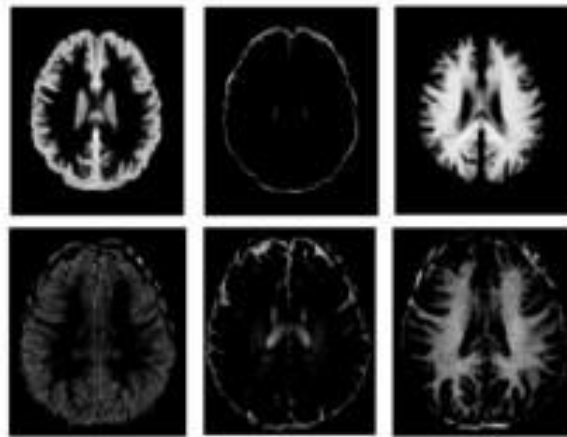


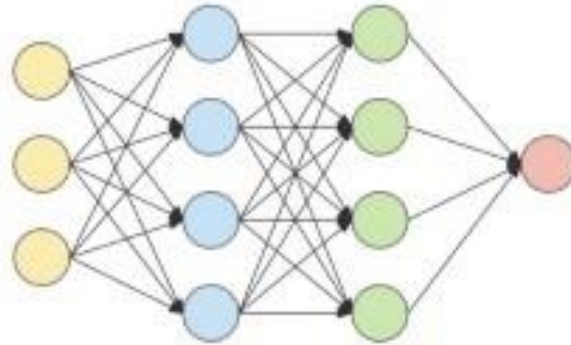
Wyjaśnialna sztuczna inteligencja

Weronika Hryniewska

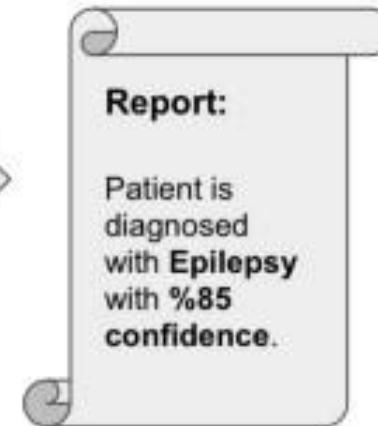
Epilepsy Detection Model with Brain MRI Data



Brain MRI data



Complex ML model

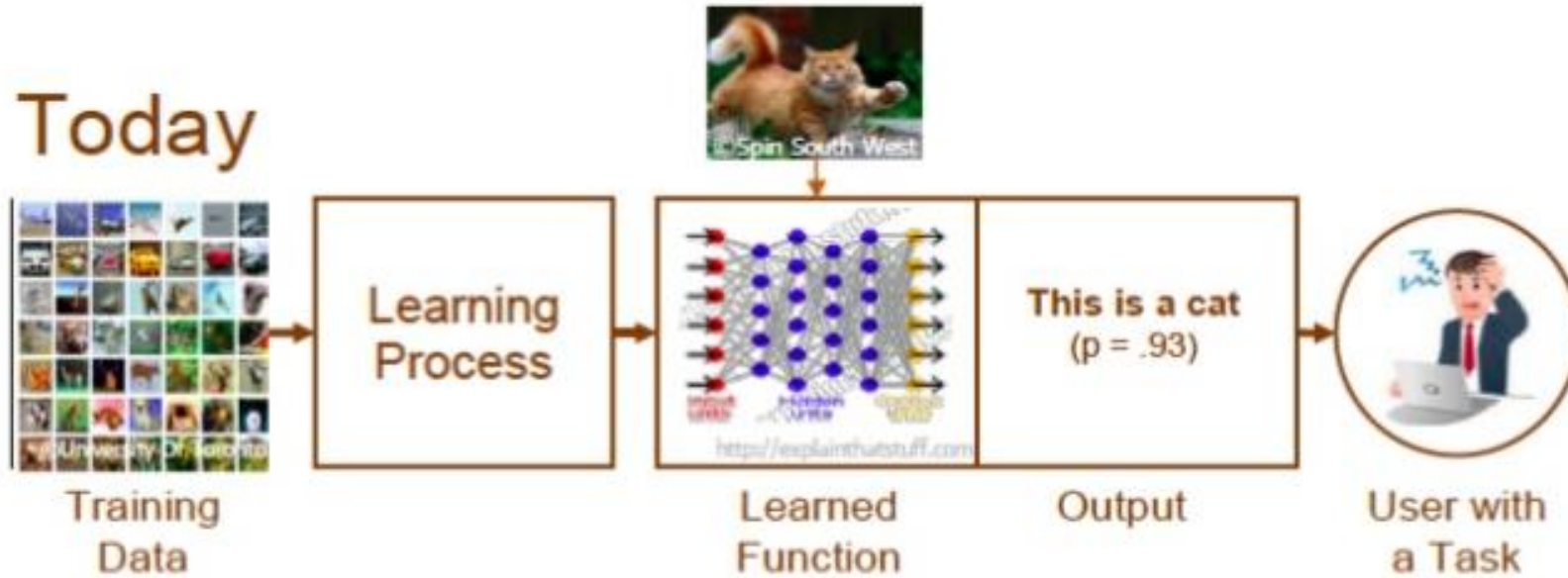


But why?!

Can I trust this prediction?

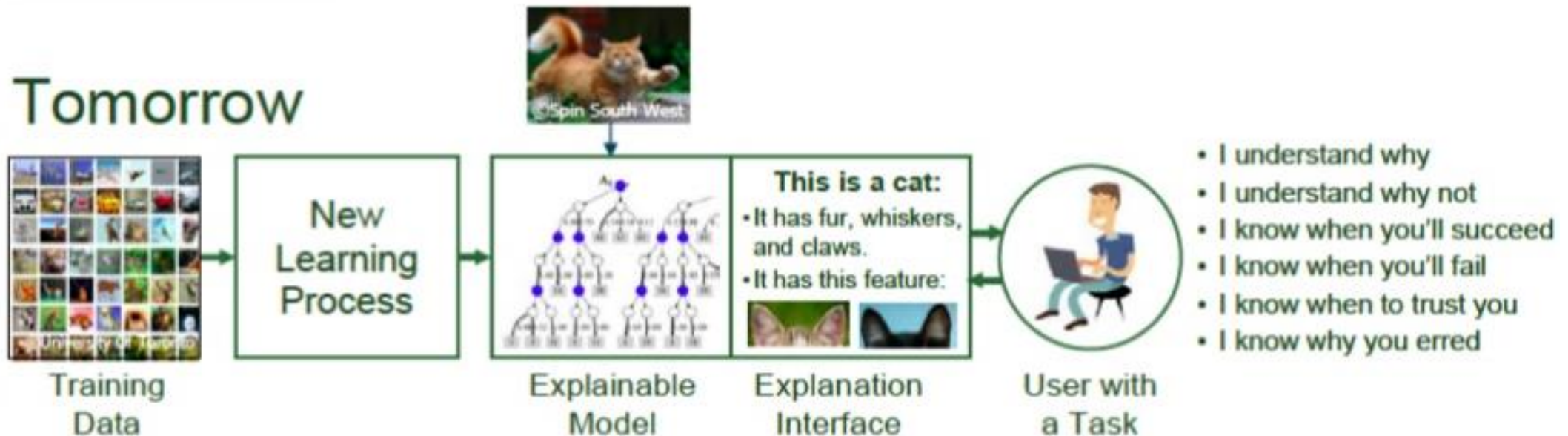
Example

Today

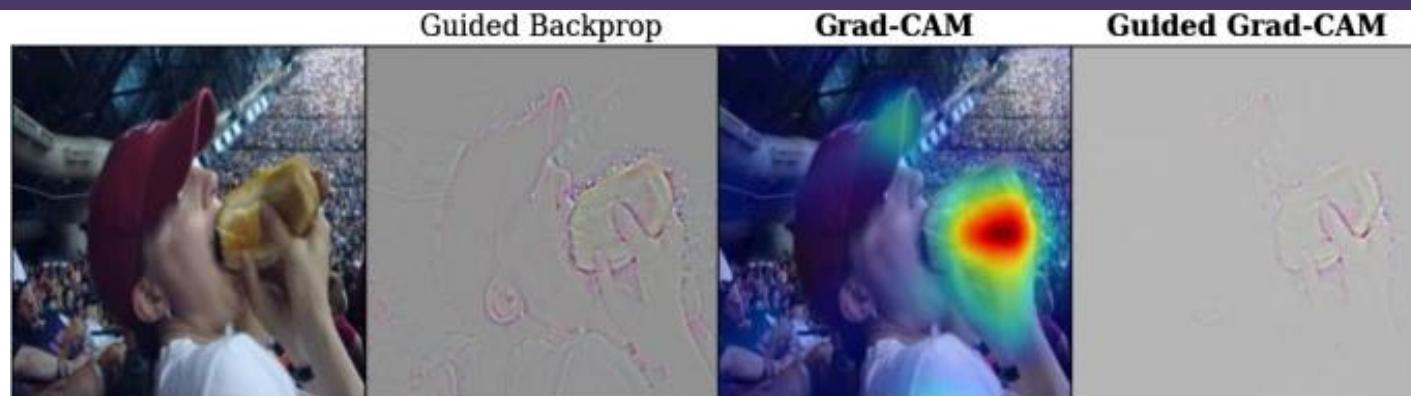


- Why did you do that?
- Why not something else?
- When do you succeed?
- When do you fail?
- When can I trust you?
- How do I correct an error?

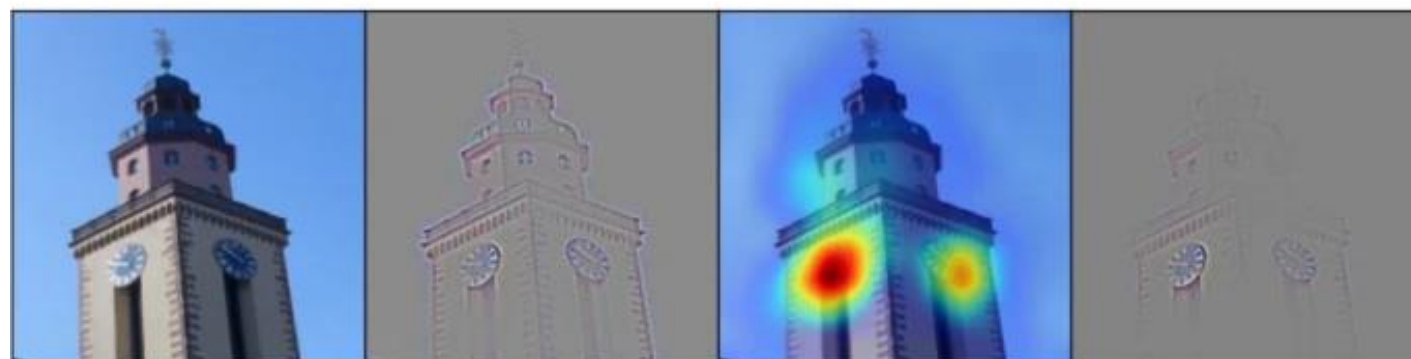
Tomorrow



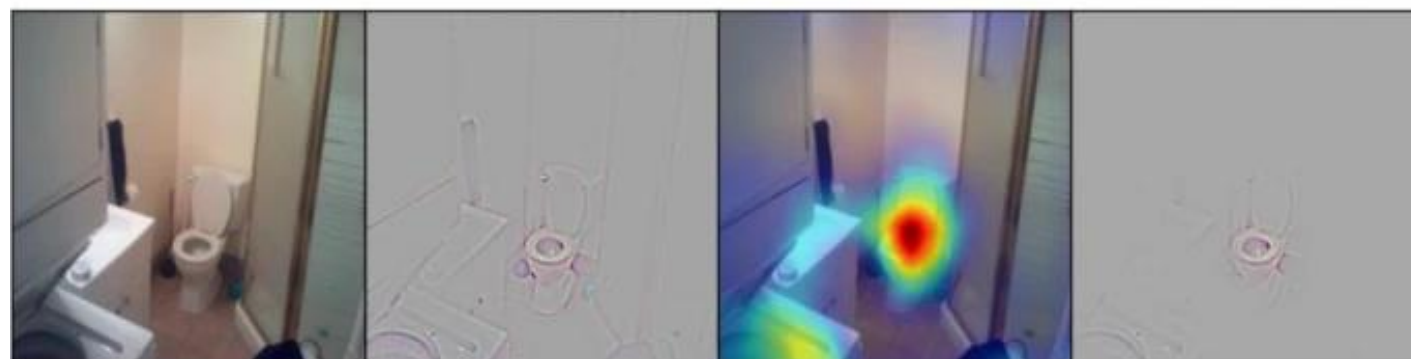
- I understand why
- I understand why not
- I know when you'll succeed
- I know when you'll fail
- I know when to trust you
- I know why you erred









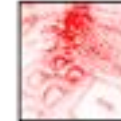

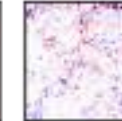
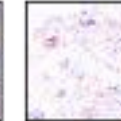
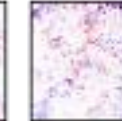

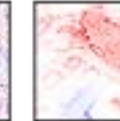
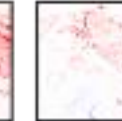









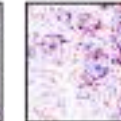

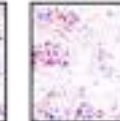

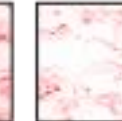








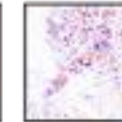
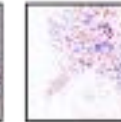
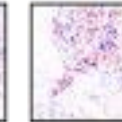
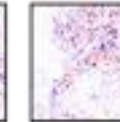
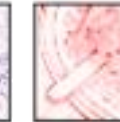
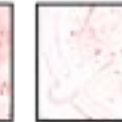









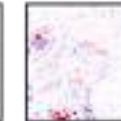

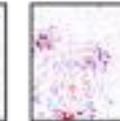














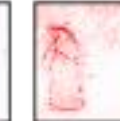









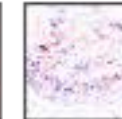
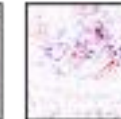
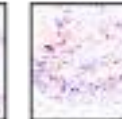
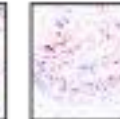
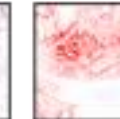
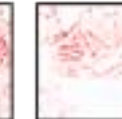
A man is holding a hot dog in his hand



