (Zeiler and Fergus 2014) occludes different segments of an input image and visualized the change in the activations of subsequent layers. Several methods follow a similar idea, but use other importance measures or occlusion strategies (Petsiuk, Das, and Saenko 2018; Ribeiro, Singh, and Guestrin 2016; Seo, Oh, and Oh 2019). More complicated works aim to generate an explanation by optimizing for a perturbed version of the image (Fong and Vedaldi 2017; Fong, Patrick, and Vedaldi 2019; Dabkowski and Gal 2017; Du et al. 2018). (Wagner et al. 2019) proposes a new adversarial defense technique which filters gradients during optimization to achieve fine-grained explanation. However, such perturbation based methods are computationally intensive and involve sophisticated model designs, which make it extremely hard to be integrated with other advance learning strategies.

Backpropagation Based Methods Backpropagation based methods (BBM) propagate an importance signal from an output neuron backwards through the layers to the input. These methods are usually fast to compute and produce fine-grained importance/relevancy maps. WBP is one of such method.

The pioneer methods in this category backpropagate a gradient to the image, and branches of studies extend this work by manipulating the gradient. These methods are discussed and compared in (Mahendran and Vedaldi 2016; Erhan et al. 2009). However, these maps are generally less concentrated (Dabkowski and Gal 2017; Fong and Vedaldi 2017) and less interpretable. Other BBMs such as Layer-wise Relevance Propagation (Bach et al. 2015), DeepLift (Shrikumar, Greenside, and Kundaje 2017) employ top-down relevancy propagation rules. DeepLift is sensitive to the reference inputs, which needs more human efforts and background knowledge to produce appealing saliency maps. The nature of depending on reference inputs limits its ability on model diagnosis and couple with learning strategies to continuously improving models’ performance. LRP decomposes the relevance, R , from a neuron, k , in the upper layer to every connected neurons, j , in the lower layer. The decomposition is distributed through gradients under the suggested implementation (Montavon et al. 2019). Our experiments on GA and CUB datasets show that vanilla LRP performs similar to gradient based methods, which is also demonstrated in (Montavon 2019). The variants of LRP use complex rules to prioritize positive or large relevance, making the saliency map visually appealing to human. However, our experiments demonstrate the unfaithfulness of LRP and its variants as they highlight spuriously correlated features (boarderlines and backgrounds). By contrast, our WBP backpropagates the the weights of through layers to compute the contributions of each input pixel, which is truly faithful to the model, and WBP tends to highlight the target objects themselves rather than the background. At the same time, the simplicity and efficiency makes WBP easily work with other advanced learning strategies for both model diagnosis and improvements during training.

Table 6: A list of commonly used LRP rules.(Montavon et al. 2019)

Rules	Formula
LRP	$R_j = \sum_k \frac{a_j w_{jk}}{\sum_{0,j} a_j w_{jk}} R_k$
LRP_ϵ	$R_j = \sum_k \frac{\epsilon + a_j w_{jk}}{\epsilon + \sum_{0,j} a_j w_{jk}} R_k$
LRP_γ	$R_j = \sum_k \frac{a_j (w_{jk} + \gamma w_{jk}^+)}{\sum_{0,j} a_j (w_{jk} + \gamma w_{jk}^+)} R_k$
$LRP_{\alpha\beta}$	$R_j = \sum_k (\alpha \frac{(a_j w_{jk})^+}{\sum_{0,j} (a_j w_{jk})^+} - \beta \frac{(a_j w_{jk})^-}{\sum_{0,j} (a_j w_{jk})^-}) R_k$
LRP_{flat}	$R_j = \sum_k \frac{1}{\sum_j 1} R_k$
LRP_{w^2}	$R_j = \sum_j \frac{w_{ij}^2}{\sum_i w_{ij}^2} R_j$
LRP_{Z^β}	$R_j = \sum_j \frac{x_i w_{ij} - l_i w_{ij}^+ - h_i w_{ij}^-}{\sum_i x_i w_{ij} - l_i w_{ij}^+ - h_i w_{ij}^-} R_j$

Activation Based Methods Methods under this category (such as CAM, Grad-CAM, guided Grad-CAM, Grad-CAM++) use a linear combination of class activation maps from convolutional layers to derive a saliency map. The main difference between them is how the linear combination weights are computed. The generation of saliency maps is easy and these methods can be coupled with advanced training strategies to improve training (Li et al. 2018). However, they fail at visualizing fine-grained evidence, which is particularly important in explaining medical classification models. Additionally, it is not guaranteed that the resulting explanations are faithful and reflect the decision making process of the model (Du et al. 2018; Selvaraju et al. 2017; Wagner et al. 2019). Grad-CAM++ (Chattopadhyay et al. 2018) proposes to introduce higher-order derivatives to capture pixel-level importance, while its high computational cost in calculating the second- and third-order derivatives makes it impractical for training purposes.

