

A APPENDIX

A.1 PROOF OF LEMMA 1

In the mean null test, according to the second condition in (9), we have

$$\begin{aligned}\Omega_{\mathbf{X}, \mathbf{X}^{\text{ref}}} &= \Omega_{\mathbf{X}_{\text{obs}}, \mathbf{X}_{\text{obs}}^{\text{ref}}} \\ \Leftrightarrow \left(I_{2n} - \frac{\boldsymbol{\eta}_{\mathcal{M}_{\mathbf{X}}} \boldsymbol{\eta}_{\mathcal{M}_{\mathbf{X}}}^{\top}}{\boldsymbol{\eta}_{\mathcal{M}_{\mathbf{X}}}^{\top} \boldsymbol{\eta}_{\mathcal{M}_{\mathbf{X}}}} \right) \begin{pmatrix} \mathbf{X} \\ \mathbf{X}^{\text{ref}} \end{pmatrix} &= \Omega_{\mathbf{X}_{\text{obs}}, \mathbf{X}_{\text{obs}}^{\text{ref}}} \\ \Leftrightarrow \begin{pmatrix} \mathbf{X} \\ \mathbf{X}^{\text{ref}} \end{pmatrix} &= \Omega_{\mathbf{X}_{\text{obs}}, \mathbf{X}_{\text{obs}}^{\text{ref}}} + \frac{\boldsymbol{\eta}_{\mathcal{M}_{\mathbf{X}}}}{\|\boldsymbol{\eta}_{\mathcal{M}_{\mathbf{X}}}\|^2} \boldsymbol{\eta}_{\mathcal{M}_{\mathbf{X}}}^{\top} \begin{pmatrix} \mathbf{X} \\ \mathbf{X}^{\text{ref}} \end{pmatrix}.\end{aligned}$$

By defining $\mathbf{a} = \mathbf{q}_{\mathbf{X}_{\text{obs}}, \mathbf{X}_{\text{obs}}^{\text{ref}}}$, $\mathbf{b} = \frac{\boldsymbol{\eta}_{\mathcal{M}_{\mathbf{X}}}}{\|\boldsymbol{\eta}_{\mathcal{M}_{\mathbf{X}}}\|^2}$, $\mathbf{z} = \boldsymbol{\eta}_{\mathcal{M}_{\mathbf{X}}}^{\top} \begin{pmatrix} \mathbf{X} \\ \mathbf{X}^{\text{ref}} \end{pmatrix}$, we obtain the result in Lemma 1.

In the global null test, according to the second condition in (9),

$$\begin{aligned}\mathcal{U}_{\mathbf{X}, \mathbf{X}^{\text{ref}}} &= \mathcal{U}_{\mathbf{X}_{\text{obs}}, \mathbf{X}_{\text{obs}}^{\text{ref}}} \\ \Leftrightarrow P_{\mathcal{M}_{\mathbf{X}}}^{\perp} \begin{pmatrix} \mathbf{X} \\ \mathbf{X}^{\text{ref}} \end{pmatrix} &= \mathcal{U}_{\mathbf{X}_{\text{obs}}, \mathbf{X}_{\text{obs}}^{\text{ref}}} \\ \Leftrightarrow (I_{2n} - P_{\mathcal{M}_{\mathbf{X}}}) \begin{pmatrix} \mathbf{X} \\ \mathbf{X}^{\text{ref}} \end{pmatrix} &= \mathcal{U}_{\mathbf{X}_{\text{obs}}, \mathbf{X}_{\text{obs}}^{\text{ref}}} \\ \Leftrightarrow \begin{pmatrix} \mathbf{X} \\ \mathbf{X}^{\text{ref}} \end{pmatrix} &= \mathcal{U}_{\mathbf{X}_{\text{obs}}, \mathbf{X}_{\text{obs}}^{\text{ref}}} + \mathcal{V}_{\mathbf{X}_{\text{obs}}, \mathbf{X}_{\text{obs}}^{\text{ref}}} \sigma^{-1} \left\| P_{\mathcal{M}_{\mathbf{X}}} \begin{pmatrix} \mathbf{X} \\ \mathbf{X}^{\text{ref}} \end{pmatrix} \right\|.\end{aligned}$$