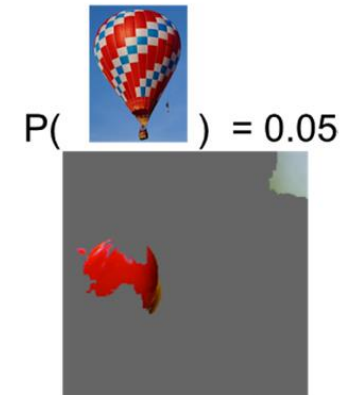
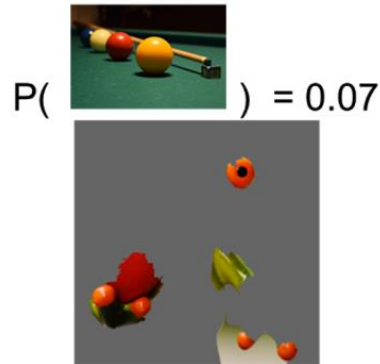
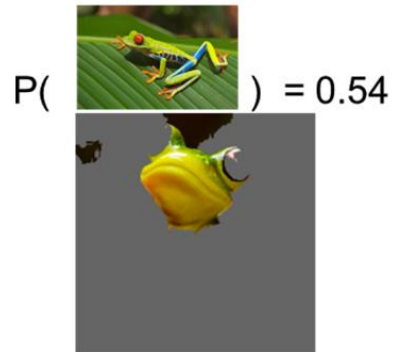
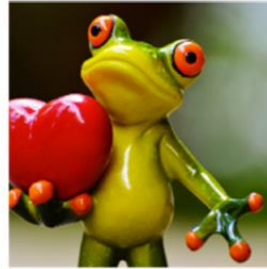
A large clock tower with a clock on the top of it



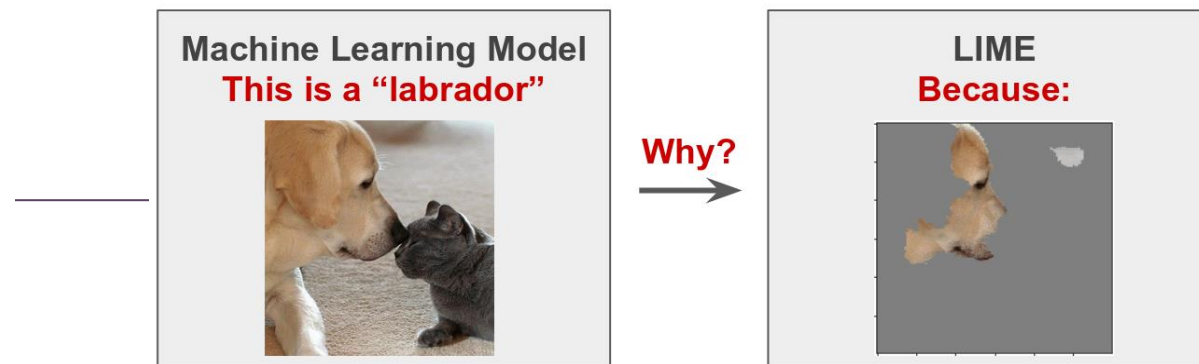
A bathroom with a toilet and a sink

	Input	Gradient	SmoothGrad	Deconvnet	Guided Backprop	PatternNet	PatternAttribution	DeepTaylor	Input * Gradient	Integrated Gradients	LRP-Z	LRP-Epsilon	LRP-PresetAFlat	LRP-PresetBFlat	
label: baseball pred: crayfish															logit: 9.90 prob: 0.18
label: bell pepper pred: bell pepper															logit: 20.36 prob: 0.98
label: ice lolly pred: ice cream															logit: 12.75 prob: 0.34
label: broom pred: broom															logit: 15.65 prob: 0.71
label: abaya pred: cloak															logit: 11.07 prob: 0.33
label: Dungeness crab pred: Dungeness crab															logit: 12.39 prob: 0.39

Local Interpretable Model-Agnostic Explanations (LIME)

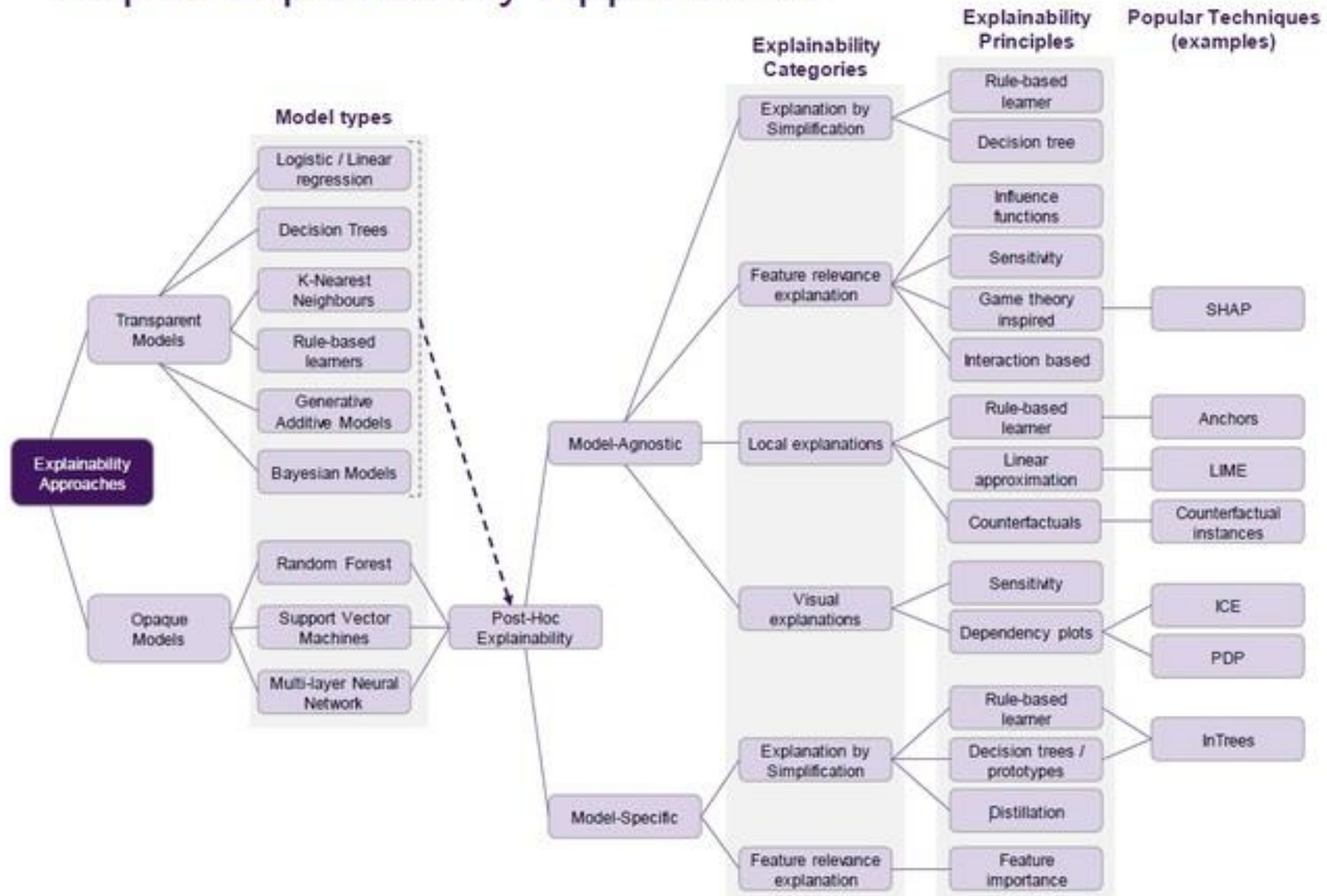


www.oreilly.com/content/introduction-to-local-interpretable-model-agnostic-explanations-lime



<https://towardsdatascience.com/interpretable-machine-learning-for-image-classification-with-lime-ea947e82ca13>

Map of Explainability Approaches

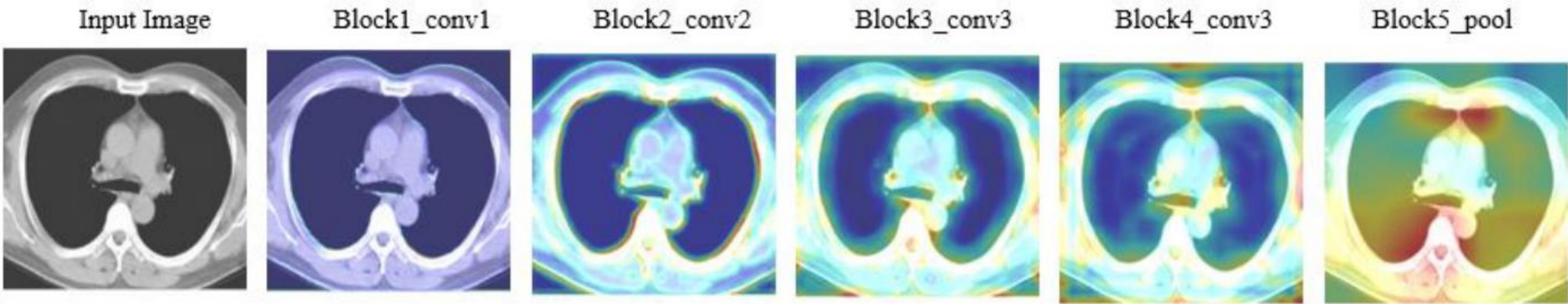


Przykłady metod XAI dla zdjęć płuc chorych na Covid-19

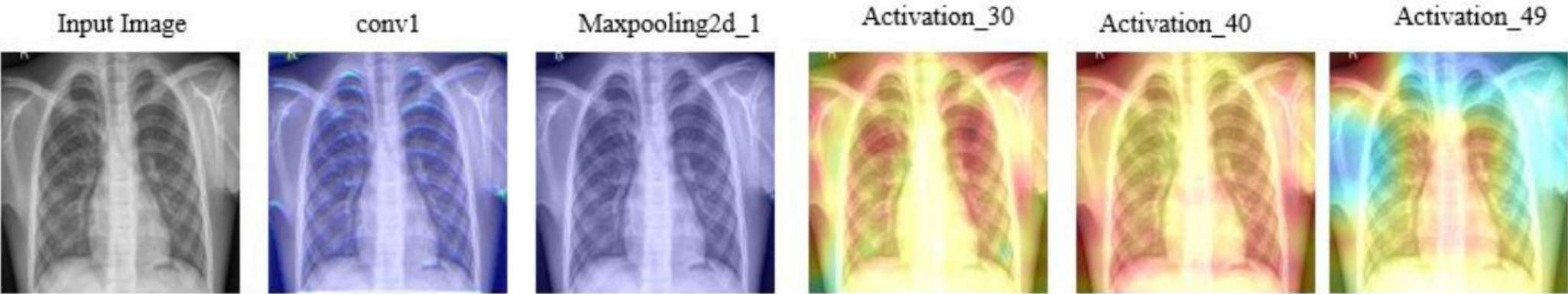
Weronika Hryniewska

Abbas, A., Abdelsamea, M. M., & Gaber, M. M. (2020). *Classification of COVID-19 in chest X-ray images using DeTraC deep convolutional neural network*. <http://arxiv.org/abs/2003.13815>

Heatmap of class activation of other patient's CT scan image on different layer acquired by VGG16



Heatmap of class activation of other patient's chest x-ray image on different layer by ResNet50



Are such distances from the marked areas acceptable?

With red arrow the radiologists marked the areas where the COVID-19 disease is manifested, as is shown from the activation maps, the proposed models rightly highlighted the areas marked by the radiologist.