Interpretable Models

Unlike the *post hoc* saliency map generation described above, an alternative approach is to train a separate module to explicitly produce model explanations (Fukui et al. 2019; Goyal et al. 2019; Chang et al. 2018; Fong and Vedaldi 2017; Shrikumar, Greenside, and Kundaje 2017). Such *post hoc* causal explanations can be generated with black-box classifiers based on a learned low-dimensional representation of the data (O’Shaughnessy et al. 2020). Related to our work is adversarial-based visual explanation method is developed in (Wagner et al. 2019), highlighting the key features in the input image for a specific prediction. Contrastive explanations are produced in (Dhurandhar et al. 2018) to justify the predictions from a deep neural network. Also in (Goyal et al. 2019) the authors generate counterfactual visual explanations that highlight what and how regions of an image would need to change in order for the model to predict a *distractor* class c' instead of the predicted class c . The main differences to our construction are two fold: (*i*) they rely on a separate module to be trained, and (*ii*) these approaches only produce explanations, but such explanations are not exploited to provide feedback for model improvement.

Striking the goal of both good explanation and good performance is more challenging. One promising direction is to inject model-dependent perturbations to the input images as strategic augmentations (Fong and Vedaldi 2017; Dabkowski and Gal 2017; Chang et al. 2018). In such examples, parts of the image have been masked and replaced with various references such as mean pixel values, blurred image regions, random noise, outputs of generative models, etc. However, these pixel-level perturbations are very costly and difficult to craft. (Wang et al. 2019b) propose new learning objectives for attention separability and cross-layer consistency, which result in improved attention discriminability and reduced visual confusion. However, it generates heatmap style attention maps, which fail in fine-resolution model explanations which is important in medical related tasks. In (Fukui et al. 2019) an additional attention branch is learned to generate attention map, and then applies the attention map to the original image or feature map; they achieve com-

pelling attention maps on natural images. However, as the attention maps are not derived directly from the classification model, there is no guarantee for their faithfulness. Further, having an additional attention network results in increased network size, which raises concerns for the risk of over-fitting, particularly on datasets with a limited sample size.

CUB Experiment Details

CUB dataset descriptions and experiemnt settings

CUB has 11,788 images of 200 bird species. To train a VGG11 network, we use 8,190 training images and validate the model on 2,311 validation images, with the accuracy are reported on 1227 testing images. The network is trained for 100 epochs with a learning rate decay of 0.1 every 30 epochs. The batch size is 32. The optimizer is a SGD with initial learning rate at 0.01.

Classification performance improvement with PPI

We compare classification performances among model trained with different objections. The baseline is VGG11 classification without PPI. Three different saliency mapping methods are tested within our PPI framework: LRP, Grad-CAM, and WBP. Top 1000 points in all saliency maps are used to generate the soft mask so the comparison is fair. During training, since only a small portion of points are used to generate the mask, the contribution from L_{con} is about 100 times smaller than other losses. To fix this imbalance, the L_{con} is weighted 100× more after the first 20 epochs. The results are shown in Table 7.

Table 7: Accuracy on CUB

Models	Accuracy
VGG 11	0.662
+PPI _{LRP}	0.680
+PPI _{Grad-CAM}	0.683
+PPI _{WBP}	0.696

Geographic Atrophy (GA) Experiment Details

GA dataset descriptions

Our GA dataset is derived from the A2A SD-OCT Study (<http://ClinicalTrials.gov> identifier NCT00734487), which was an ancillary observational prospective study of a subset of eyes from the AREDS2 conducted at four sites (National Eye Institute, Duke Eye Center, Emory Eye Center, and Devers Eye Institute) (Leuschen et al. 2013). In this experiment (with our new dataset, that we will make public), each OCT volume image consists of 100 scans, each of which being a 512 × 1000 pixel image (Boyer et al. 2017). 1,085 OCT images are collected from 275 subjects during 5 years. An example of 3D OCT images is shown in Figure. 7.