By defining $\mathbf{a} = \mathcal{U}_{\mathbf{X}_{\text{obs}}, \mathbf{X}_{\text{obs}}^{\text{ref}}}$, $\mathbf{b} = \mathcal{V}_{\mathbf{X}_{\text{obs}}, \mathbf{X}_{\text{obs}}^{\text{ref}}}$, $\mathbf{z} = \sigma^{-1} \left\| P_{\mathcal{M}_{\mathbf{X}}} \begin{pmatrix} \mathbf{X} \\ \mathbf{X}^{\text{ref}} \end{pmatrix} \right\|$, we obtain the result in Lemma 1.

A.2 EXAMPLES OF PIECEWISE LINEAR FUNCTIONS

Examples of piecewise linear components in a trained CNN with $\mathbf{X} \in \mathbb{R}^2$ are provided as follows:

ReLU: Consider f is ReLU function. Then, $K(f) = 4$, $\boldsymbol{\psi}_k = (0 \ 0)^{\top}$ for any $k \in [4]$,

$$\begin{aligned}\Psi_1^f &= \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix}, \mathcal{P}_1^f = \left\{ \mathbf{X} : \begin{matrix} X_1 < 0, \\ X_2 < 0 \end{matrix} \right\}, \quad \Psi_2^f = \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix}, \mathcal{P}_2^f = \left\{ \mathbf{X} : \begin{matrix} X_1 < 0, \\ X_2 \geq 0 \end{matrix} \right\}, \\ \Psi_3^f &= \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix}, \mathcal{P}_3^f = \left\{ \mathbf{X} : \begin{matrix} X_1 \geq 0, \\ X_2 < 0 \end{matrix} \right\}, \quad \Psi_4^f = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}, \mathcal{P}_4^f = \left\{ \mathbf{X} : \begin{matrix} X_1 \geq 0, \\ X_2 \geq 0 \end{matrix} \right\}.\end{aligned}$$

This can be similarly extended to the case of Leaky ReLU.

Max-pooling: Consider $f(\mathbf{X}) = \max\{X_1, X_2\}$. Then, it is represented as a piecewise linear function with $K(f) = 2$, $\boldsymbol{\psi}_k = (0)$ for any $k \in [2]$,

$$\Psi_1^f = \begin{pmatrix} 1 & 0 \end{pmatrix}, \mathcal{P}_1^f = \{\mathbf{X} : X_1 \geq X_2\}, \quad \Psi_2^f = \begin{pmatrix} 0 & 1 \end{pmatrix}, \mathcal{P}_2^f = \{\mathbf{X} : X_1 < X_2\}.$$

Convolution and matrix-vector multiplication: In a neural network, the multiplication results between the weight matrix and the output of the previous layer and its summation with the bias vector is a linear function. In a CNN, the convolution operation is obviously a linear function.

Upsampling: Consider f is the upsampling operation on $\mathbf{X} \in \mathbb{R}^2$, then it can be represented as a piecewise linear function with $K(f) = 1$, $\boldsymbol{\psi}_1 = (0 \ 0 \ 0 \ 0)^{\top}$,

$$\Psi_1^f = \begin{pmatrix} 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \end{pmatrix}^{\top}, \quad \mathcal{P}_1^f = \mathbb{R}^2.$$

Sigmoid and hyperbolic tangent: If there is any specific demand to use non-piecewise linear activation functions, we can apply a piecewise-linear approximation approach to these functions.

