

Automatic Diagnosis of Pulmonary Embolism Using an Attention-guided Framework: A Large-scale Study

Luyao Shi¹, Deepta Rajan², Shafiq Abedin²,
Manikanta Srikar Yellapragada³, David Beymer²,
and Ehsan Dehghan²

¹Yale University

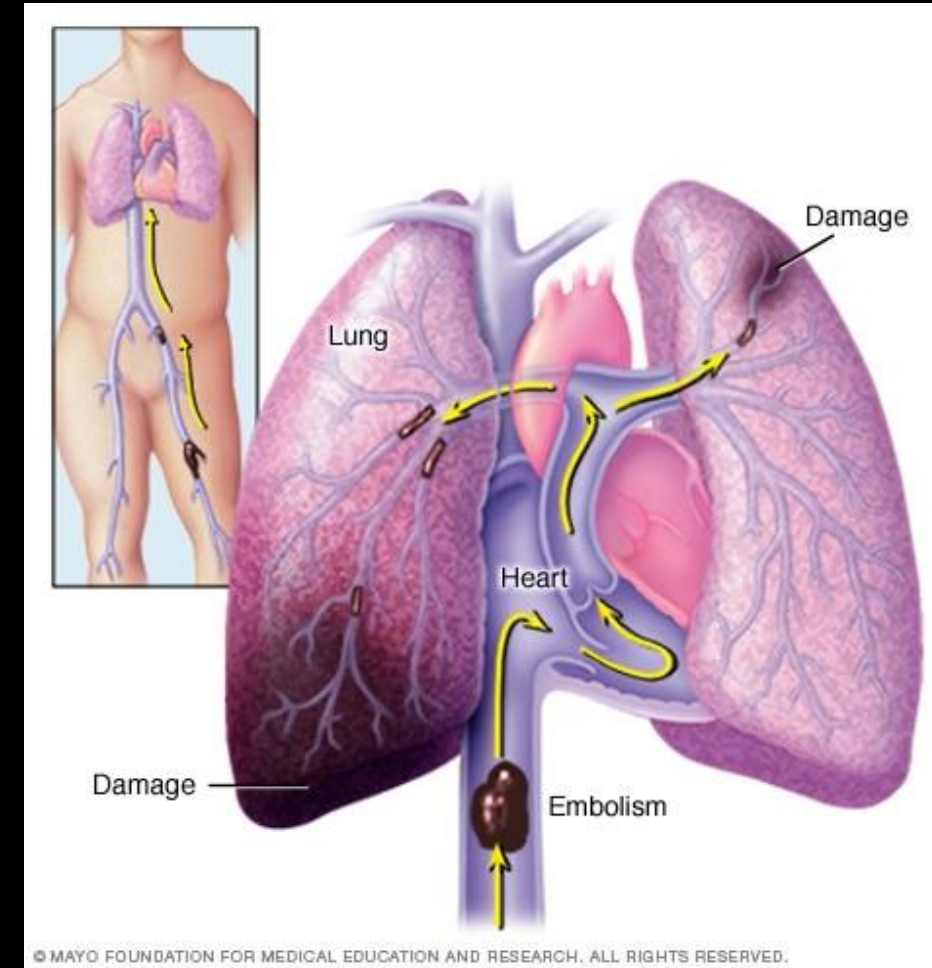
²IBM Research

³New York University



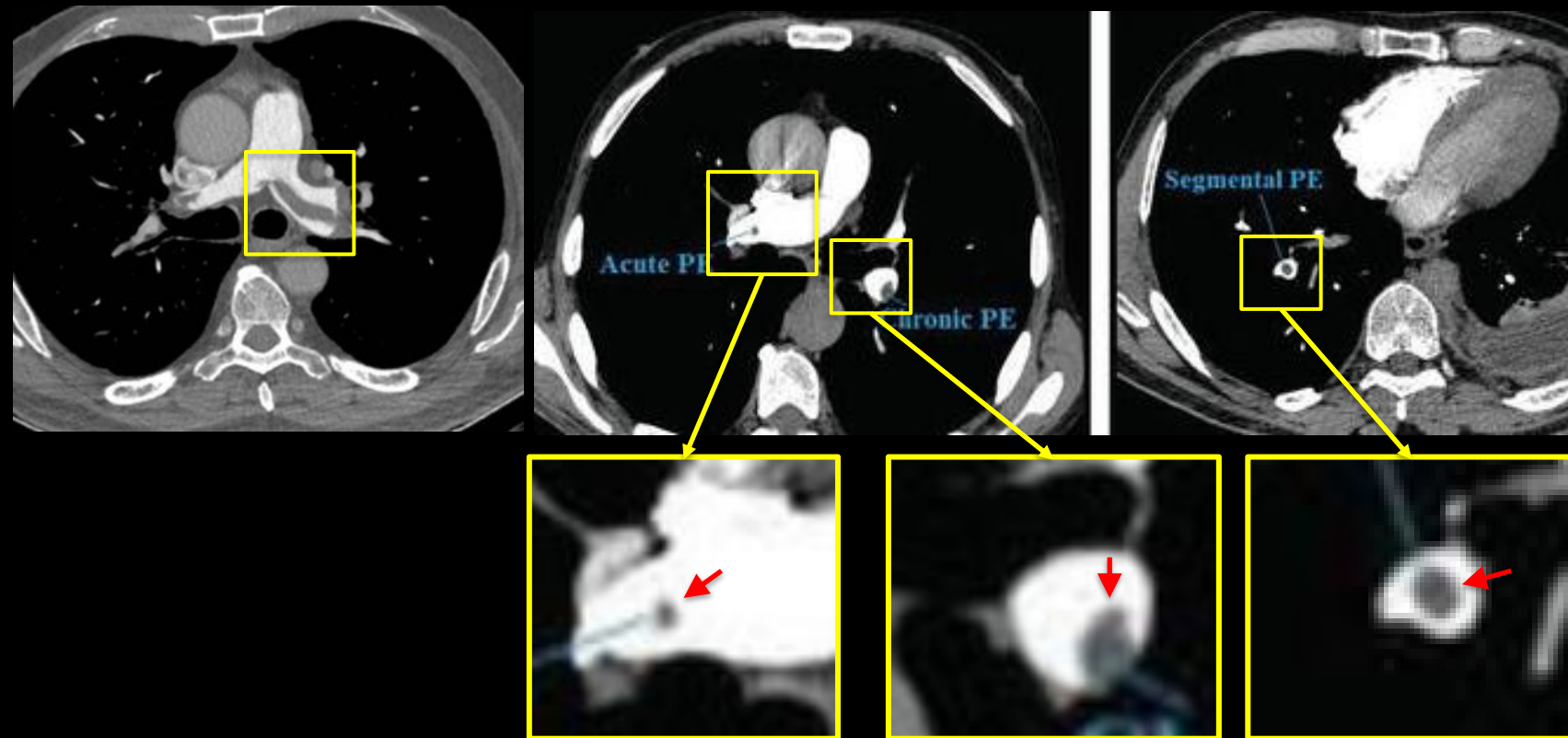
About pulmonary embolism (PE)

- **Causes:** a clump of material, most often a blood clot, gets wedged into an artery in patients' lungs. These blood clots most commonly come from the deep veins of patients' legs.
- **Mortality:**
 - About 100,000 deaths/year in US.
 - 1 of 4 people who have a PE die without warning.
 - 10 to 30% of people will die within one month of diagnosis.
- Prompt recognition of the diagnosis and immediate initiation of therapeutic action is important.