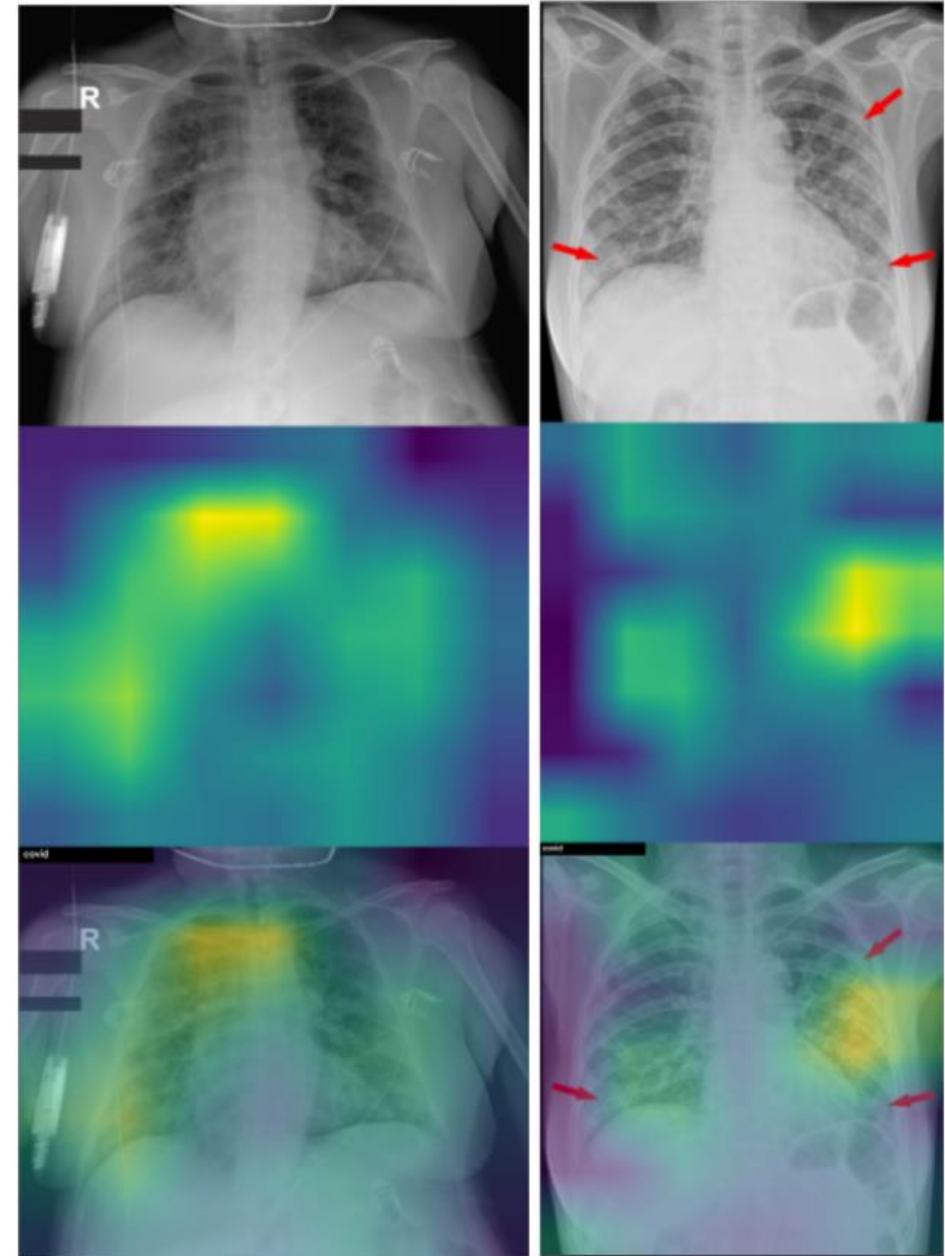
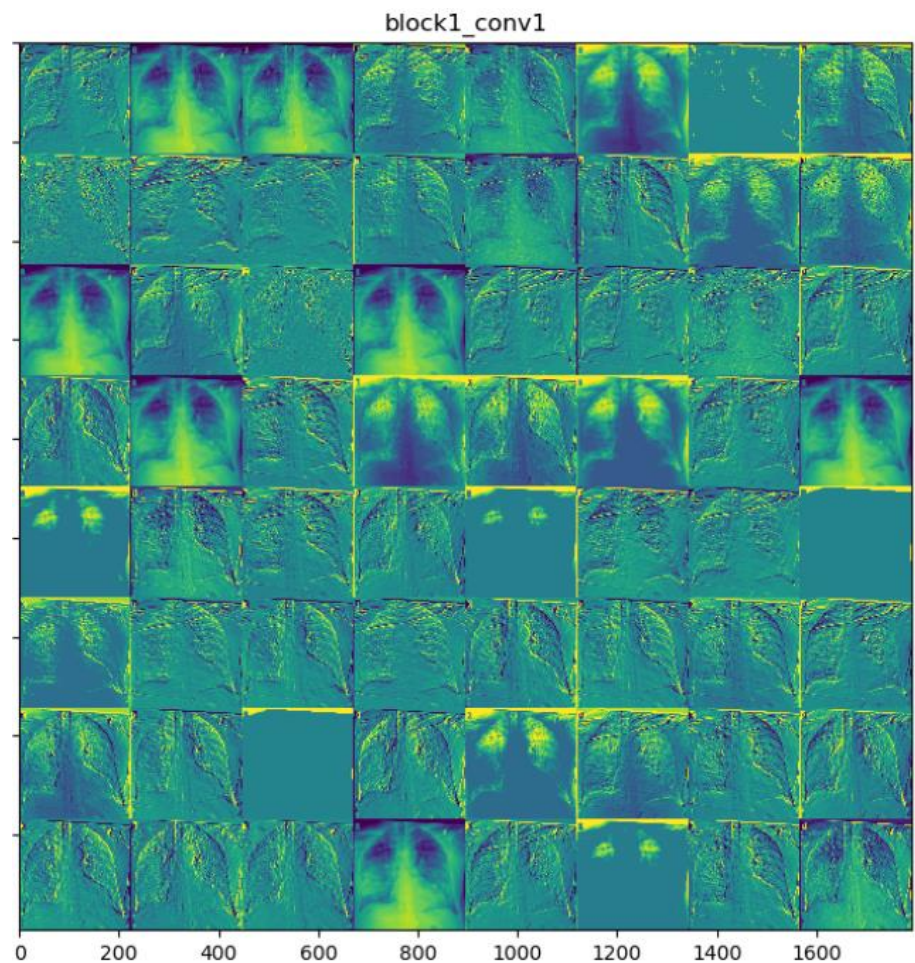
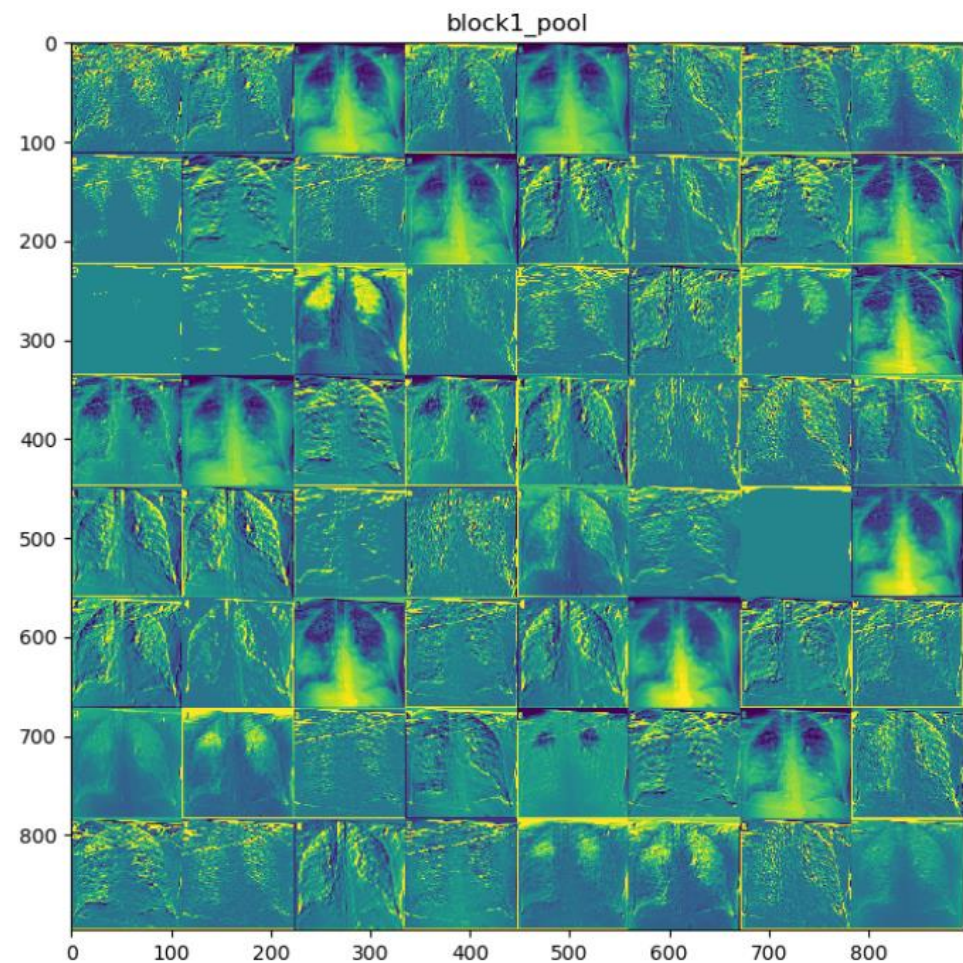


Fig. 9. Examples of COVID-19 model activation maps.



(a) Feature Map at Convolution Layer 1



(b) Feature Map at pooling Layer

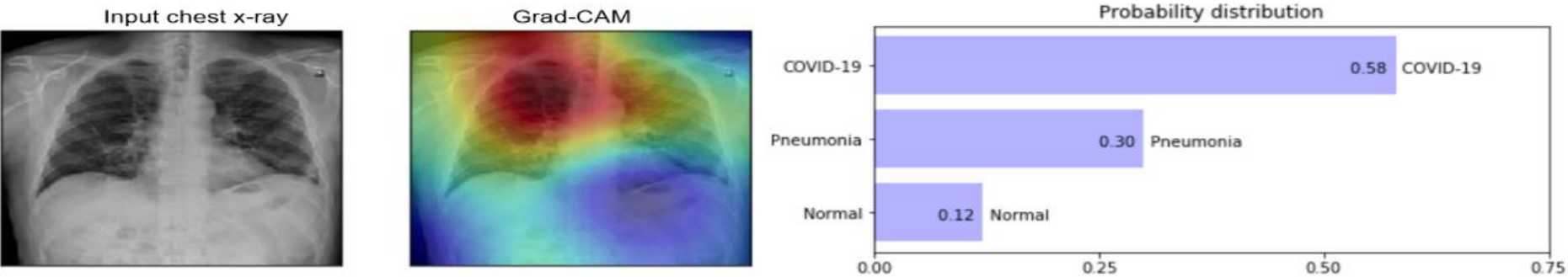


Figure 7: The input chest x-ray classification, decision visualization with Grad-CAM and explanation

Are such big areas relevant?

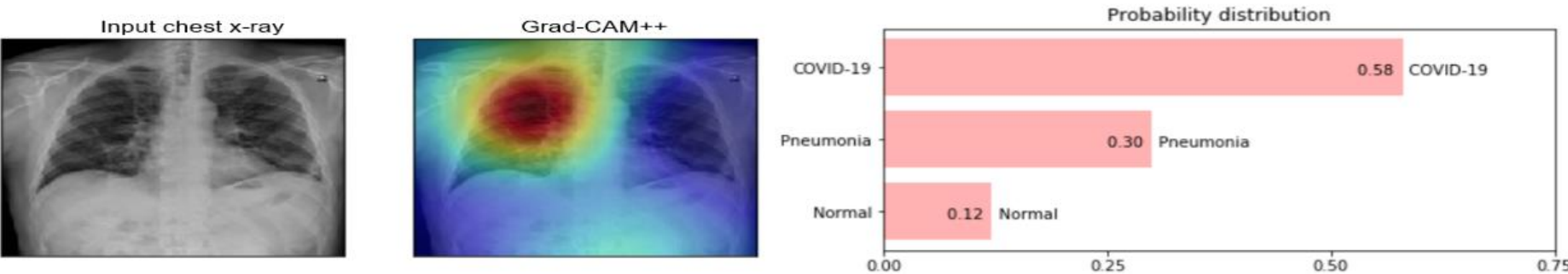


Figure 8: The input chest x-ray classification, decision visualization with Grad-CAM++ and explanation

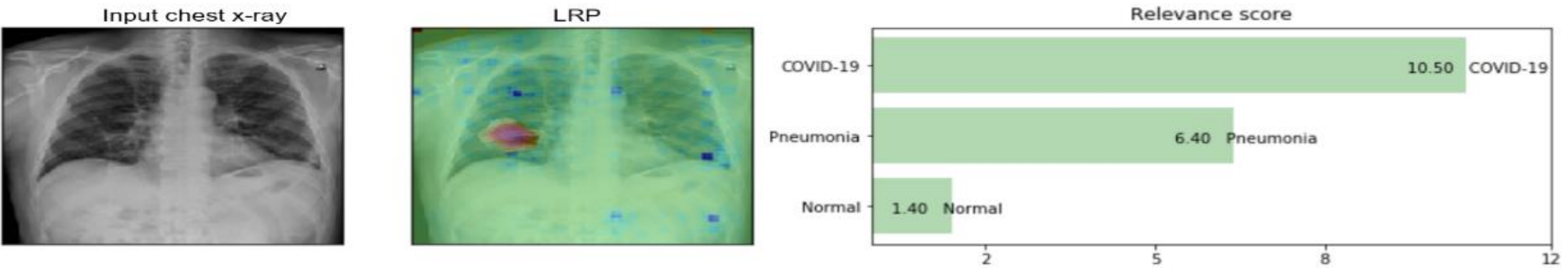
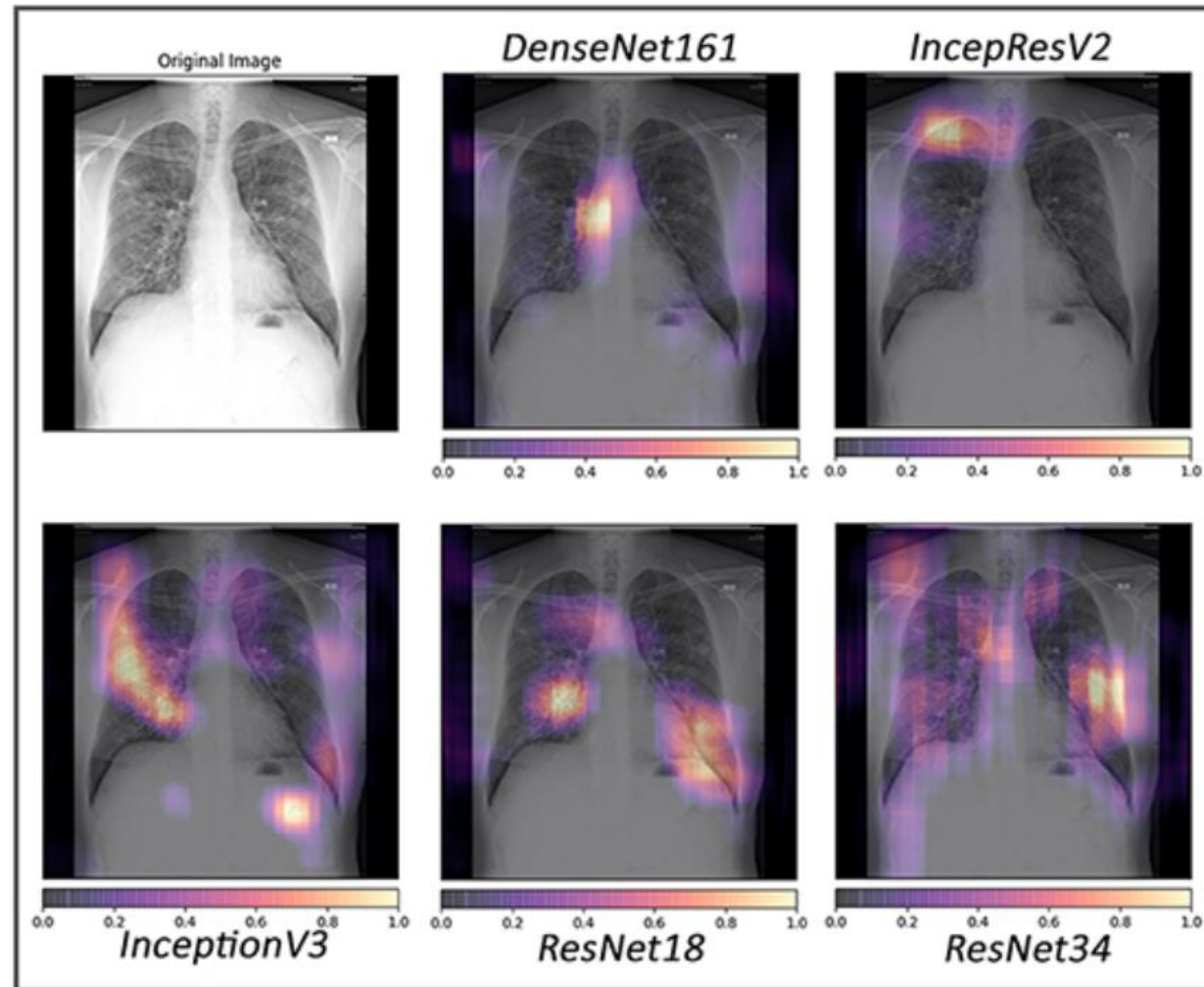