A.3 PROOF OF LEMMA 2

At f_1 , given a real value z , the input is $\beta^{f_0} + \gamma^{f_0} z = \mathbf{a}_{1:n} + \mathbf{b}_{1:n} z$. By checking the value of this input, we can easily obtain the polytope

$$\{\Delta_{k_1}^{f_1}(\beta^{f_0} + \gamma^{f_0} z) \leq \delta_{k_1}^{f_1}\}, \quad k_1 \in [K(f_1)],$$

that $\beta^{f_0} + \gamma^{f_0} z$ belongs to. Based on the obtained polytope, we can calculate the interval $[L_{k_1}^{f_1}, U_{k_1}^{f_1}]$,

$$L_{k_1}^{f_1} = \max_{j: (\Delta_{k_1}^{f_1} \gamma^{f_0})_j < 0} \frac{(\delta_{k_1}^{f_1})_j - (\Delta_{k_1}^{f_1} \beta^{f_0})_j}{(\Delta_{k_1}^{f_1} \gamma^{f_0})_j} \quad \text{and} \quad U_{k_1}^{f_1} = \min_{j: (\Delta_{k_1}^{f_1} \gamma^{f_0})_j > 0} \frac{(\delta_{k_1}^{f_1})_j - (\Delta_{k_1}^{f_1} \beta^{f_0})_j}{(\Delta_{k_1}^{f_1} \gamma^{f_0})_j}.$$

Moreover, based on the obtained polytope, we can easily obtain $\Psi_{k_1}^{f_1}$ and $\psi_{k_1}^{f_1}$, $k_1 \in [K(f_1)]$. Therefore, the output of the first layer at z can be defined as

$$\begin{aligned} f_1(z) &= \Psi_{k_1}^{f_1}(\beta^{f_0} + \gamma^{f_0} z) + \psi_{k_1}^{f_1} \\ &= \beta^{f_1} + \gamma^{f_1} z, \end{aligned}$$

where $\beta^{f_1} = \Psi_{k_1}^{f_1} \beta^{f_0} + \psi_{k_1}^{f_1}$ and $\gamma^{f_1} = \Psi_{k_1}^{f_1} \gamma^{f_0}$. Next, we input $\beta^{f_1}, \gamma^{f_1}$ to f_2 .

At the 2nd layer, similarly, the input is $\beta^{f_1} + \gamma^{f_1} z$. By checking the value of this input, we can easily obtain the polytope

$$\{\Delta_{k_2}^{f_2}(\beta^{f_1} + \gamma^{f_1} z) \leq \delta_{k_2}^{f_2}\}, \quad k_2 \in [K(f_2)],$$

that $\beta^{f_1} + \gamma^{f_1} z$ belongs to. Based on the obtained polytope, we can calculate the interval $[L_{k_2}^{f_2}, U_{k_2}^{f_2}]$,

$$L_{k_2}^{f_2} = \max_{j: (\Delta_{k_2}^{f_2} \gamma^{f_1})_j < 0} \frac{(\delta_{k_2}^{f_2})_j - (\Delta_{k_2}^{f_2} \beta^{f_1})_j}{(\Delta_{k_2}^{f_2} \gamma^{f_1})_j} \quad \text{and} \quad U_{k_2}^{f_2} = \min_{j: (\Delta_{k_2}^{f_2} \gamma^{f_1})_j > 0} \frac{(\delta_{k_2}^{f_2})_j - (\Delta_{k_2}^{f_2} \beta^{f_1})_j}{(\Delta_{k_2}^{f_2} \gamma^{f_1})_j}.$$

Moreover, based on the obtained polytope, we can easily obtain $\Psi_{k_2}^{f_2}$ and $\psi_{k_2}^{f_2}$, $k_2 \in [K(f_2)]$. Therefore, the output of the first layer at z can be defined as

$$\begin{aligned} f_2(z) &= \Psi_{k_2}^{f_2}(\beta^{f_1} + \gamma^{f_1} z) + \psi_{k_2}^{f_2} \\ &= \beta^{f_2} + \gamma^{f_2} z, \end{aligned}$$

where $\beta^{f_2} = \Psi_{k_2}^{f_2} \beta^{f_1} + \psi_{k_2}^{f_2}$ and $\gamma^{f_2} = \Psi_{k_2}^{f_2} \gamma^{f_1}$.