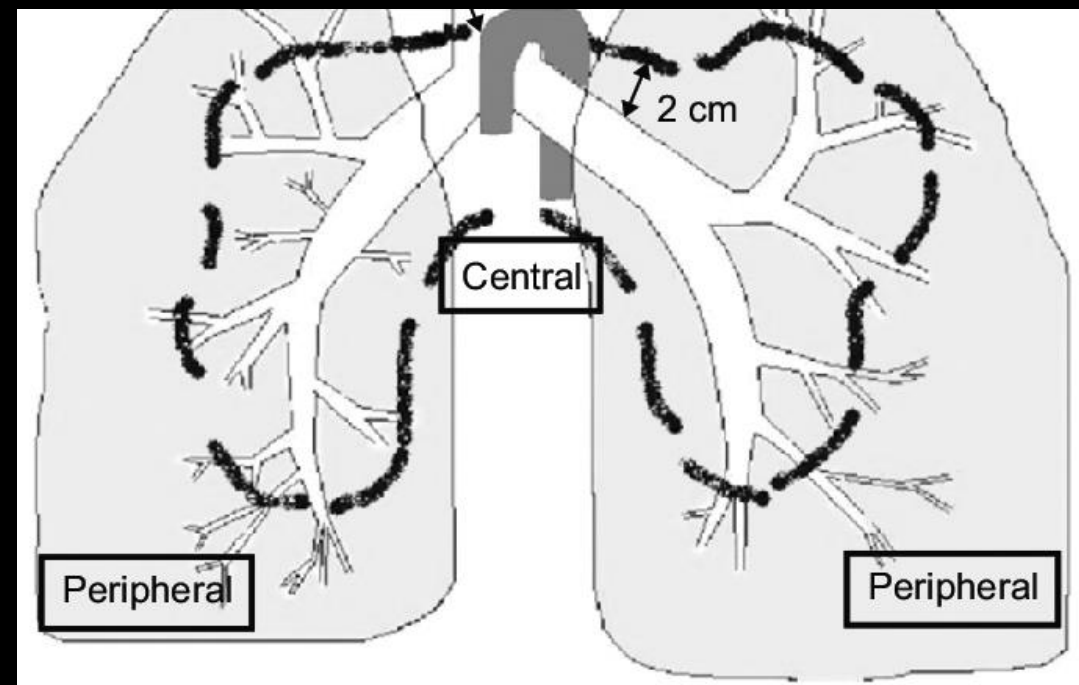


About pulmonary embolism (PE)

- **Contrast Enhanced Chest CT** is the preferred method of diagnostic imaging in patients with a clinical risk score indicative of PE.
- PE can be visualized as perfusion defects.



- **Challenges:**
 - Increased probability of **false-positive** findings when the lesions involve peripheral pulmonary vascular regions.
 - Confounding factors:
 - Poorly filled vein with contrast media
 - Impacted bronchi or parenchymal disease
 - Lymphoid tissues around the vessels
 - Respiratory/cardiac motion artifacts
 - Image noise
 - PE detection/exclusion is quite **time-consuming** and dependent on the experience of the radiologist.
- **GOAL:** A deep learning-based computer-aided diagnosis (CAD) platform to detect PE with high accuracy.



In-Hye Jung et al. Clinical outcome of fiducial-less Cyber Knife radiosurgery for stage I non-small cell lung cancer.(2015).



Training with
pixel-level annotated data



End-to-end training with
patient-level labels



Hybrid Training

Pros

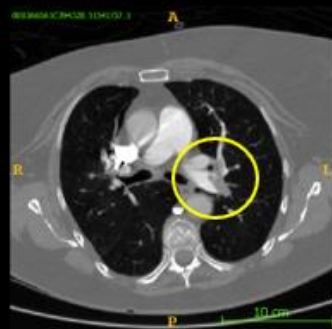
- Better interpretability
- Higher accuracy

- Largely available training data
- Better scalability

Cons

- Time consuming for radiologists
- Limited availability
- Less scalability

- Less interpretability
- Potentially worse performance



pixel-level annotated data



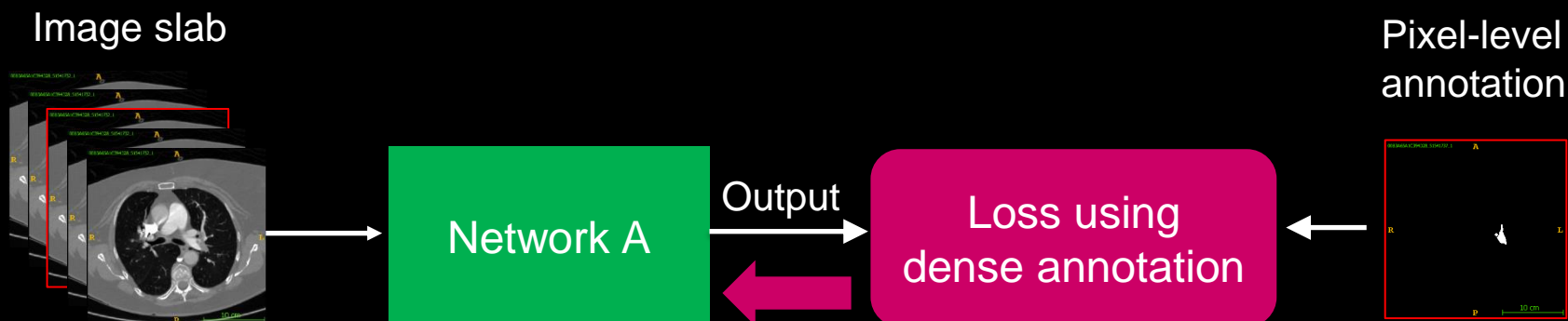
patient-level label

→ PE or not?