Image differences between 4 sites

Images in GA dataset are collected from 4 different sites, hereafter denoted as A, B, C, and D respectively. There are 315 images (101 positive samples) from site A, 334 images (73 positive samples) from site B, 260 images (131 positive samples) from site A, and 176 images (59 positive samples) from site D. We show typical example images from 4 sites separately in Figure. 8. As the dataset is collected during 7 years, some images in site D are of smaller image size as they are sampled with different type of machine. We paddle these images by repeating the left and right areas, as show in the right bottom example.

Multi-view CNN Variation

We use a variant of the multi-view CNN model (Su et al. 2015) to process the 3D OCT inputs, and use it as our baseline solution. The architecture of this model is outlined in Figure 9. For each slice, the model feed it into a CNN network, and get the feature f_i of slice i ($f_i = CNN(x_i)$), followed by a fully connected layer and a Sigmoid activation to get a probability score $p_i = sigmoid(FC_1(f_i))$. We observe that slices in different slices contributes differently to the identification of GA, which motivates us to implement a location-aware view pooling, illustrated in the right part of Fig. 9. Each slice is assigned to a position id, ranging from 1 to 100. The model first uses an embedding layer to embed the position id to a six dimension position feature vector e_i . Then, we combine the feature vector f_i extracted from the slice image with the corresponding e_i together. The combined feature vector is fed into a fully connected layer to get the logit score a_i .

$$a_i = FC_2([f_i, e_i]) \quad (16)$$

To reduce computational burden during training, we randomly sample 10 out of the 100 slices (with an abuse of notation, denoted by a_1, a_2, \dots, a_{10}) and send them into a Softmax function to get the attention weights for the 10 sampled slices, using the following equation

$$w_i = \frac{\exp((ReLU(a_i) + \delta)/\tau)}{\sum_{k=1}^{10} \exp((ReLU(a_k) + \delta)/\tau)} \quad (17)$$

Here δ is a trainable bias term parameter, initialized to a high value to stabilize the training, and gradually attenuated to a small number during training. τ is the temperature parameter, which is set to a small value to sharpen the attention weight, which helps us to find out the most important slices for GA diagnosis. The get final predicted probability of GA (GA score) for an image x at inference time, we compute the weighted summation of the probabilities w_i of all 100 slices $GA = \sum_i w_i p_i$.

Experiment settings

The CNN network is an Inception_v3, which is pre-trained on ImageNet. For training all models, we use the Adam optimizer with a learning rate of 5×10^{-5} with a learning rate decay of 0.5 every 10 epochs for the pre-trained CNN network, and the Adam optimizer with a learning rate of 5×10^{-3} with a learning rate decay of 0.2 every 10 epochs for the

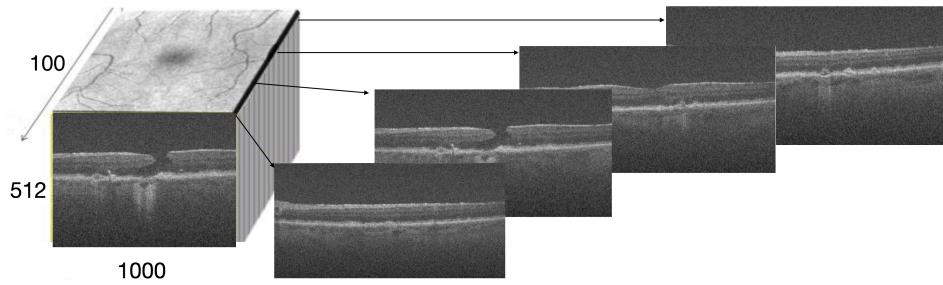


Figure 7: Illustration of a 3D OCT image example.