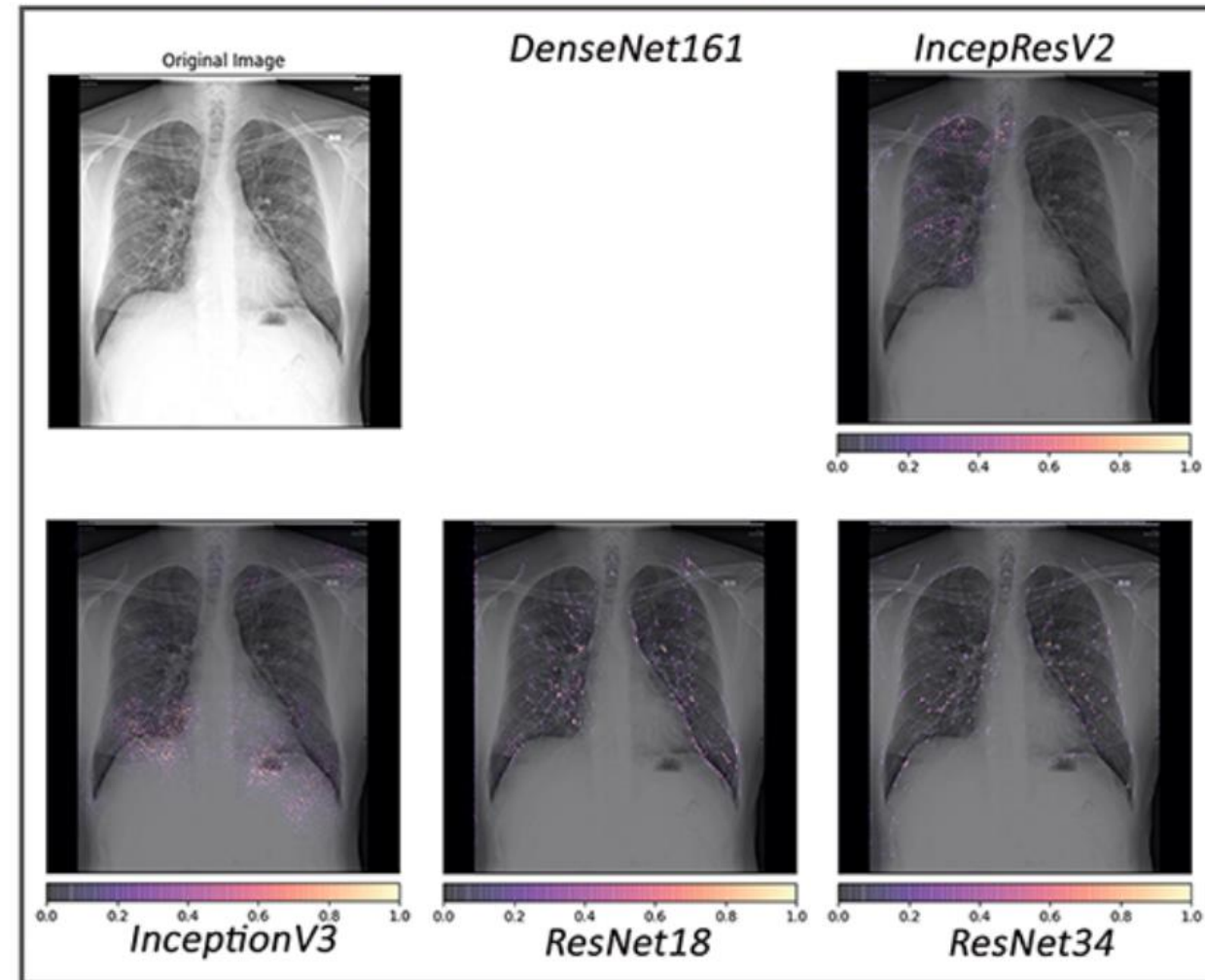
Figure 9: The input chest x-ray classification, decision visualization with LRP and explanation

Which visualization method is most useful?