A.4 EXPERIMENTAL DETAILS.

Methods for comparison. We compared our proposed method with the following approaches:

- Naive: the classical z -test is used to calculate the naive p -value.
- Bonferroni: the number of all possible hypotheses are considered to account for the selection bias. The p -value is computed by $p_{\text{bonferroni}} = \min(1, p_{\text{naive}} * 2^n)$
- Over-conditioning (OC): additionally conditioning on the observed activeness and inactiveness of all the nodes. The limitation of this method is its low statistical power due to over-conditioning.

Network for experiments. In all the experiments, we used the network structure shown in Fig. 8. We generate 1000 images from $N(\mathbf{0}, I_n)$ as negative samples, and 1000 images from $N(\mathbf{s}, I_n)$ as positive samples. We set $s_i = 0, \forall s_i \in \mathcal{S}$, and $s_i \neq 1, \forall s_i \notin \mathcal{S}$, where \mathcal{S} are the set of indexes whose pixels have signals and the location is randomly determined and $|\mathcal{S}| = n/4$ for each image size n . We trained the network with these images and used Adam (Kingma & Ba, 2014) as an optimizer.

Experiment for robustness. Additionally, we confirm the robustness of the proposed method in terms of FPR control by applying our method to data following Laplace distribution (Laplace), skew normal (Skew) distribution and t distribution (t-dist). The results are shown in Fig. 9. Our method still maintains good performance in terms of FPR control.

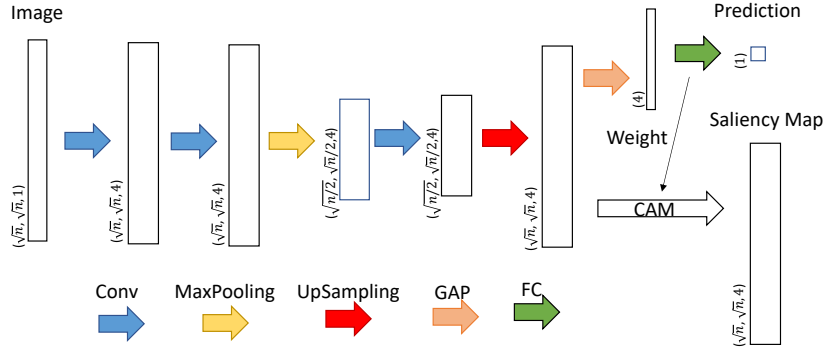


Figure 8: Network structure.

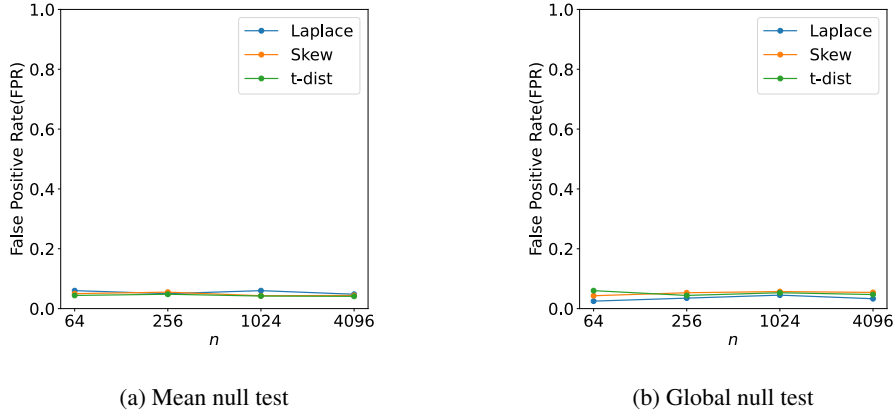


Figure 9: False positive rate of the proposed method when data are nonnormal

Experimental setting on brain image dataset. We examine the brain image dataset extracted from the dataset used in Buda et al. (2019), which includes 941 images without tumors (C1) and 939 images with tumors (C2). We selected 50 images from C1 as reference images and used 841 images from C1 and 889 images from C2 for DNN training. The remaining images from C1 and C2 are used for demonstrating the advantages of the proposed selective p -value.