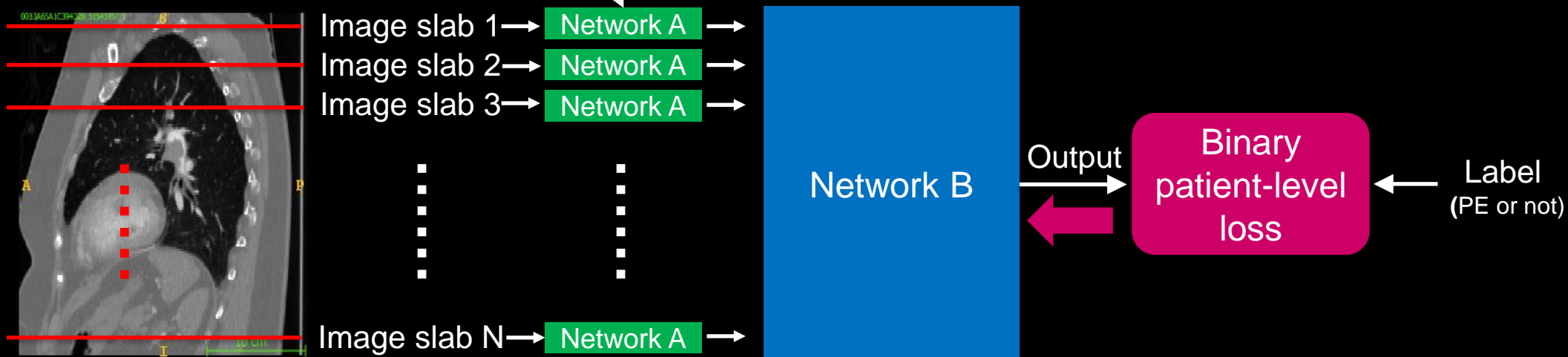
Hybrid Training Overview

Stage 1:
Training with
pixel-level annotated data



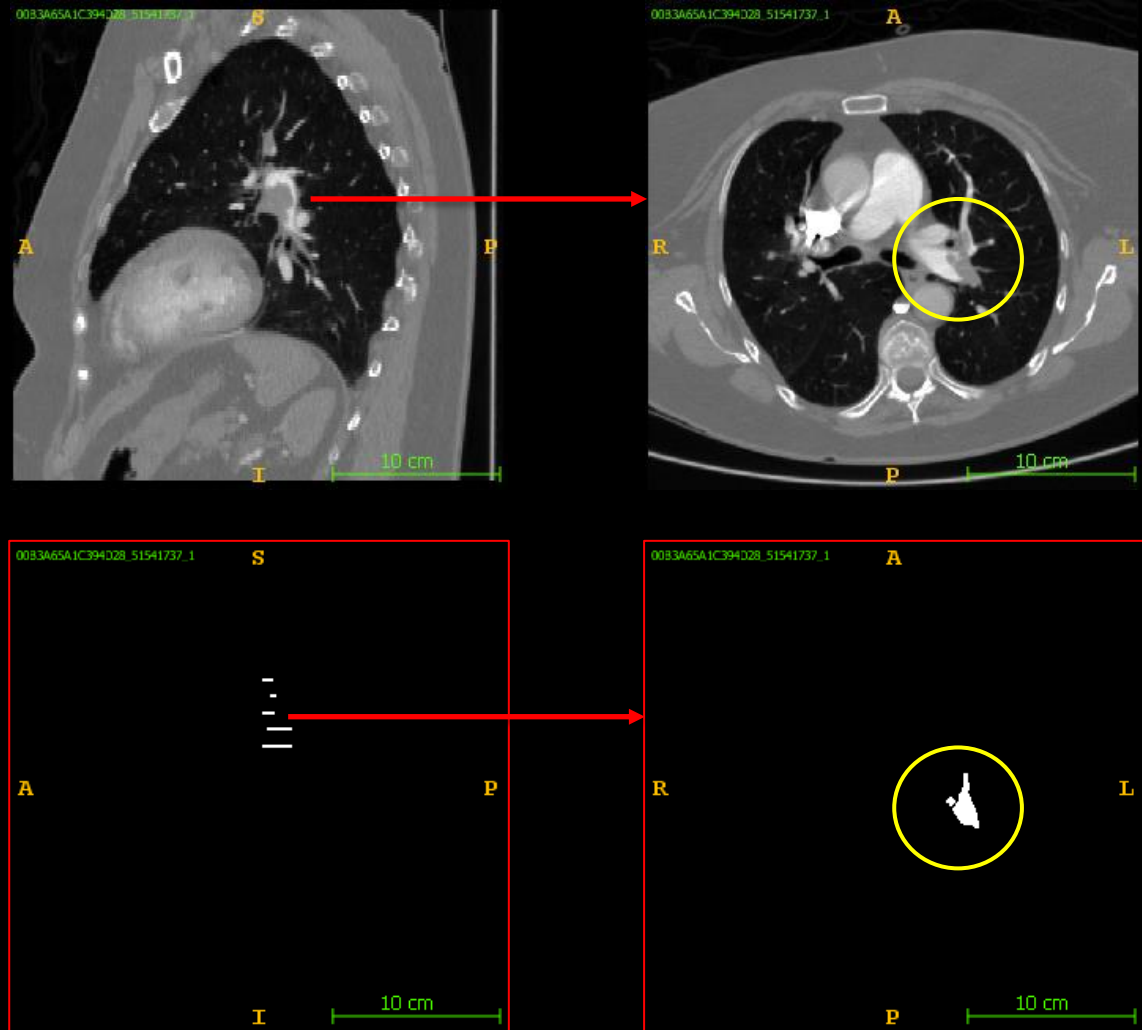
Feature encoder
(weight fixed)

Stage 2:
Training with
patient-level labels



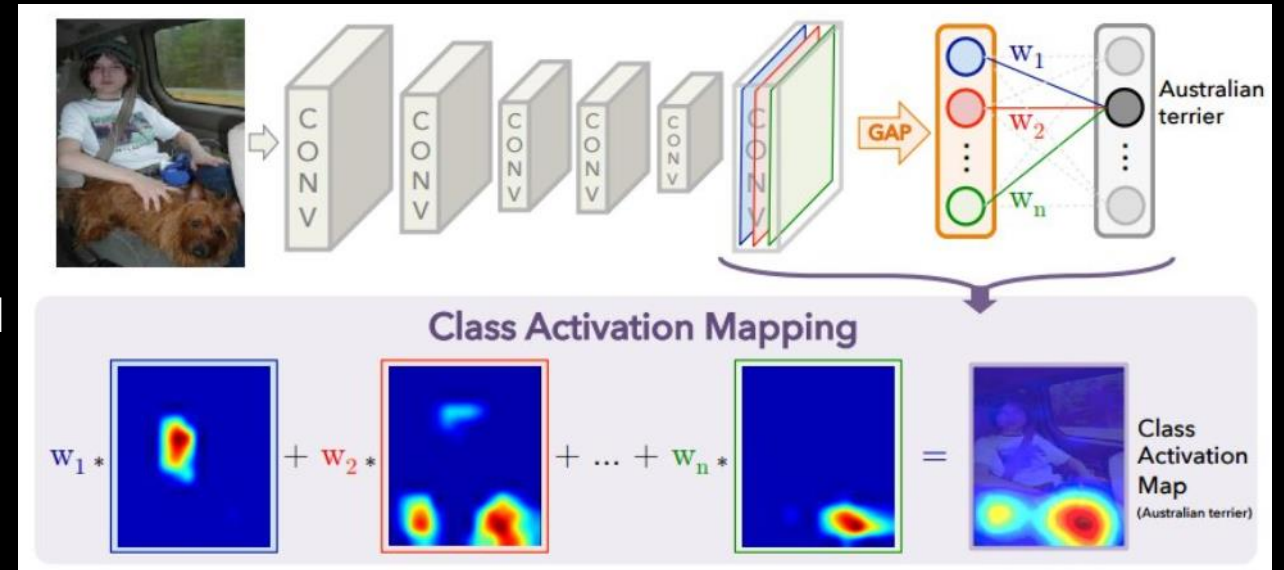
Stage1: training with pixel-level annotated data

- Pixel-level annotations every 10mm
- **Goal:** train an image encoding network that focus its attention on PE