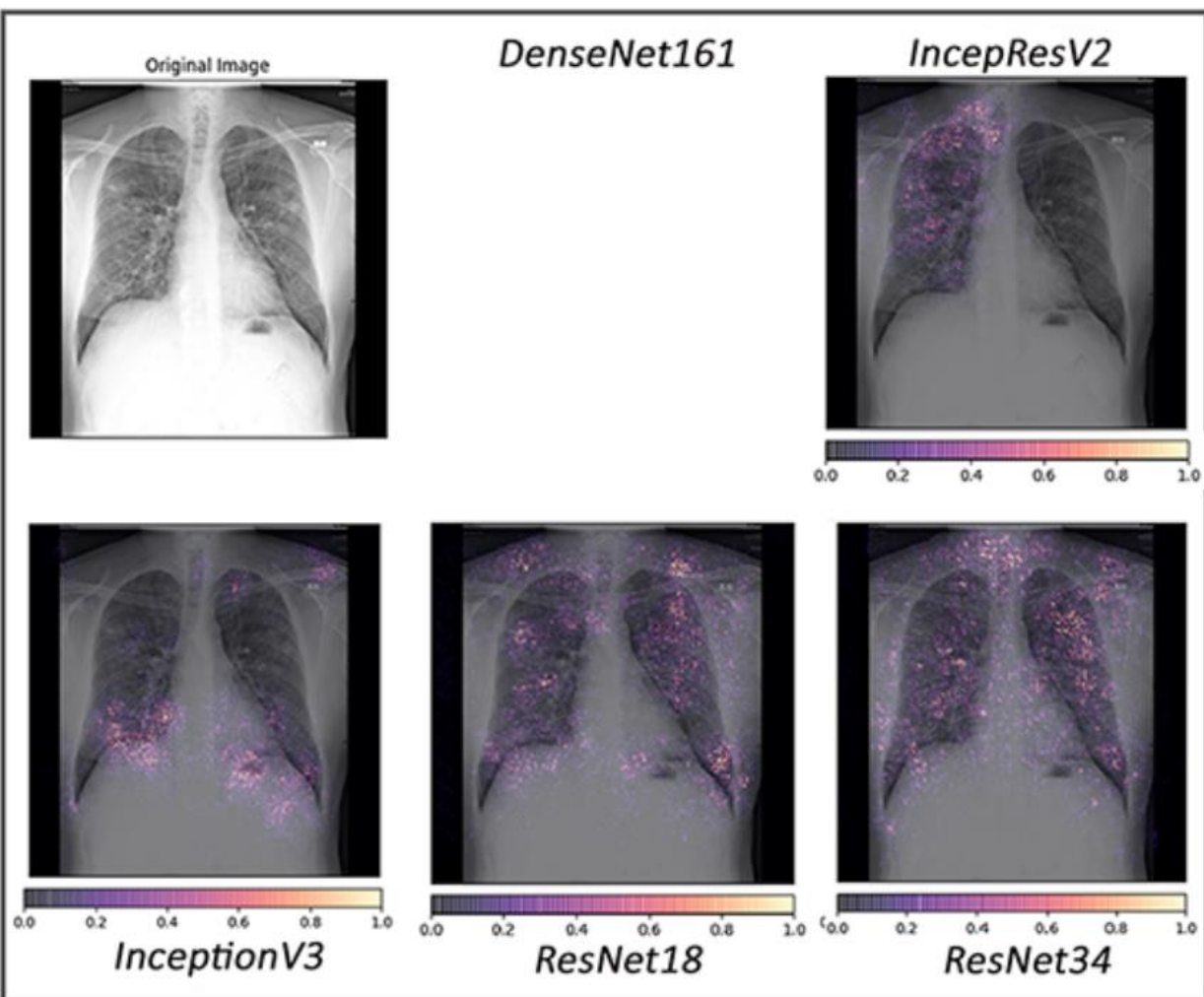
Occlusion



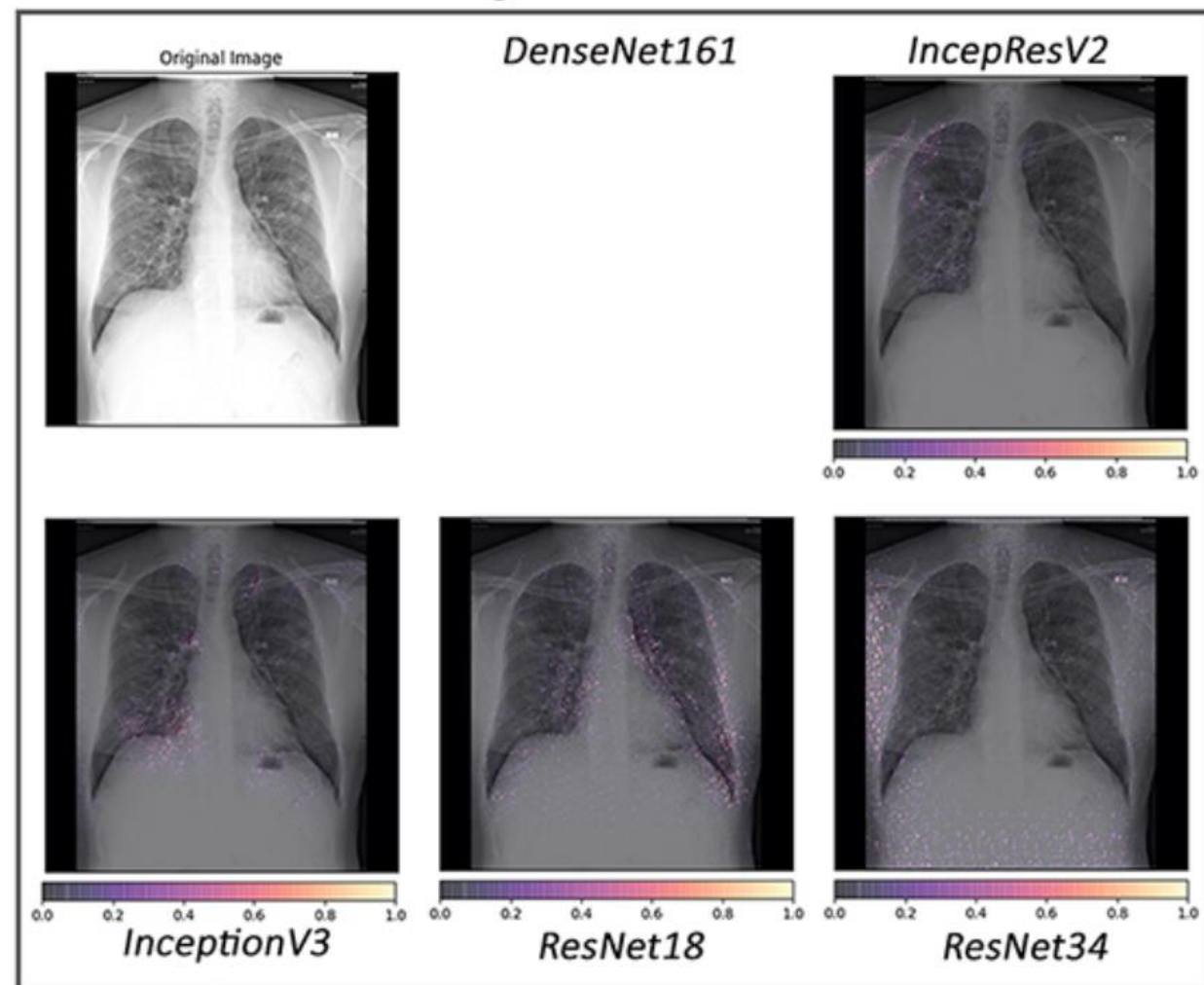
Guided backpropagation



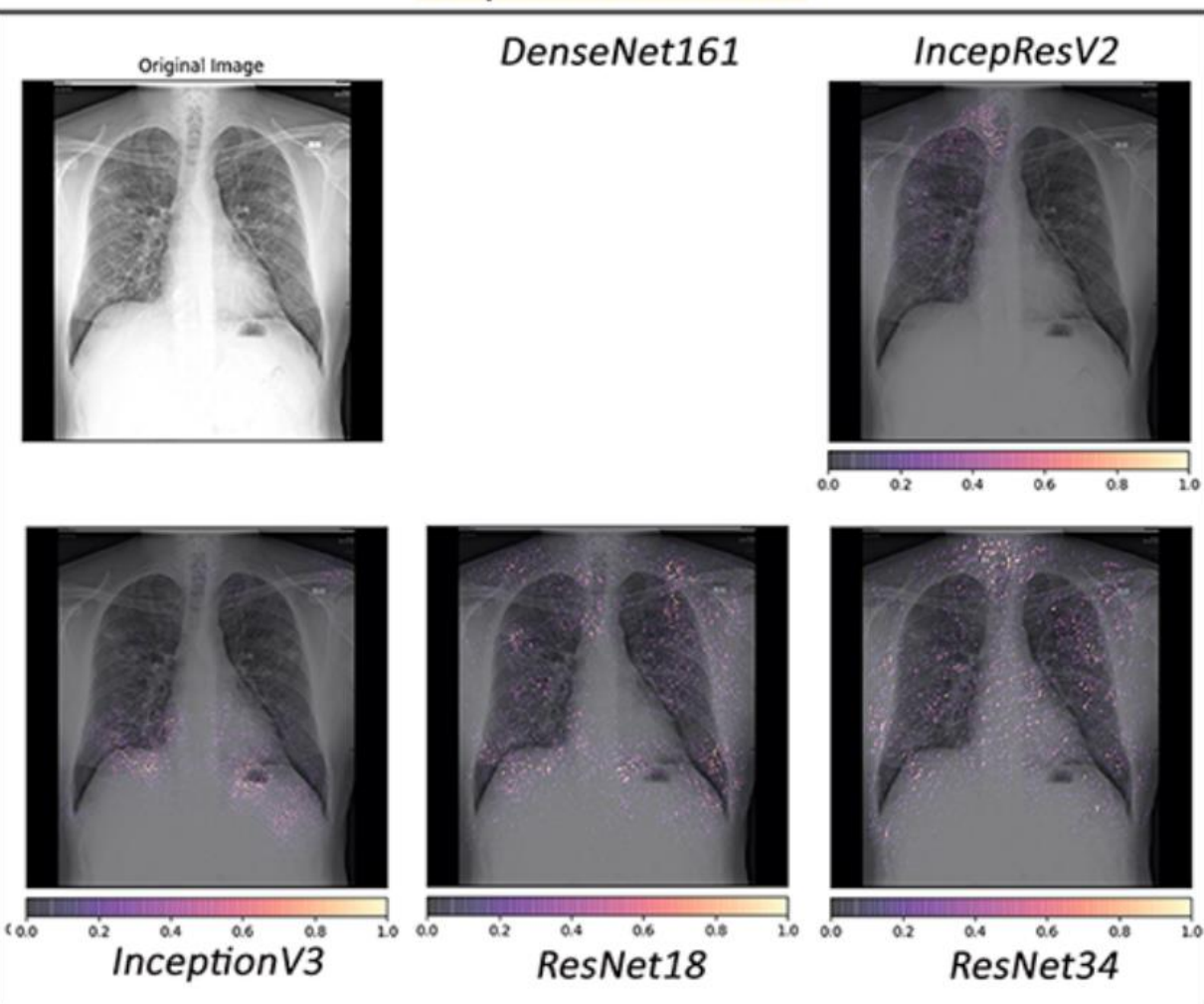
Saliency



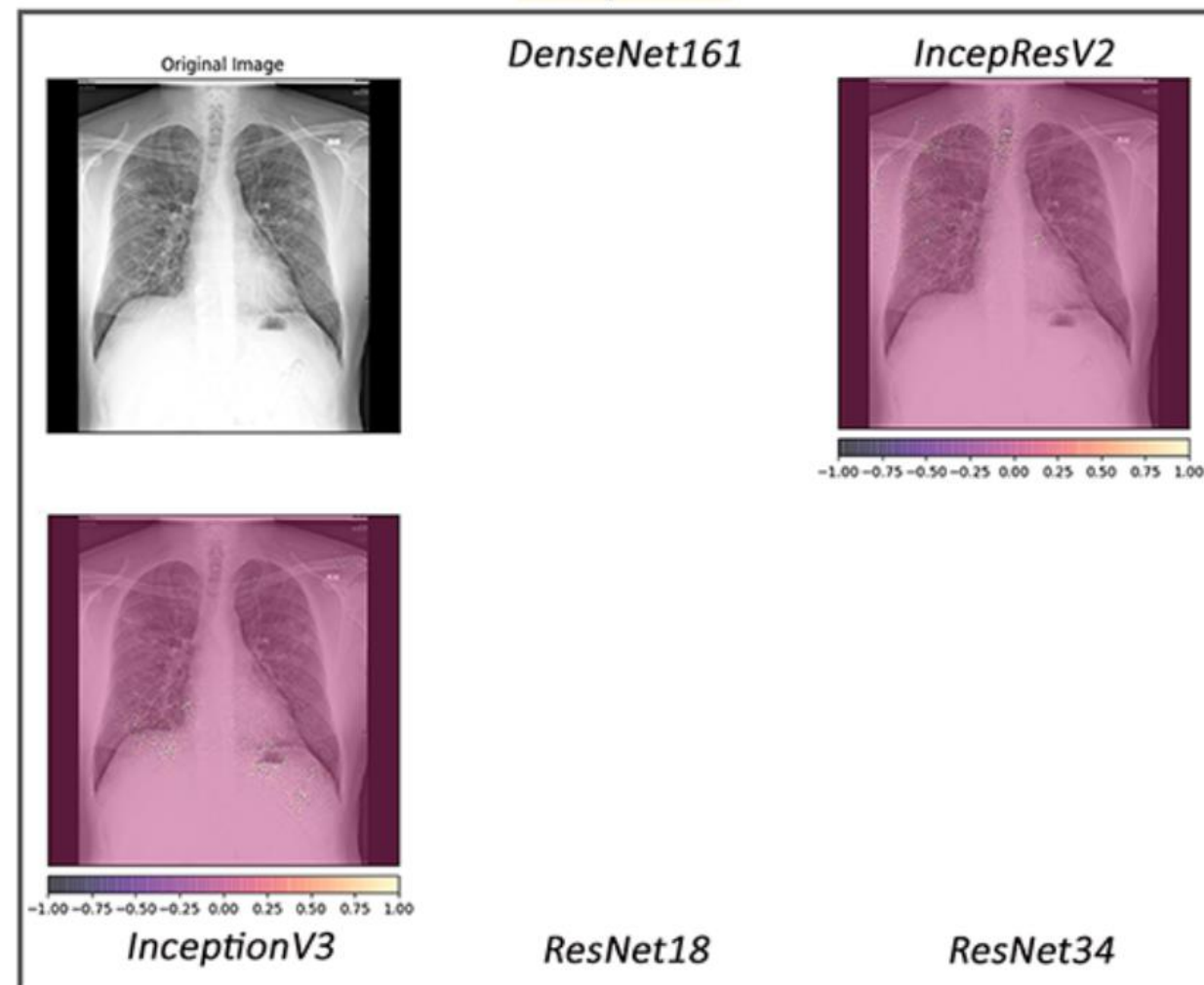
Integrated Gradients



Input X Gradient



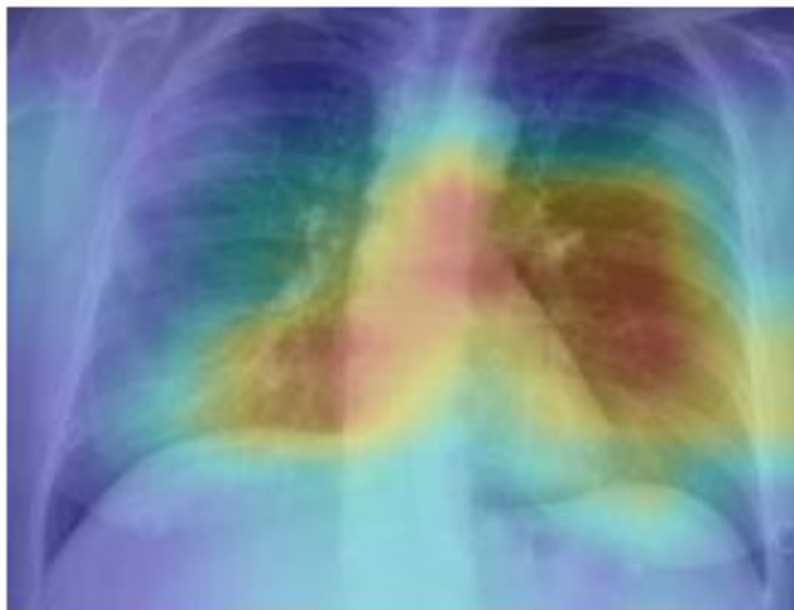
DeepLIFT



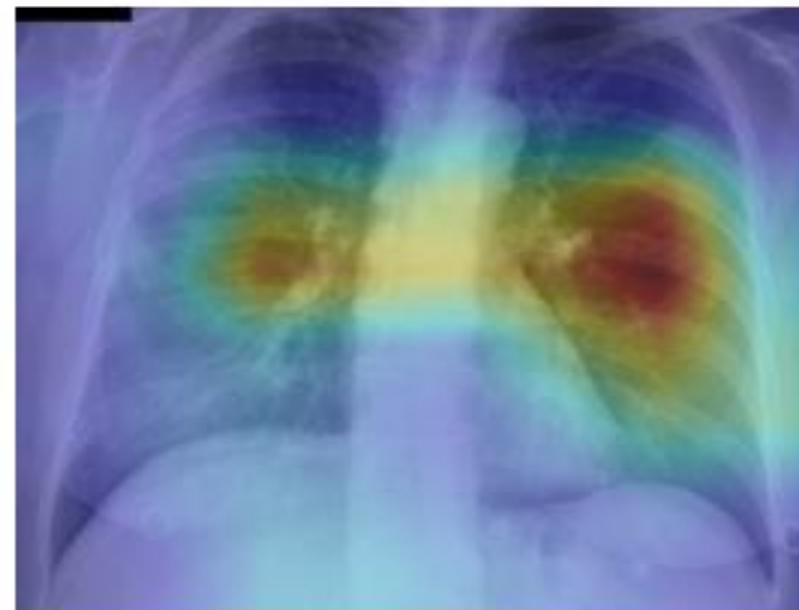
Is the area of changes related to the severity of the disease?



(a) Original Positive (Mild)



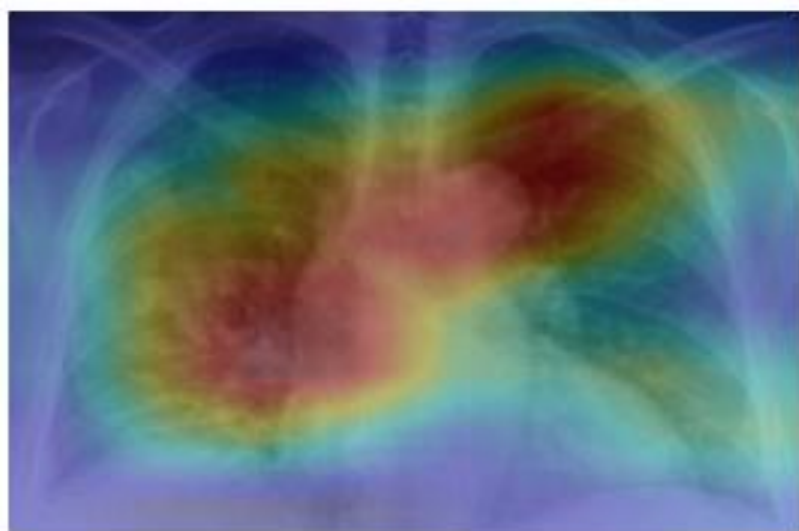
(b) why positive



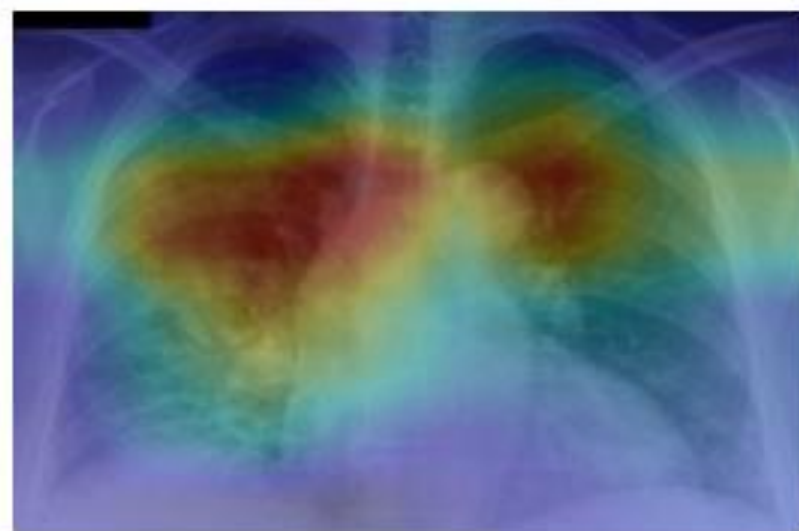
(c) why negative



(a) Original Positive (Moderate)



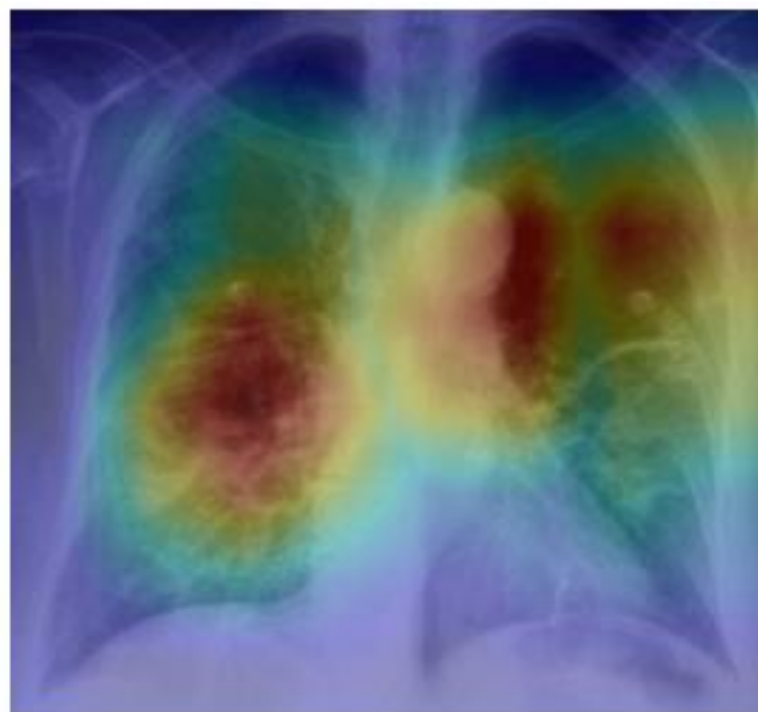
(b) why positive



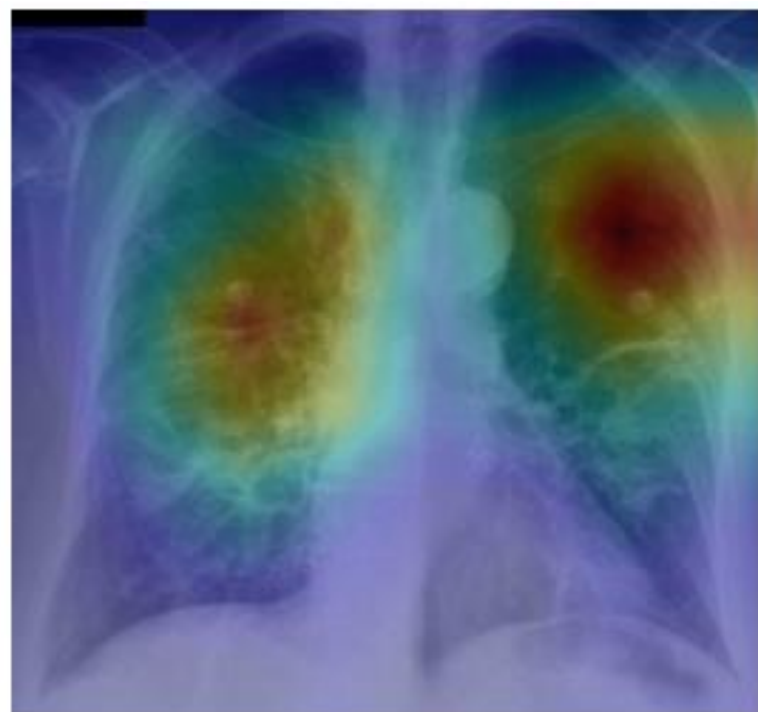
(c) why negative



(a) Original Positive (Severe)



(b) why positive



(c) why negative



(a) Original Negative



(b) why positive



(c) why negative

Ouyang, X., Huo, J., Xia, L., Shan, F., Liu, J., Mo, Z., Yan, F., Ding, Z., Yang, Q., Song, B., Shi, F., Yuan, H., Wei, Y., Cao, X., Gao, Y., Wu, D., Wang, Q., & Shen, D. (2020). Dual-Sampling Attention Network for Diagnosis of COVID-19 from Community Acquired Pneumonia. *IEEE Transactions on Medical Imaging*, 39(XX), 1–1. <https://doi.org/10.1109/tmi.2020.2995508>

