

# TOWARD TRUSTWORTHY HEALTHCARE AI: ATTENTION-BASED FEATURE LEARNING FOR COVID-19 WITH CHEST RADIOGRAPHY

7/22/2022

Kai Ma, Pengcheng Xi, Karim Habashy, Ashkan Ebadi, Stéphane Tremblay, Alexander Wong

Faculty of Engineering, University of Waterloo

Digital Technologies Research Centre, National Research Council Canada



# OVERVIEW

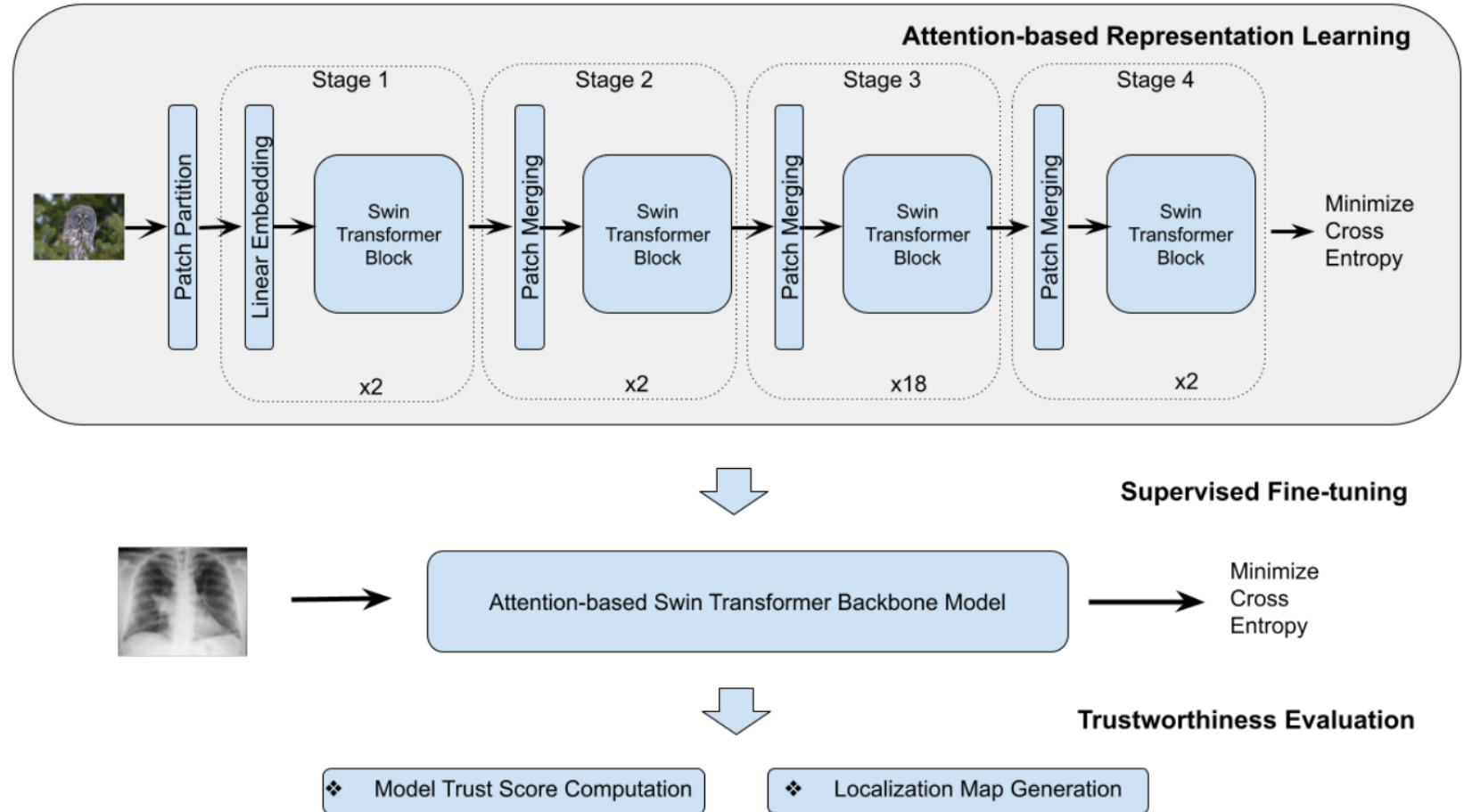
- Evaluation of CNNs and Transformer deep learning architectures
  - Traditional CNN approach vs. Attention mechanism
  - Representation learning capability for chest X-ray (CXR) classification
- Model performance — Precision, Sensitivity
- Model trustworthiness
  - Trust Score — Wong et. al., 2020
  - Visualization — AblationCAM

# MOTIVATION

- Importance of trust in medical AI
  - Explainability of results and model confidence are crucial
  - Quantifiable metric
- Transformer's potential architectural advantages
  - CNNs have inductive biases
  - Attention mechanism — analogous to imaging evaluation by doctors
- Fair evaluation of localization maps
  - Existing literature uses visualization techniques that are specific to architecture

# METHODOLOGY - MODEL ARCHITECTURE

- Models:
  - ResNet-50
  - DenseNet-121
  - Swin Transformer



# METHODOLOGY – DATASET, TRUST QUANTIFICATION

- Dataset: COVIDx V9B
  - Combination of data repositories (RSNA, etc)
  - 10% of training set sampled for validation
- Trust score computation:
  - Penalize undeserved confidence, reward well-placed confidence
- AblationCAM
  - Localization maps through ablation analysis

Table 1. Data split for COVIDx V9B

SPLIT	NEGATIVE	POSITIVE	TOTAL
TRAIN	13,992	15,950	30,482
TEST	200	200	400

$$Q_z(x, y) = \begin{cases} C(y | x)^\alpha, & \text{if } x \in R_{y=z|M} \\ (1 - C(y | x))^\beta, & \text{if } x \in R_{y \neq z|M}, \end{cases}$$

# RESULTS – PERFORMANCE & TRUST SCORE

Table 3. Precision scores on the unseen COVIDx V9B test split. The best results in each class are bolded.