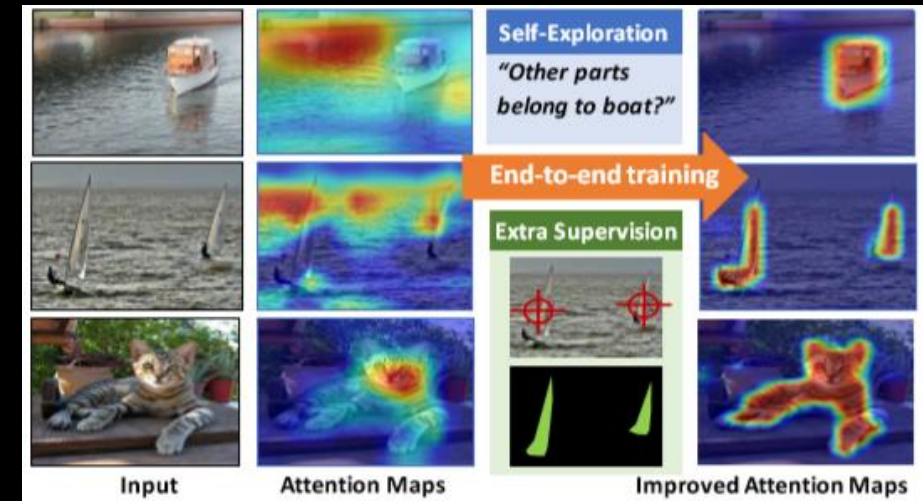
Attention map

- **Class activation map (CAM):** indicates the discriminative image regions used by the CNN to identify a particular class.



Zhou et al. Learning Deep Features for Discriminative Localization. (2016)

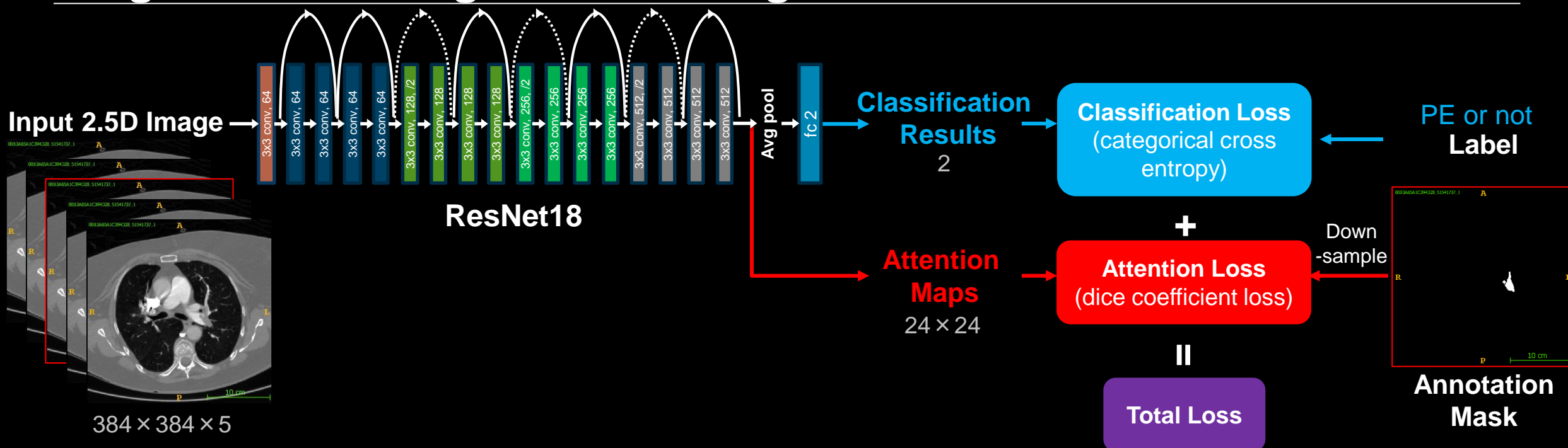
- **Guided attention inference networks (GAIN):** supervise the attention maps while training the network.



Li et al., Tell Me Where to Look: Guided Attention Inference Network. (2018)



Stage1: attention-guided training



- Resample volumetric images (bilinear interpolation): slice thickness [0.5mm, 5mm] \rightarrow 2.5mm
- 10,388 slabs (5 slices) of annotated pairs from 1,670 positive volumetric images
- Same amount of negative slabs randomly sampled from 593 negative volumetric images
- Image cropped to center 384×384 , [-1024HU, 500HU] \rightarrow [0, 255]
- 80% training, 20% validation
- Training epochs: 100 (save the model with the highest val. acc.)



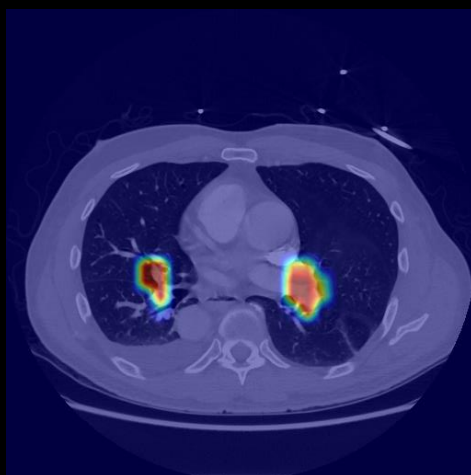
Stage1: results on the validation set

Example Attention Maps

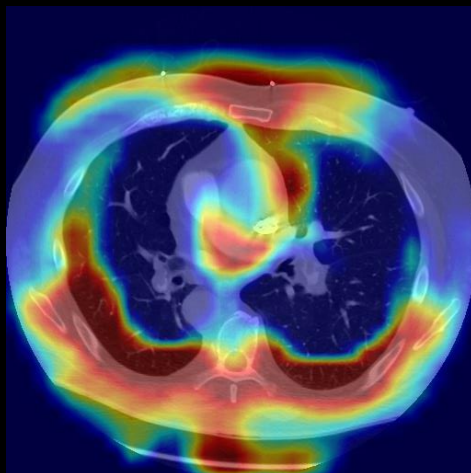
Annotation Mask
(down-sampled)