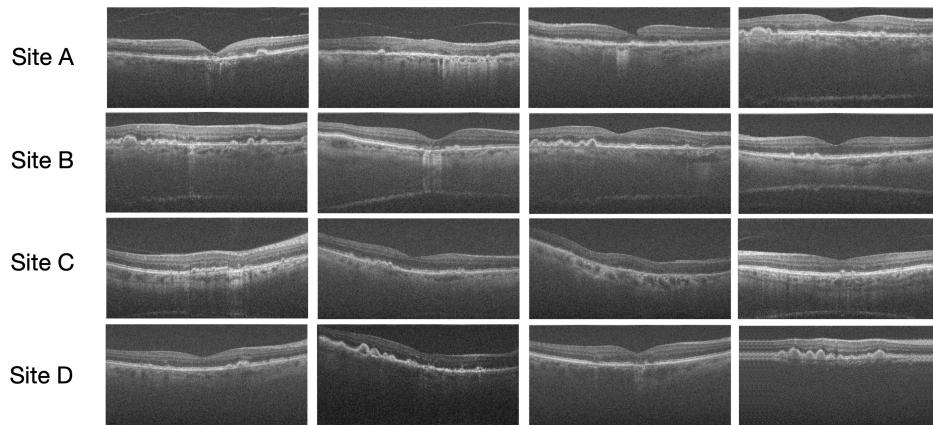


Figure 8: OCT slice examples from 4 site.

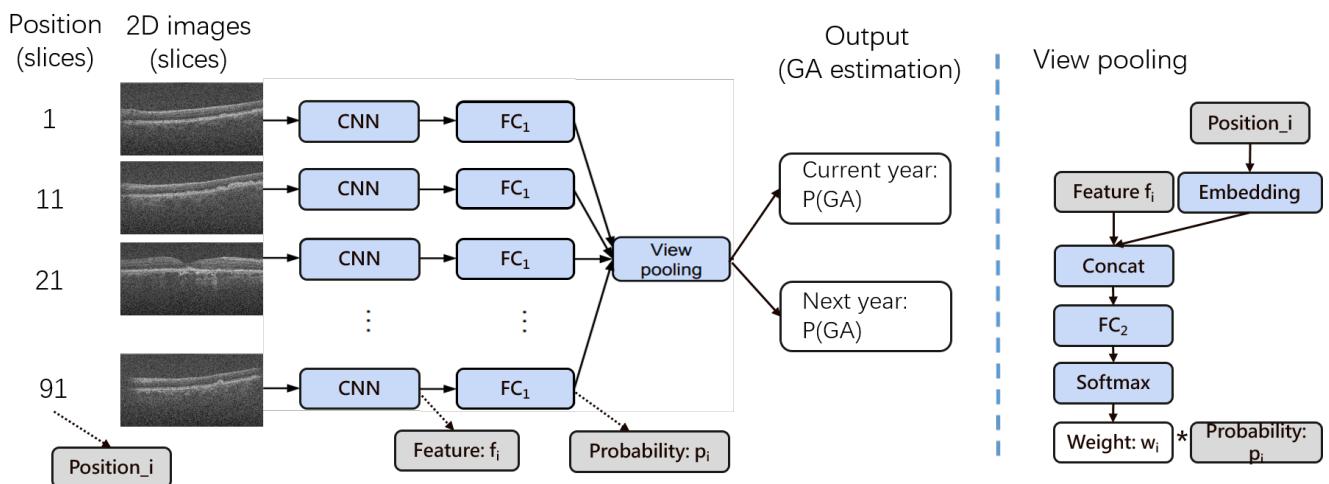


Figure 9: Illustration of multi-view CNN based 3D OCT image classification model.

other layers in the model. We use a batch size of 2 because of the large size of 3D OCT images and our multi-view CNN model (will be illustrated in). Random horizontal flips, and Gaussian noise are used for data augmentations during training.

LIDC Experiment Details

LIDC dataset description and experiment settings

We also test the proposed method on a public medical CT scan dataset LIDC-IDRI (Armato III et al. 2011). We follow the settings in (Kohl et al. 2018; Selvan and Dam 2020) that crops the original images into 128×128 patches centered on a lesion for which at least one radiologist has annotated. In our experiment, we focus on the classification task of predicting the presence of a lesions, which is consistent with the setup of (Selvan and Dam 2020). There are four radiologists annotates each patch with both lesion label and lesion mask. A patch in the dataset is labeled as positive if more than two (*i.e.* ≥ 3) radiologists have annotated presence of a lesion, otherwise negative. The ground-truth mask is the pixel-level union set of the four masks. We use Inception-v3 (Szegedy et al. 2016) as our base model for both standard classification and PPI-enhanced training with various saliency mapping schemes. To match the receptive field of an Inception-v3 model, we resize the input patches to 299×299 . For training all models, we use the Adam optimizer with a learning rate of 10^{-4} with a learning rate decay of 0.3 every 10 epochs after epoch 50, and a batch size of 64. Random horizontal flips, vertical flips and rotations within 20 degrees are used for data augmentations during training.