MODEL	NEGATIVE	POSITIVE
RESNET (200 EPOCHS)	<b>0.952</b>	<b>1.000</b>
DENSENET (200 EPOCHS)	0.948	0.995
SWIN-B (30 EPOCHS)	0.926	<b>1.000</b>
SWIN-B (50 EPOCHS)	0.935	<b>1.000</b>
SWIN-B (100 EPOCHS)	0.930	<b>1.000</b>
SWIN-B (200 EPOCHS)	<b>0.952</b>	<b>1.000</b>

Table 4. Sensitivity scores on the unseen COVIDx V9B test split. The best results in each class are bolded.

MODEL	NEGATIVE	POSITIVE
RESNET (200 EPOCHS)	<b>1.000</b>	<b>0.950</b>
DENSENET (200 EPOCHS)	0.995	0.945
SWIN-B (30 EPOCHS)	<b>1.000</b>	0.920
SWIN-B (50 EPOCHS)	<b>1.000</b>	0.930
SWIN-B (100 EPOCHS)	<b>1.000</b>	0.925
SWIN-B (200 EPOCHS)	<b>1.000</b>	<b>0.950</b>

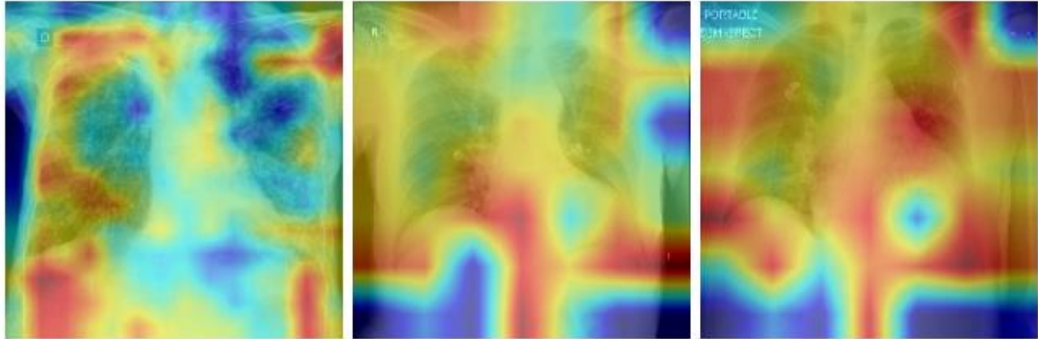
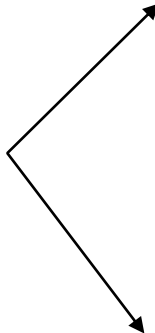
Table 5. Trust scores calculated from each experiment on the positive class. The best result is bolded.

MODEL	TRUST SCORE
RESNET (200 EPOCHS)	0.923
DENSENET (200 EPOCHS)	0.922
SWIN-B (30 EPOCHS)	0.943
SWIN-B (50 EPOCHS)	0.959
SWIN-B (100 EPOCHS)	0.954
SWIN-B (200 EPOCHS)	<b>0.963</b>

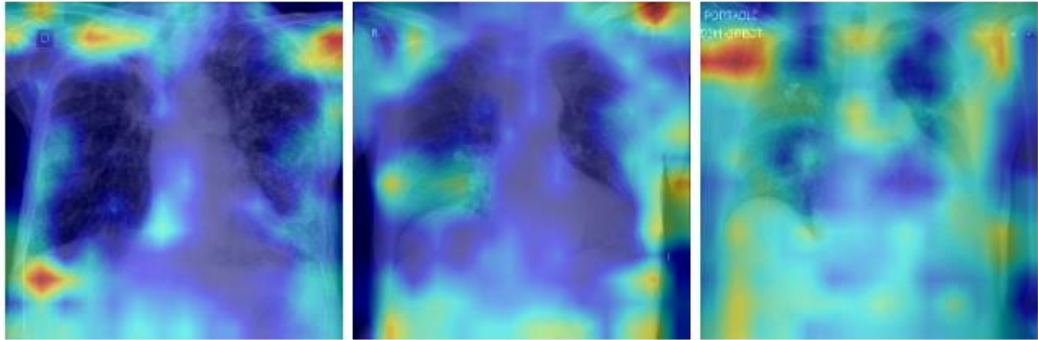
# RESULTS - VISUALIZATION



(a) Original chest radiographs for positive COVID-19 samples



(b) Swin-B 200-epoch Ablation-CAM



(c) ResNet-50 200-epoch Ablation-CAM

# WHAT'S NEXT?

- Comparison against newer, leading CNNs and different Transformer models
  - ConvNeXt, NFNet
  - Swin-L, CoCa, MaxViT
- Further validation of results on other datasets
  - SIIM-FISABIO-RSNA, 3-class COVIDx
- Exploration of data-hungriness of Transformers
- Departure from model evaluation to introducing more novelty



# THANKS FOR LISTENING!

Feel free to ask any questions!