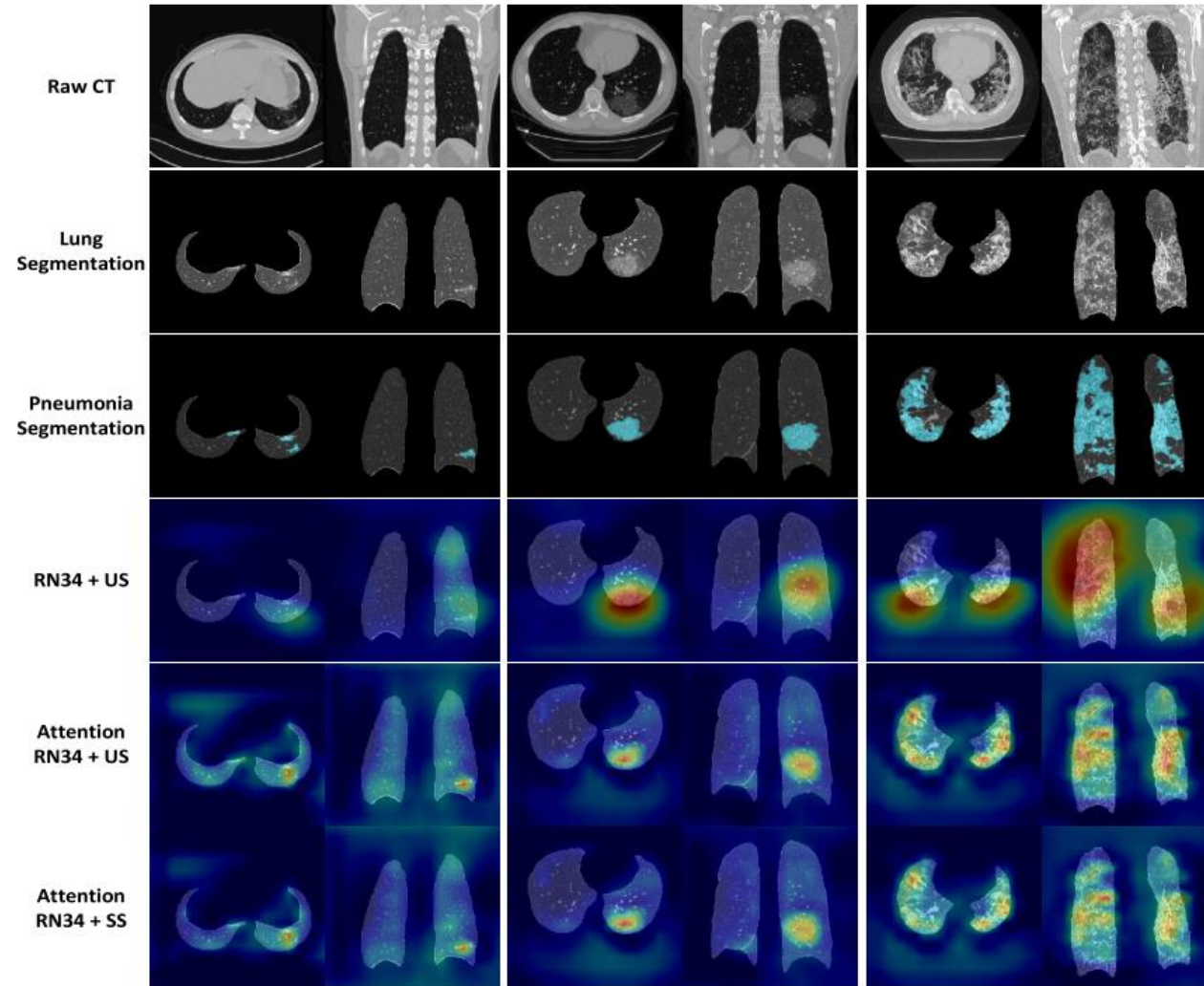


Fig. 5. Visualization results of our methods on three COVID-19 cases from small-infection group (< 0.005), median-infection group ($0.005 - 0.030$) and large-infection group (> 0.030) of the test set are shown from left to right, respectively. For each case, we show the visualization results in both axial view and coronal view. We show the original images (1st row), and the segmentation results of the lung and pneumonia infection regions (2nd and 3rd rows) by the VB-Net toolkit [10]. For the attention results, we show the Grad-CAM results of “RN34 + US” (4th row), and the attention maps obtained by our proposed attention module of “Attention RN34 + US” and “Attention RN34 + SS” models (5th and 6th rows).

Cohen, J. P., Dao, L., Morrison, P., Roth, K., Bengio, Y., Shen, B., Abbasi, A., Hoshmand-Kochi, M., Ghassemi, M., Li, H., & Duong, T. Q. (2020). *Predicting COVID-19 Pneumonia Severity on Chest X-ray with Deep Learning*. 8(December 2019). <http://arxiv.org/abs/2005.11856>

The extent of lung involvement by ground glass opacity or consolidation for each lung (right lung and left lung separately) was scored as:

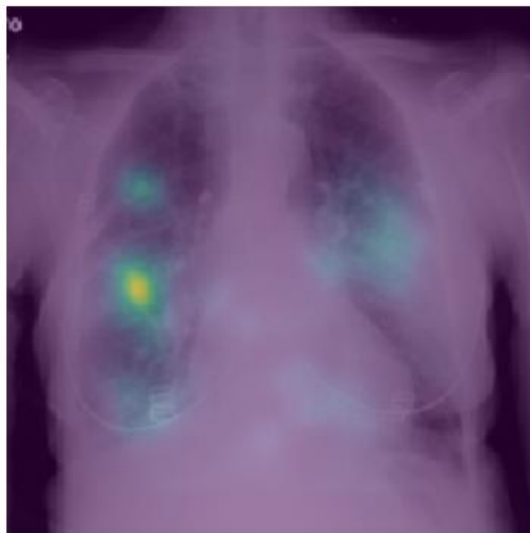
- 0 = no involvement;
- 1 = <25% involvement;
- 2 = 25-50% involvement;
- 3 = 50- 75% involvement;
- 4 = >75% involvement.

The total extent score ranged from 0 to 8 (right lung and left lung together).

The degree of opacity for each lung (right lung and left lung separately) was scored as:

- 0 = no opacity;
- 1 = ground glass opacity;
- 2 = consolidation;
- 3 = white-out.

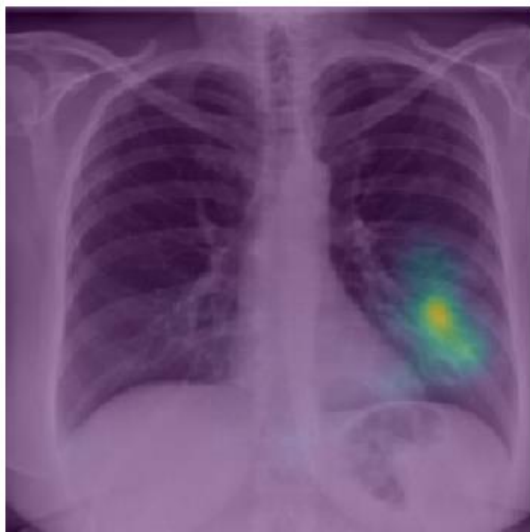
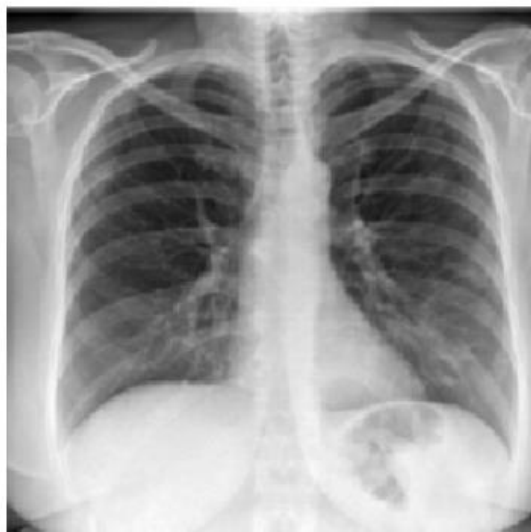
The total opacity score ranged from 0 to 6 (right lung and left lung together).



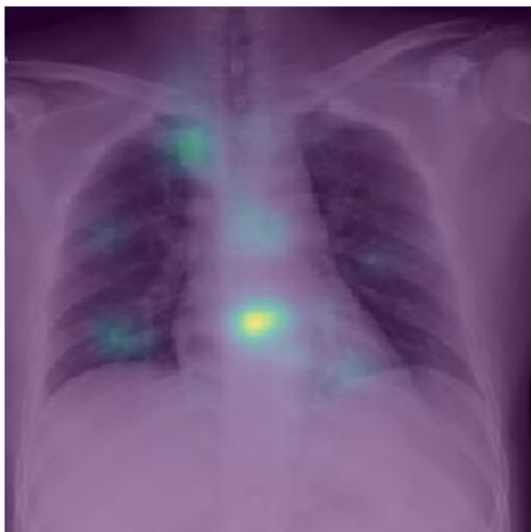
(a) Geographic Extent Score: 5, Predicted: 5.3



(b) Geographic Extent Score: 0, Predicted: -0.8



(c) Geographic Extent Score: 2, Predicted: 0.62



(d) Geographic Extent Score: 0, Predicted: 1.05

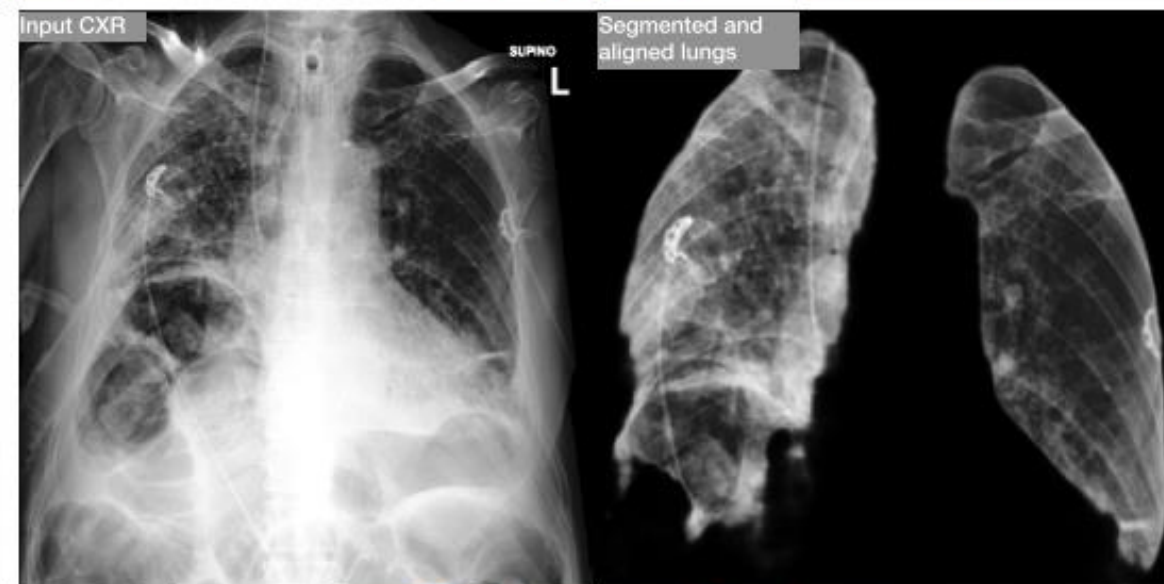
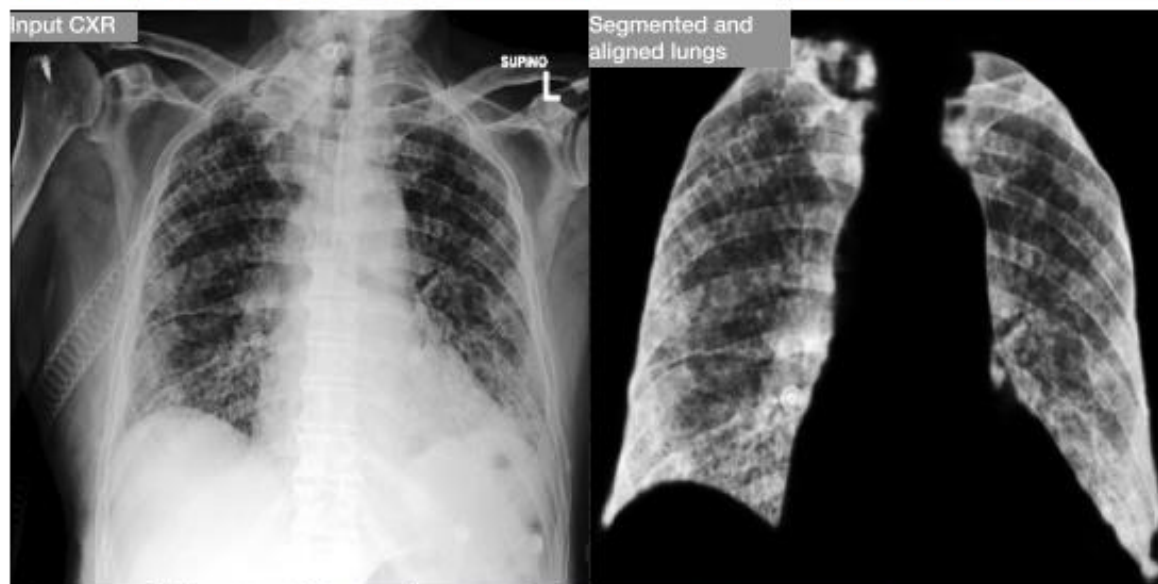
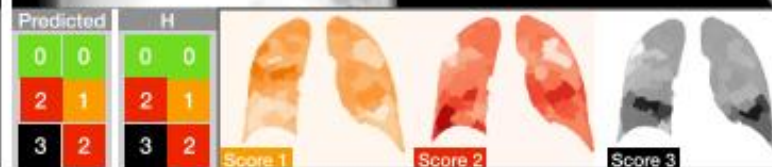
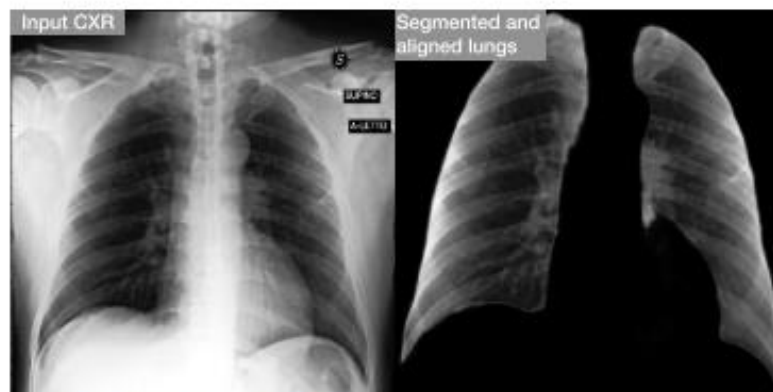
Signoroni, A., Savardi, M., Benini, S., Adami, N., Leonardi, R., Gibellini, P., Vaccher, F., Ravanelli, M., Borghesi, A., Maroldi, R., & Farina, D. (2020). *End-to-end learning for semiquantitative rating of COVID-19 severity on Chest X-rays*. 1–22.
<http://arxiv.org/abs/2006.04603>

According to it, lungs in anteroposterior (AP) or posteroanterior (PA) views, are subdivided into six zones, three for each lung:

- Upper zones (A and D): above the inferior wall of the aortic arch;
- Middle zones (B and E): below the inferior wall of the aortic arch and above the inferior wall of the right inferior pulmonary vein (i.e., the hilar structures);
- Lower zones (C and F): below the inferior wall of the right inferior pulmonary vein (i.e., the lung bases).