Attention Map

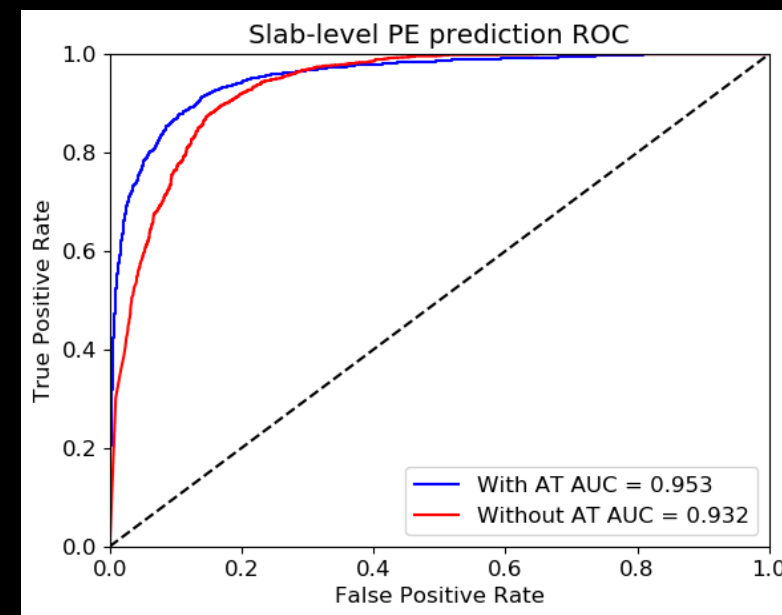
With
Attention Training



Without
Attention Training



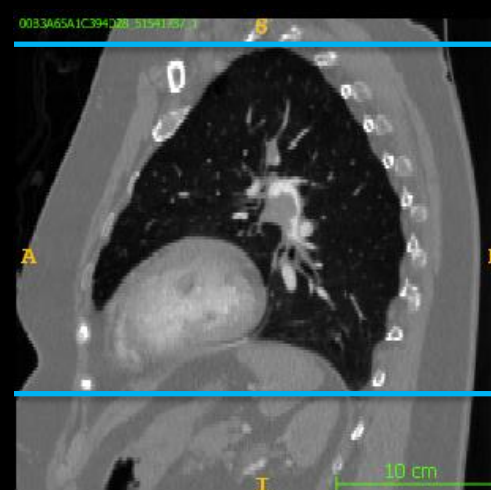
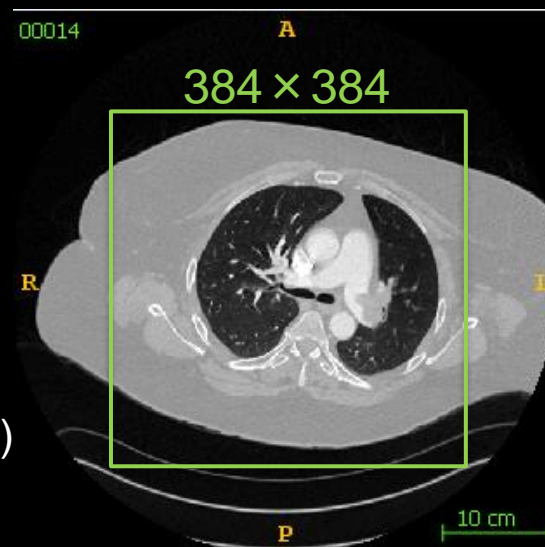
ResNet's slab-level PE prediction result on the validation data



Stage2: training data and image pre-processing

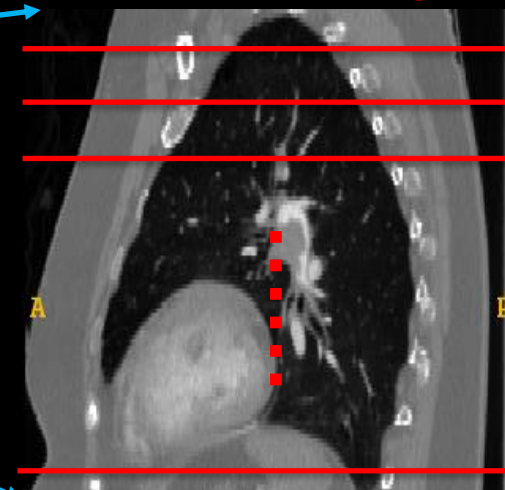
- **Data-preprocessing:**

- Image cropped to center 384×384 , $[-1024\text{HU}, 500\text{HU}] \rightarrow [0-255]$
- Identify lung regions using lung mask (produced by in-house lung segmentation method) resize to 200 slices, then sample **50 slabs**
- $? \times 512 \times 512 \rightarrow 50 \times 384 \times 384 \times 5$