More results on brain image dataset. Additional results are shown in Figs. 10, 11, 12 and 13

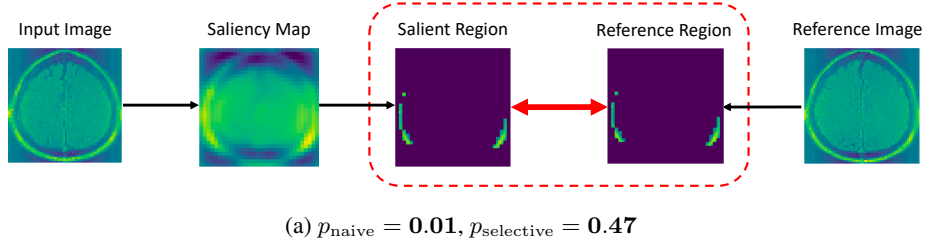


Figure 10: Inference on salient regions for images without tumor (mean null test).

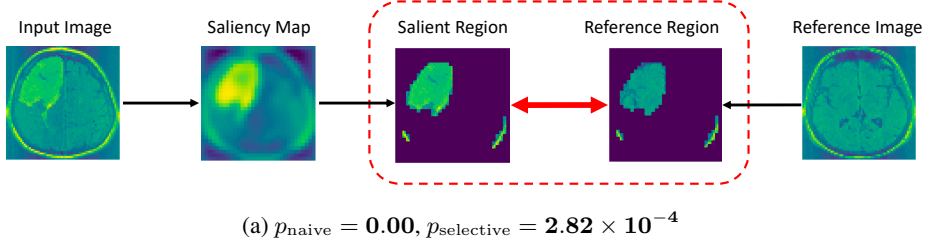


Figure 11: Inference on salient regions for images where there exists a tumor (mean null test).

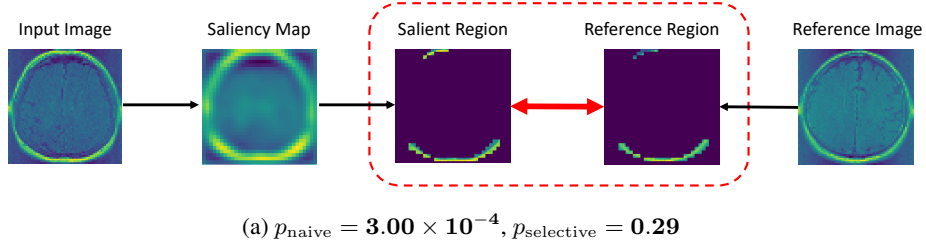


Figure 12: Inference on salient regions for images without tumor (global null test).

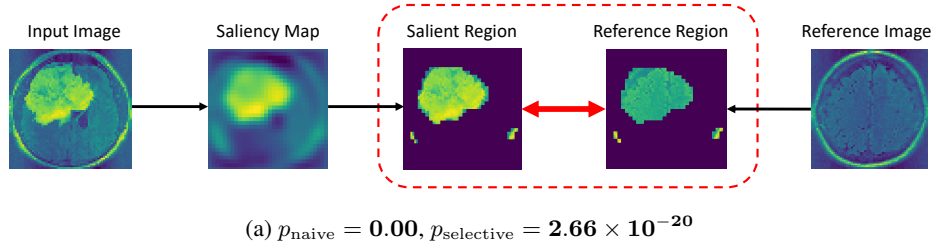


Figure 13: Inference on salient regions for images where there exists a tumor (global null test).