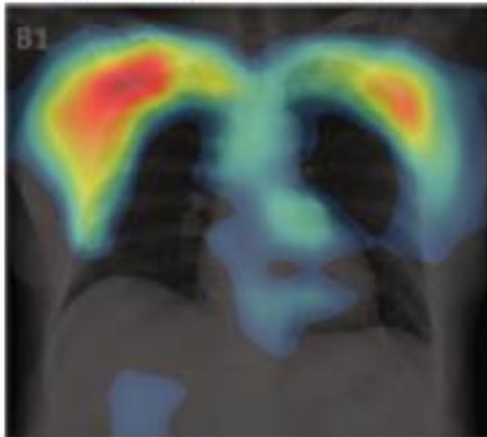
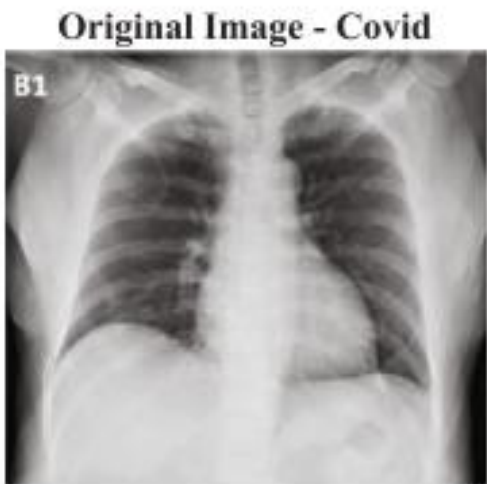
For each zone, a score (ranging from 0 to 3) is assigned, based on the detected lung abnormalities:

- 0: no lung abnormalities
- 1: interstitial infiltrates
- 2: interstitial (dominant), and alveolar infiltrates
- 3: interstitial, and alveolar (dominant) infiltrate



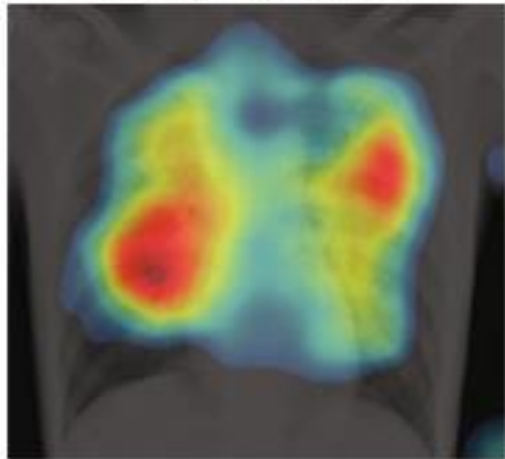
Ucar, F., & Korkmaz, D. (2020). COVIDiagnosis-Net: Deep Bayes-SqueezeNet based diagnosis of the coronavirus disease 2019 (COVID-19) from X-ray images. *Medical Hypotheses*, 140(April), 109761. <https://doi.org/10.1016/j.mehy.2020.109761>

Covid, 0.99999
Normal, 4.6985e-06
Pneumonia, 5.4529e-07



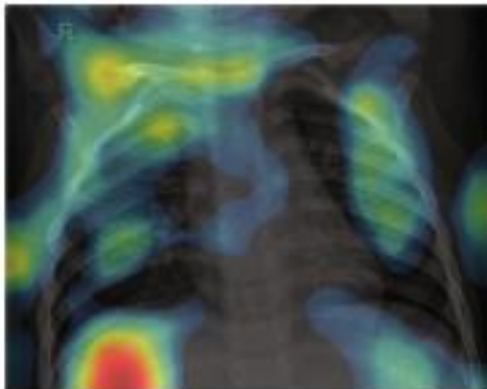
(a)

Original Image - Pneumonia



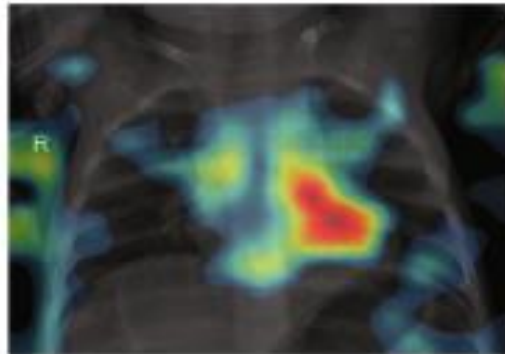
(c)

Normal, 1
Pneumonia, 4.5248e-08
Covid, 9.5175e-11



(b)

Original Image - Normal

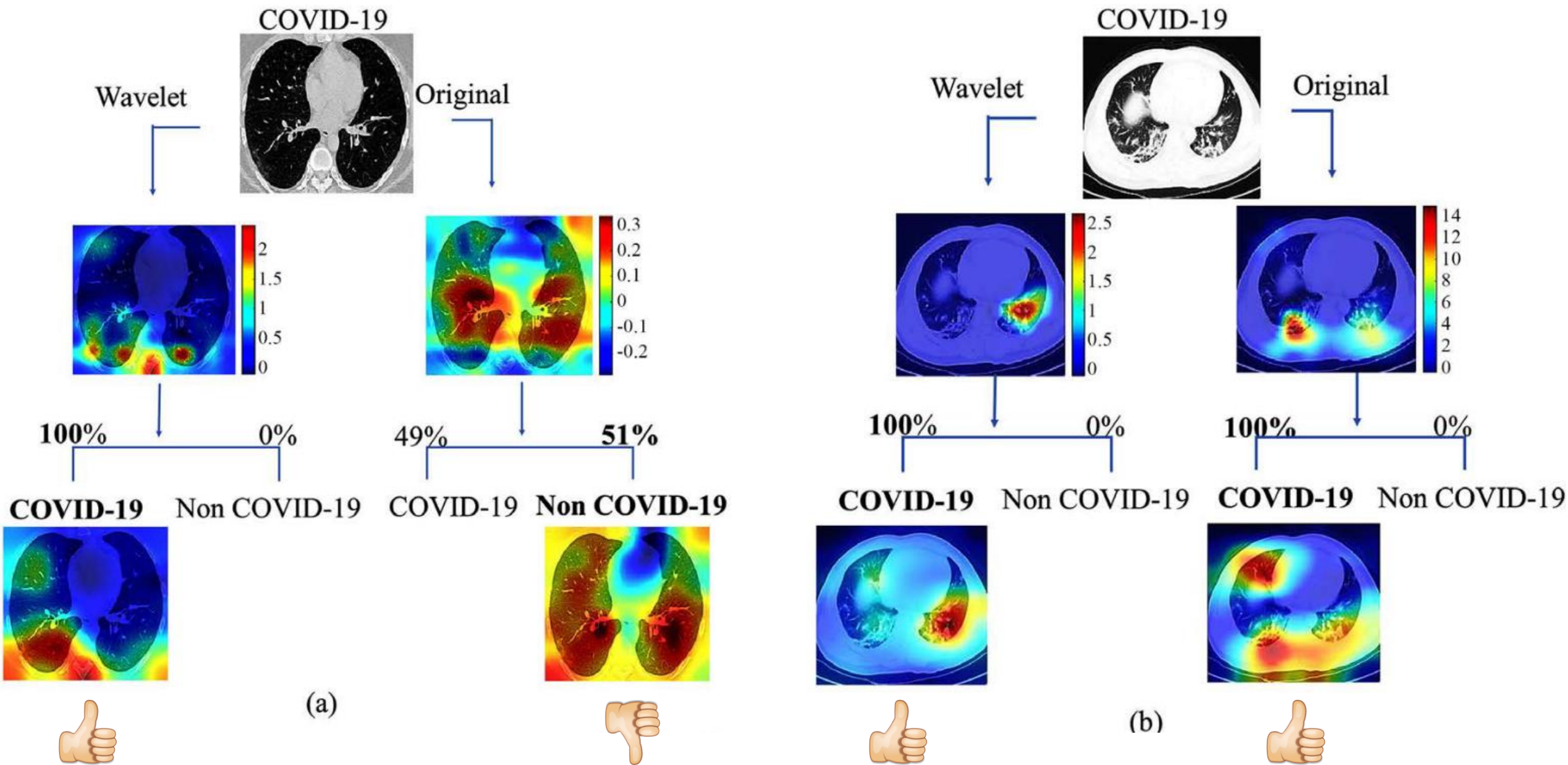


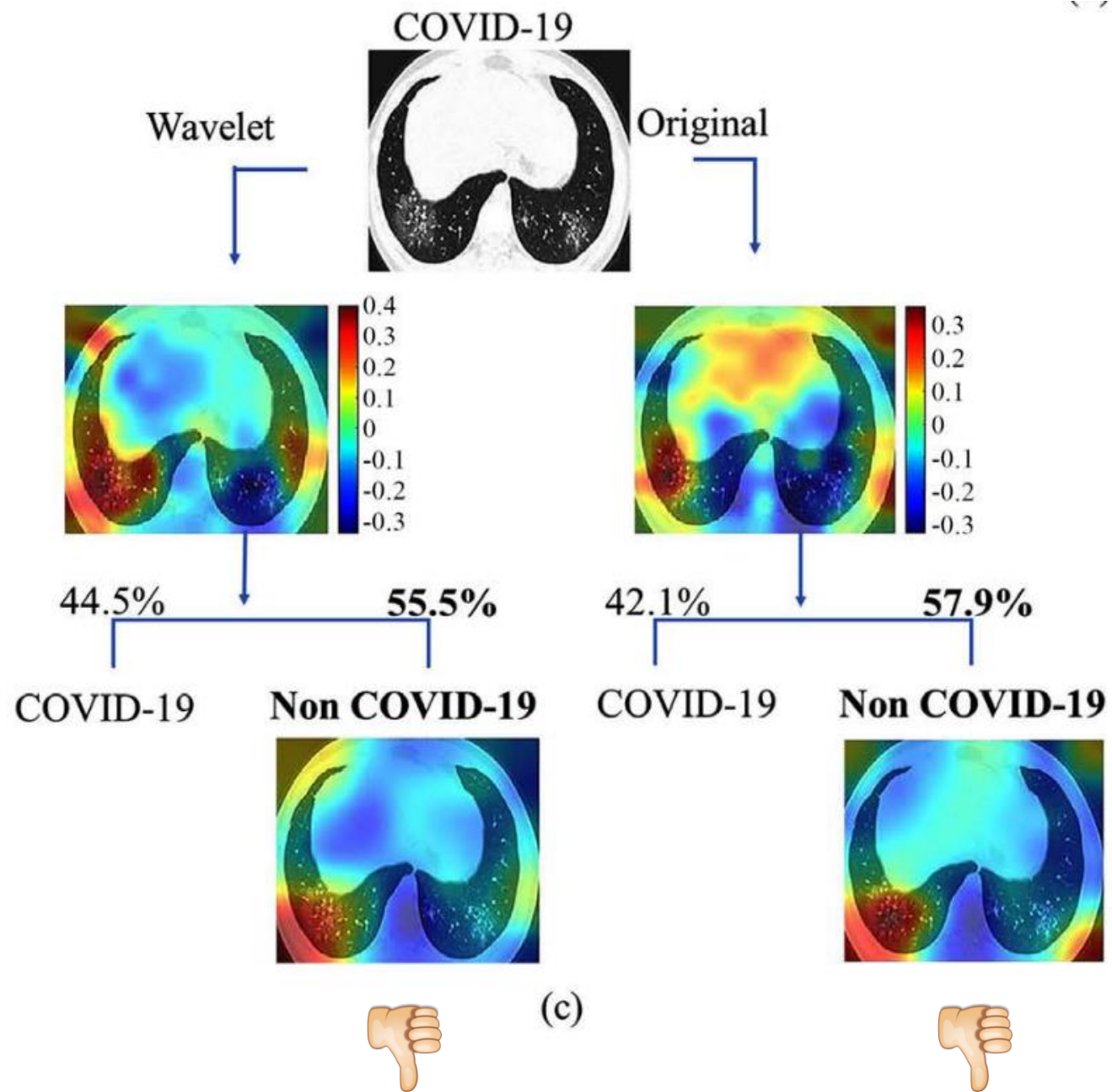
(d)

Pneumonia, 0.99998
Normal, 2.1445e-05
Covid, 1.7446e-08

Pneumonia, 0.74539
Normal, 0.25461
Covid, 3.4892e-06

Matsuyama, E. (2020). A Deep Learning Interpretable Model for Novel Coronavirus Disease (COVID-19) Screening with Chest CT Images. *Journal of Biomedical Science and Engineering*, 13(7), 140–152.