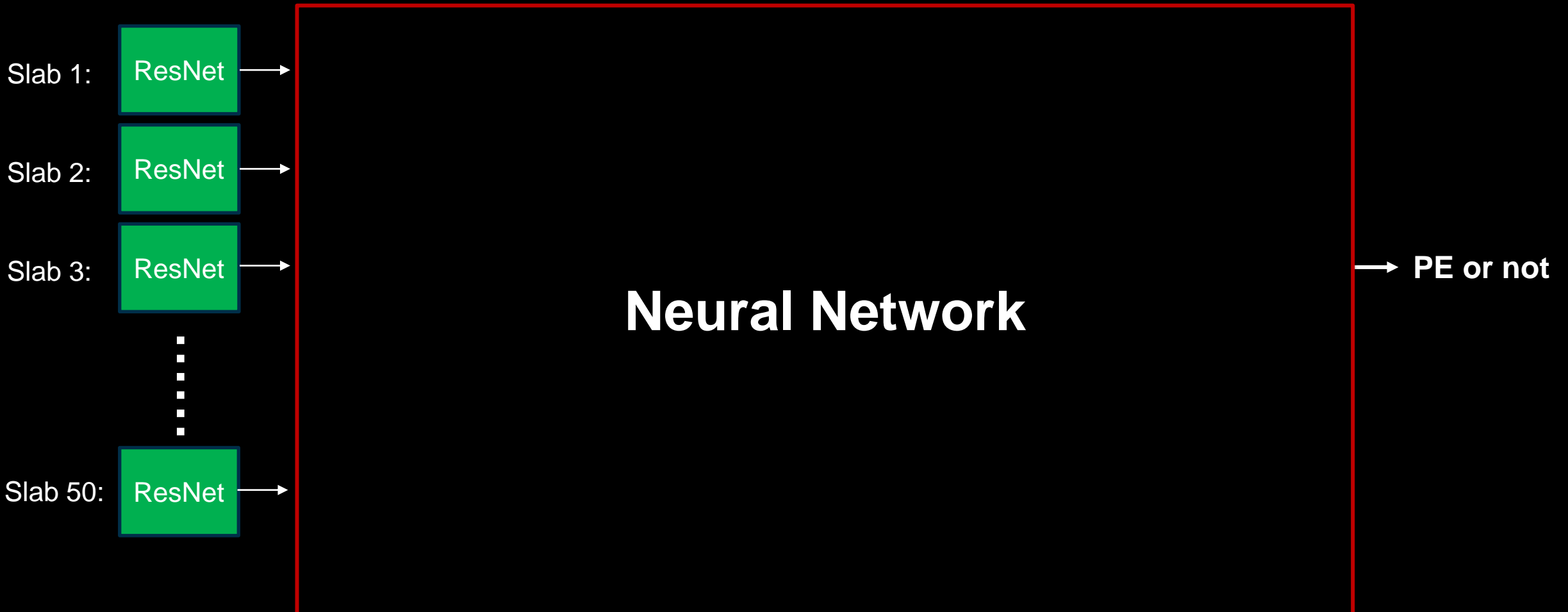


Resize

Slab Sampling

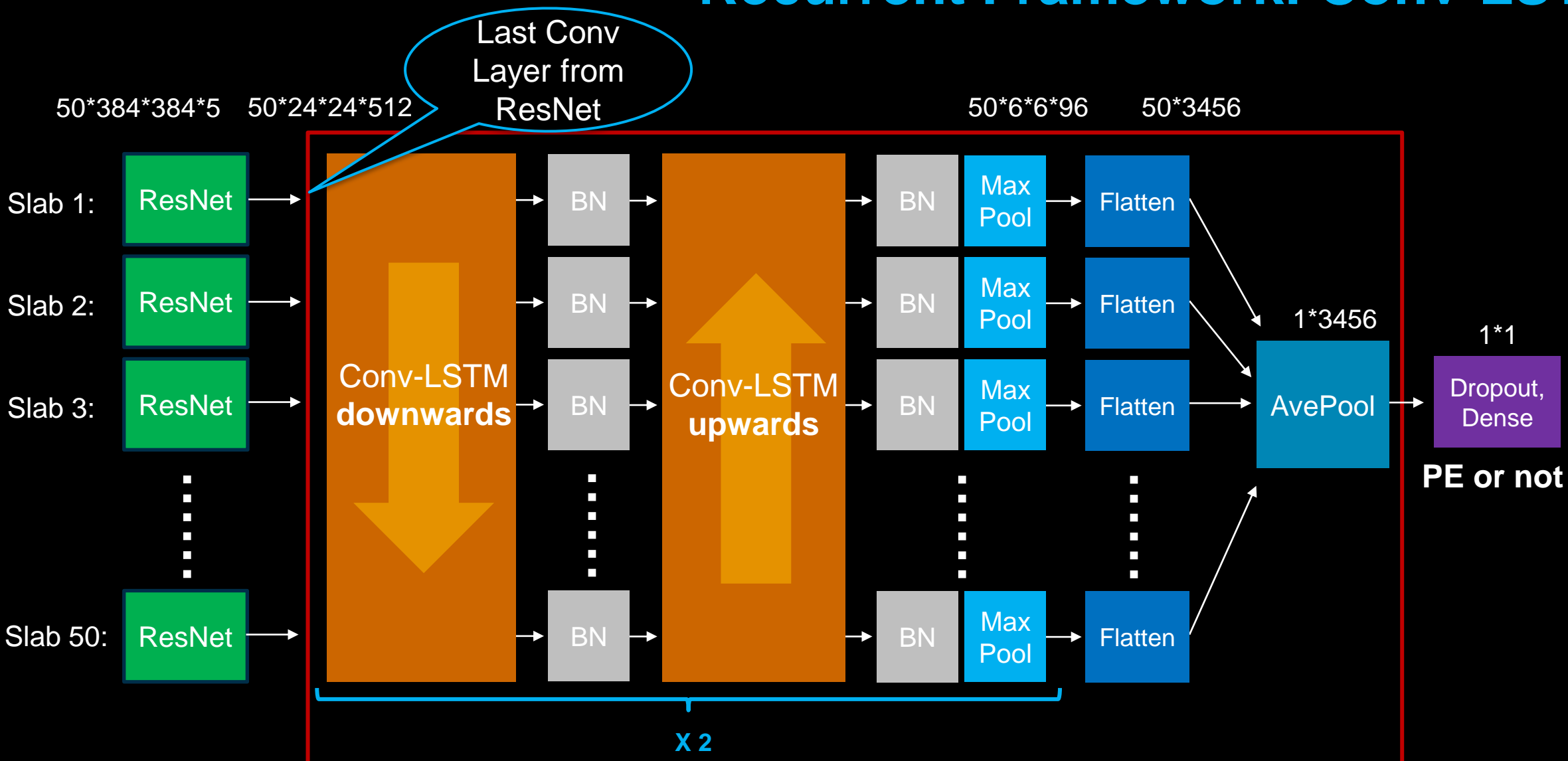


Stage2: training with patient-level labeled data

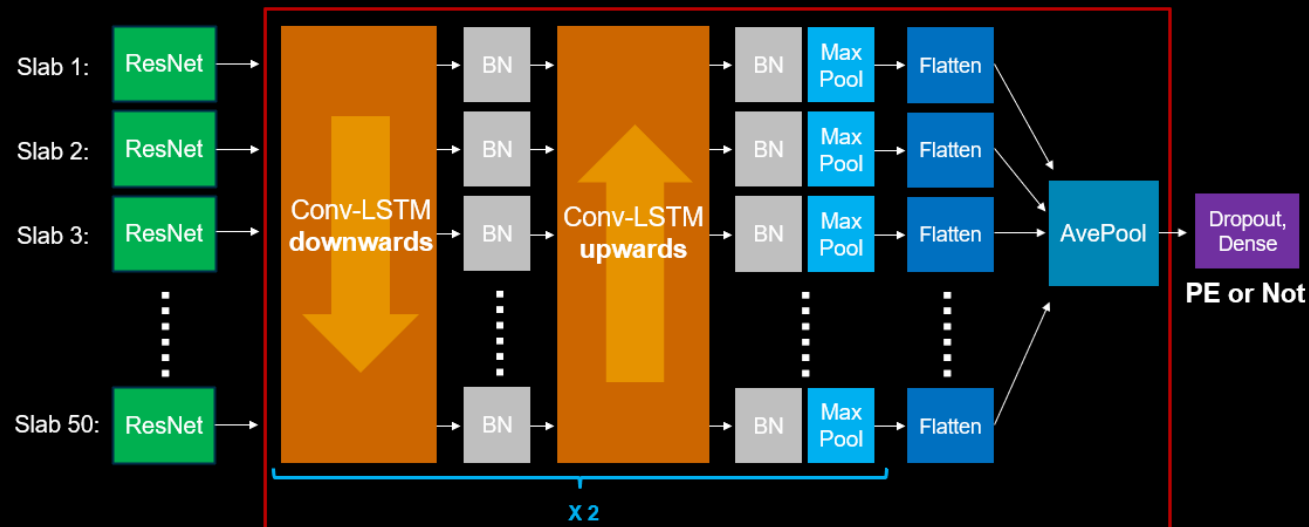
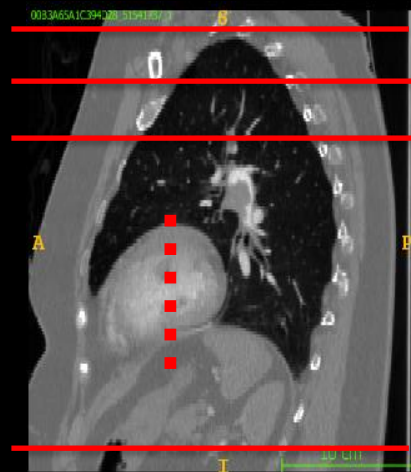


Stage2: training with patient-level labeled data

Recurrent Framework: Conv-LSTM



Stage2: training parameters



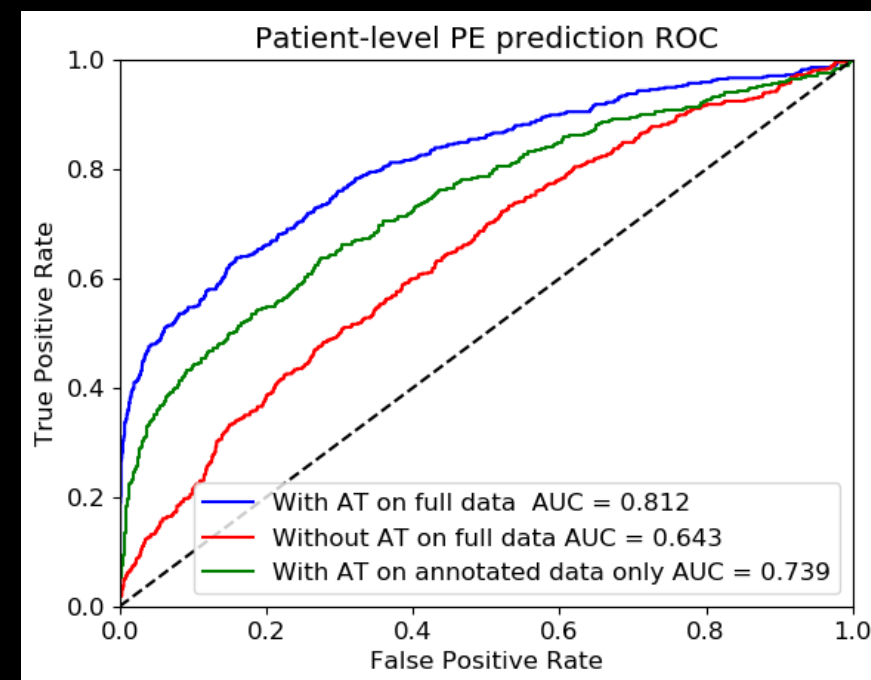
- Classification loss: Binary cross entropy (BCE)
- Optimizer: Adam optimizer
- Learning rate: 10^{-4}
- Training epochs: 50 (save the model with the highest val. acc.)



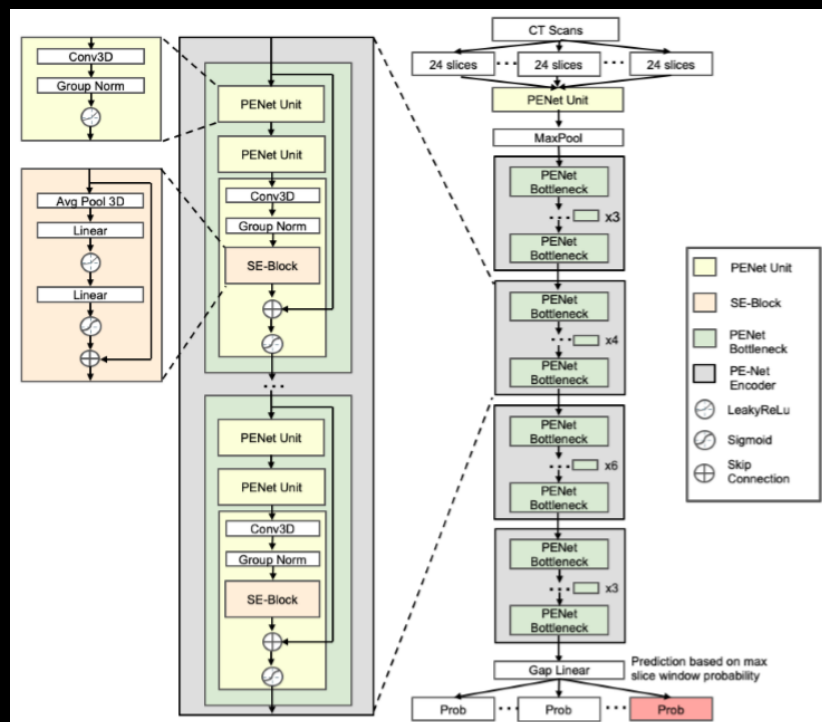
Stage2: patient-level inference results on testing data

- **Training data:**
 - Annotated Studies: 1670+, 593-
 - Labeled volumetric images: 4186+, 4603-
 - 80% training, 20% validation
- **Testing data (2160 total):** 517+, 1643-

	Stage 1 Data	Stage 1 Loss	Stage 2 Data	AUC
Scenario 1	1670+, 593-	Atten. Loss + Cls. Loss	1670+, 593-	0.739
Scenario 2	1670+, 593-	Cls. Loss	1670+, 593- & 4186+, 4603-	0.643
Scenario 3	1670+, 593-	Atten. Loss + Cls. Loss	1670+, 593- & 4186+, 4603-	0.812

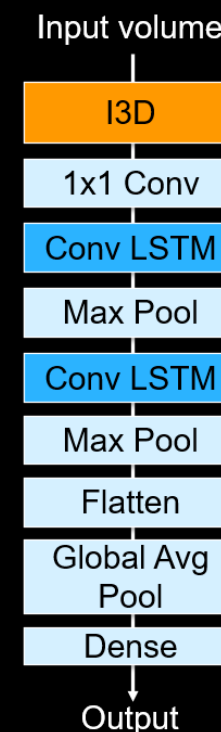


PENet



- Training data was labeled on a slice level for the presence/absence of a PE

3D CNN



- Starts with an I3D model (3D CNN pretrained on video action recognition dataset)
- Demonstrated success in acute aortic syndrome detection
- Trained only on our patient-level labeled PE data



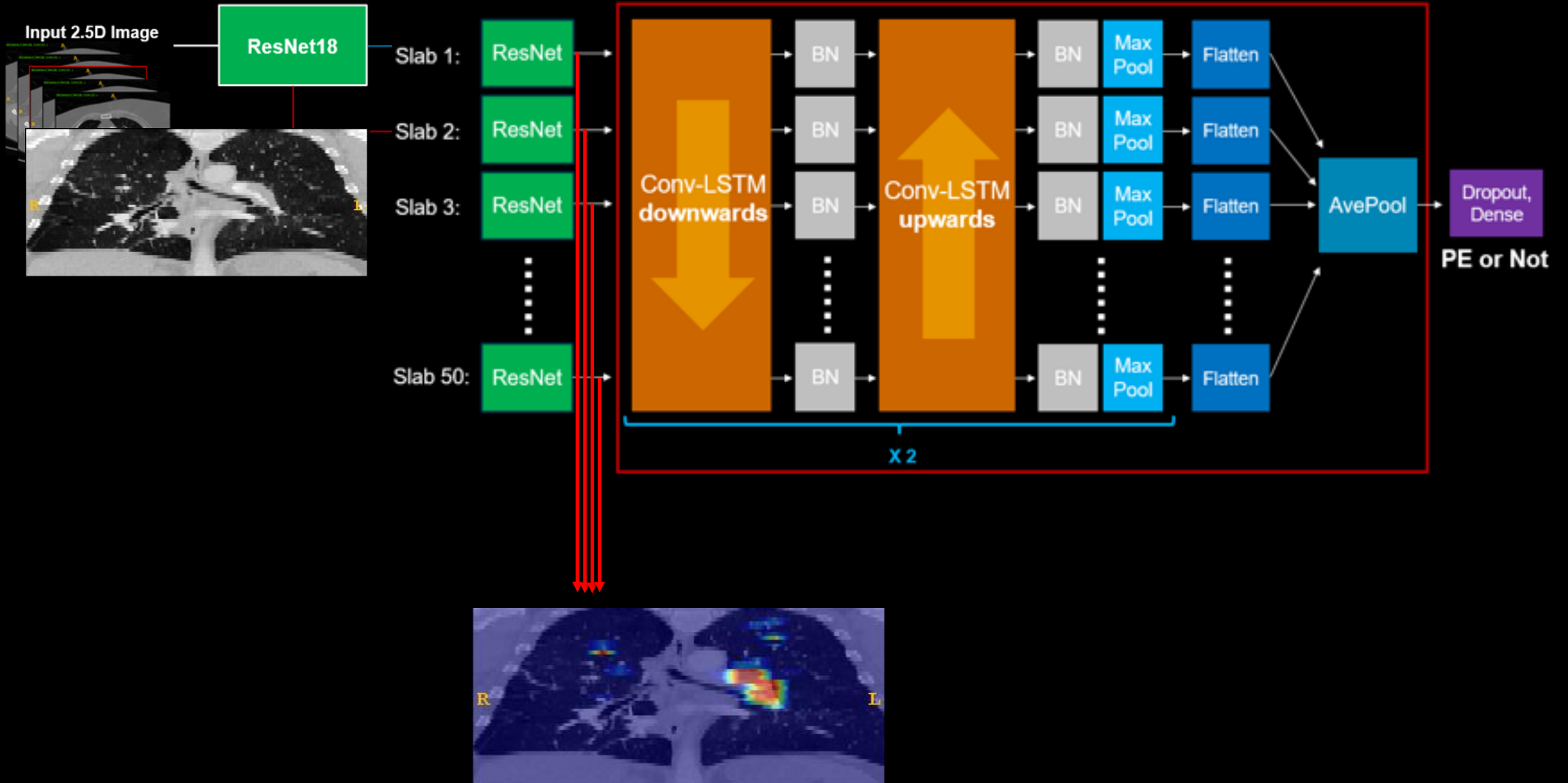
Approach	Testset description		AUC	Accuracy
	Size	Clinical sites		
PENet (int.)	198	Single	0.79	0.74
PENet (ext.)	227	Single	0.77	0.67
3D CNN	2160	Multiple	0.787	0.727
Proposed	2160	Multiple	0.812	0.781

Mixed protocols:

- Contrast-enhanced Chest CT vs. CT pulmonary angiogram
- Different dose levels (noise level)
- Different image reconstruction kernels
- Slice thickness: 0.5mm-5mm

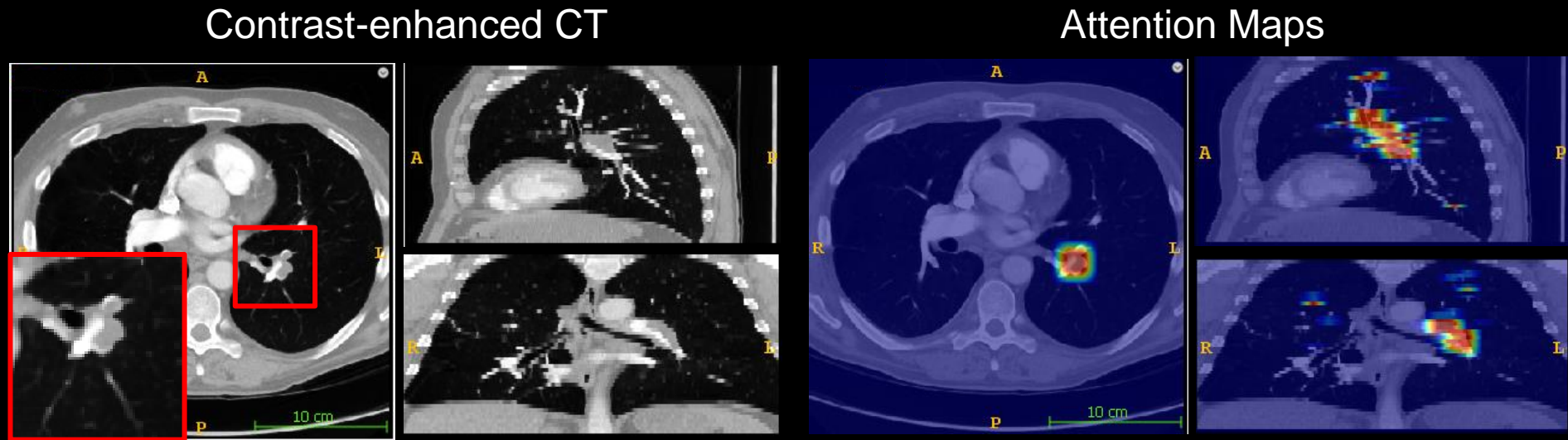


Auxiliary output – PE localization

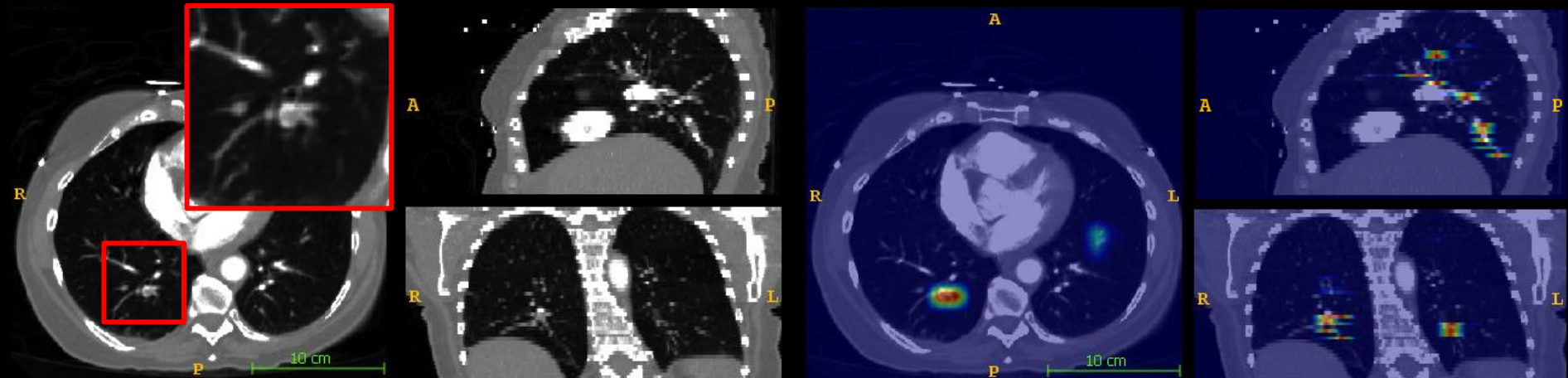


Auxiliary output – PE localization

Example 1



Example 2



- Using more efficient network structures (e.g. DenseNet) to replace ResNet18.
- In Stage1, the weights of classification loss and attention loss can be optimized (currently 1:1).
- Fully end-to-end training where the weights of ResNet can also be updated.



- We introduced a deep learning model to detect PE on volumetric contrast-enhanced CT scans using a 2-stage hybrid training strategy
 - Training with attention loss on pixel-level annotated data improves the network's localization ability
 - Second-stage convolution-LSTM networks reduce false positives on patient-level prediction
- Our evaluation involves the largest number of patient studies among all the research studies on automatic PE detection.
- Achieved state-of-the-art PE detection, while providing attention maps for radiologists as references.
- Applicable to other detection problems where the availability of volumetric imaging data exceeds radiologists' capacity to manually delineate ground truth.



Thank you!

Q&A