Albert, N. (2020). Evaluation of Contemporary Convolutional Neural Network Architectures for Detecting COVID-19 from Chest Radiographs. *ArXiv.Org*. <https://arxiv.org/abs/2007.01108>

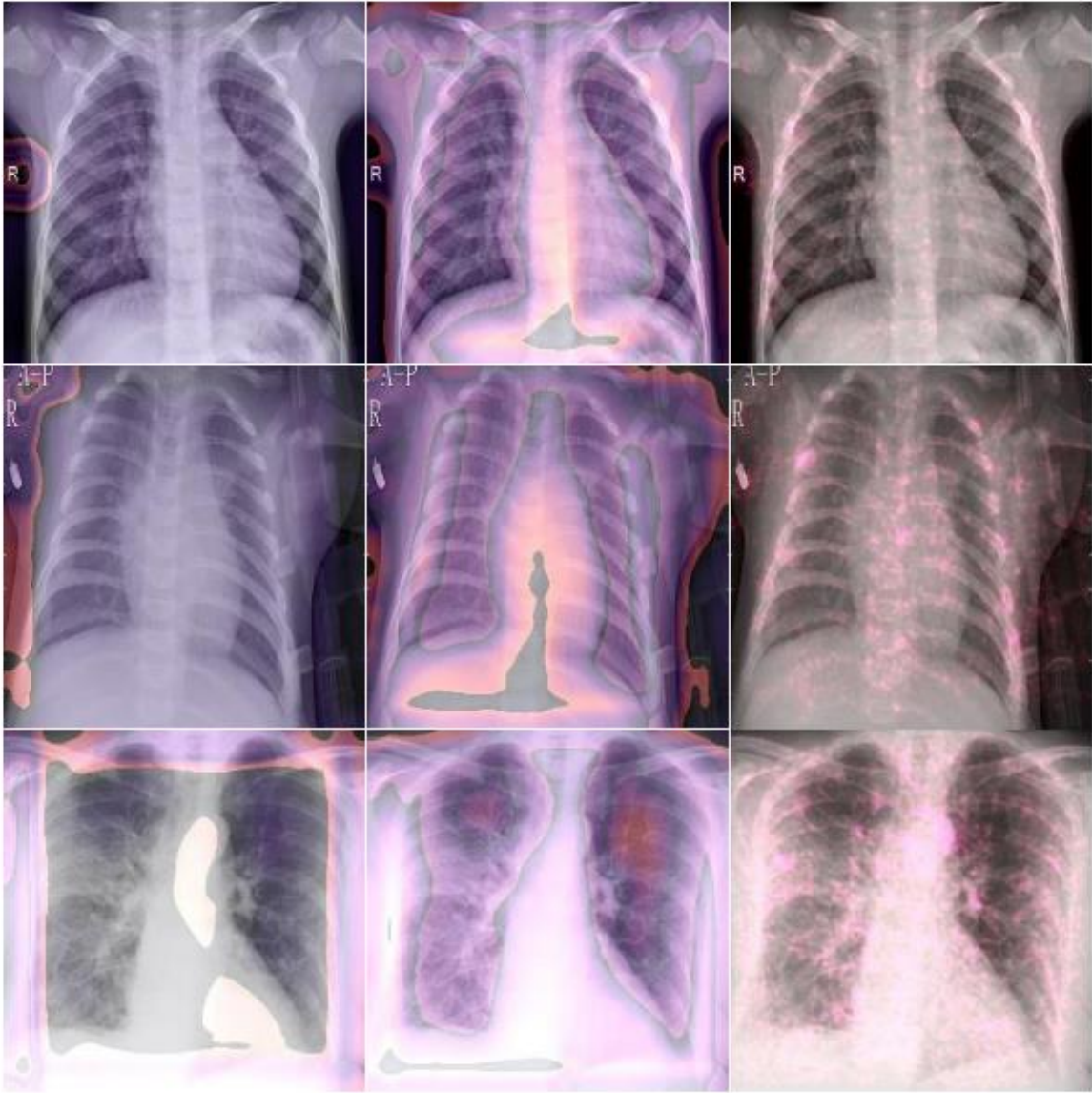


Figure 13. Rows, Top to Bottom: Normal, Bacteria, COVID.
Columns, Left to Right: ImageNet DenseNet, ImageNet ResNet, Efficient model

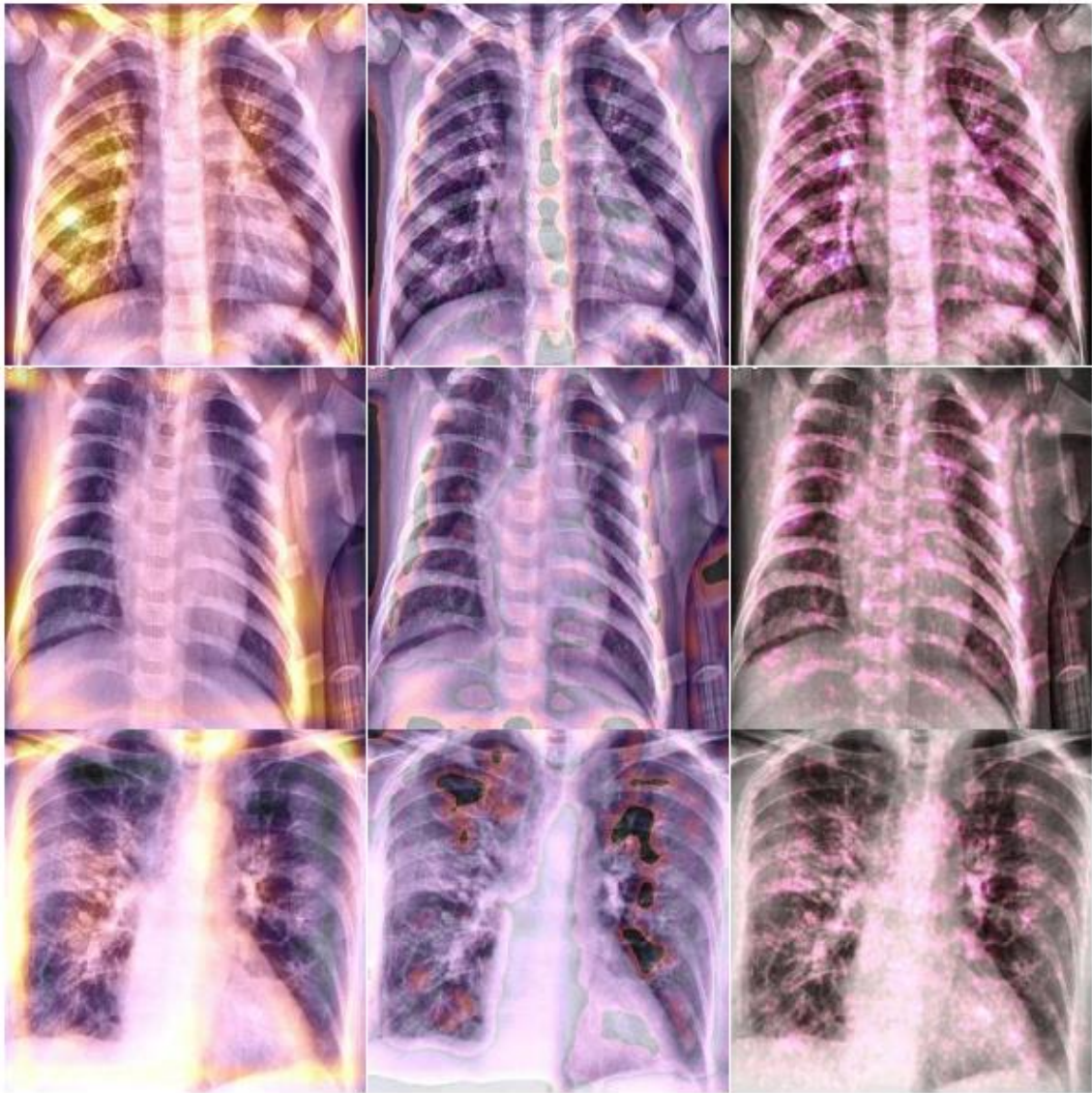
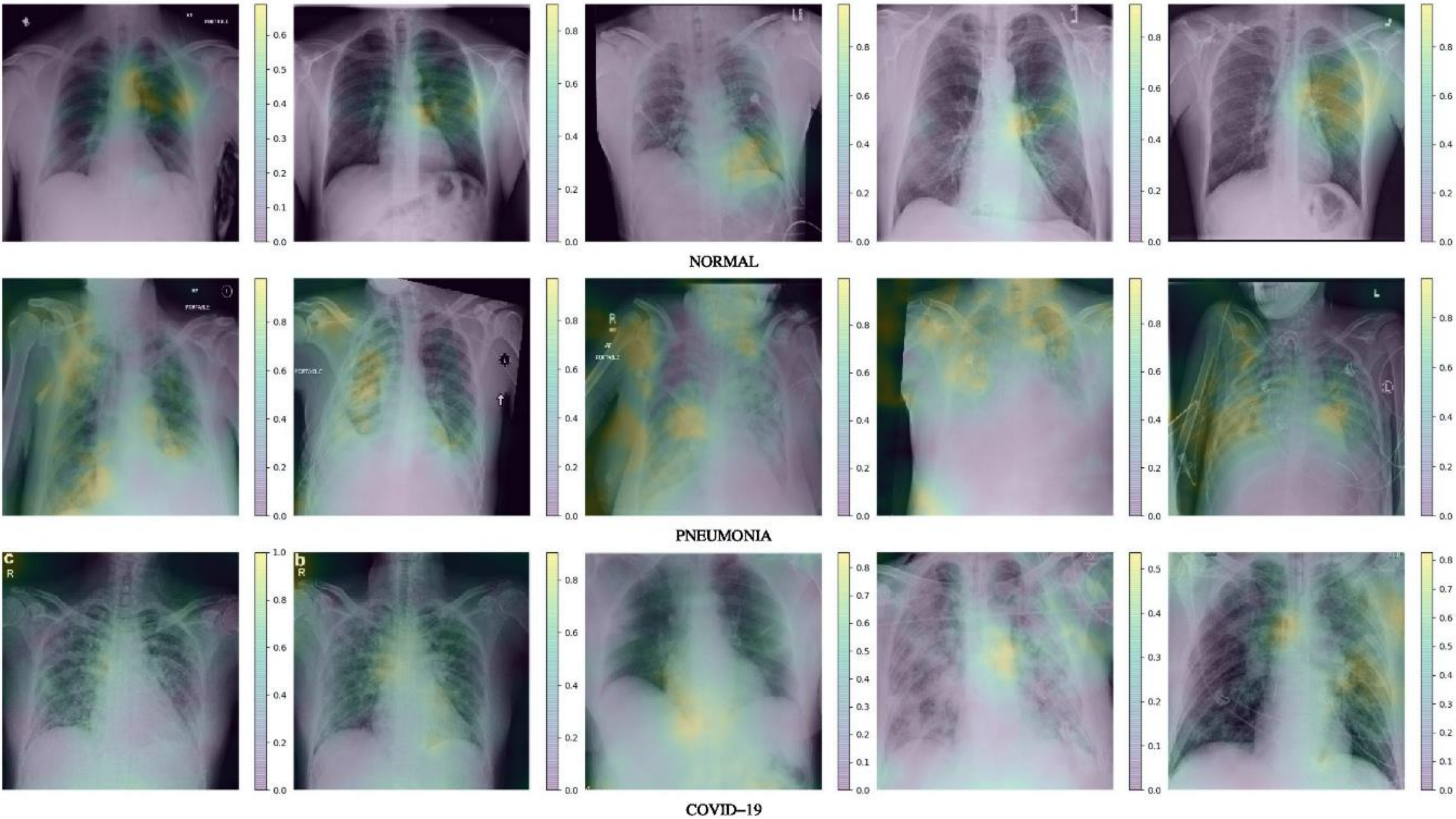


Figure 14. Rows, Top to Bottom: Normal, Bacteria, COVID.
Columns, Left to Right: ImageNet DenseNet, ImageNet ResNet, Efficient model

Khobahi, S., Agarwal, C., & Soltanian, M. (2020). CoroNet: A Deep Network Architecture for Semi-Supervised Task-Based Identification of COVID-19 from Chest X-ray Images. *MedRxiv*, 2020.04.14.20065722.
<https://doi.org/10.1101/2020.04.14.20065722>



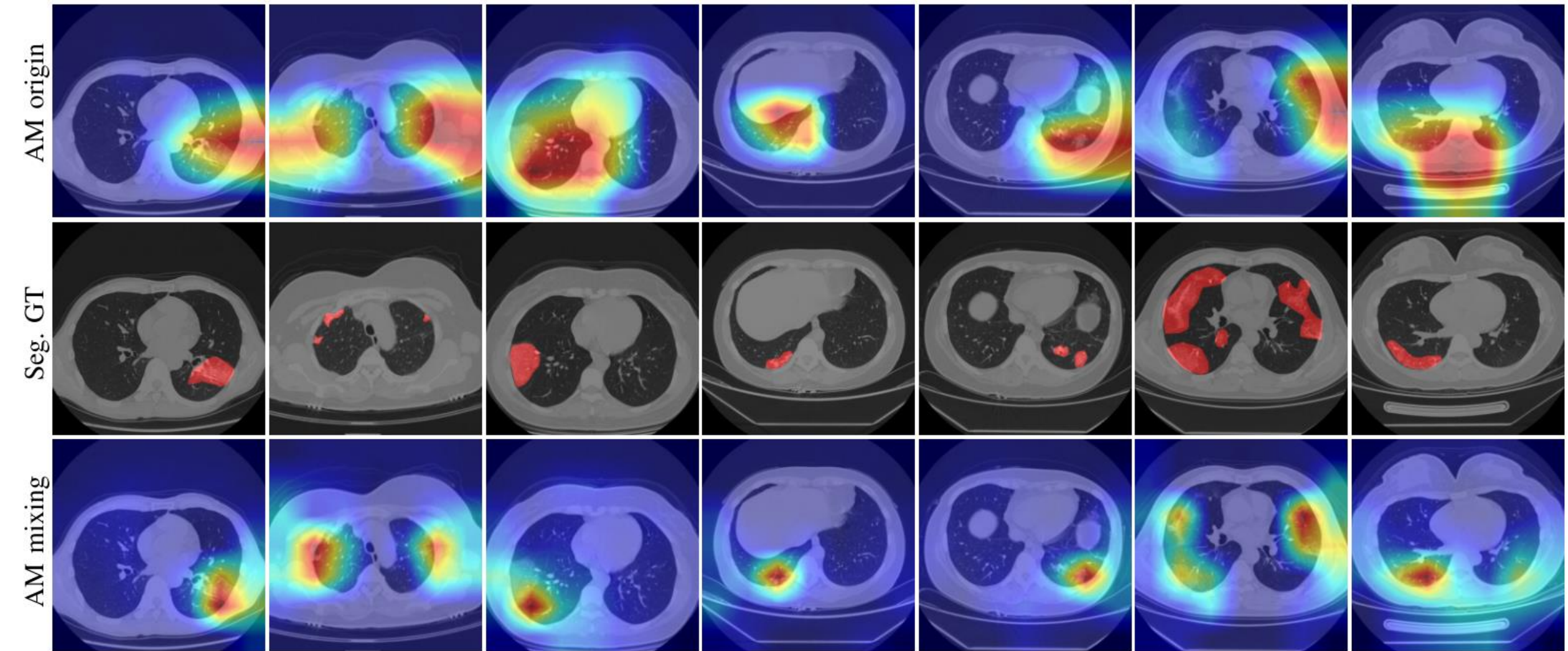


Figure 8. **Visualizations of activation mapping (AM).** AM origin (mixing) means the AM of models trained without (with) image mixing technique [49].

Ahsan, M. M., Gupta, K. D., Islam, M. M., Sen, S., Rahman, M. L., & Hossain, M. S. (2020). *Study of Different Deep Learning Approach with Explainable AI for Screening Patients with COVID-19 Symptoms: Using CT Scan and Chest X-ray Image Dataset.* <http://arxiv.org/abs/2007.12525>

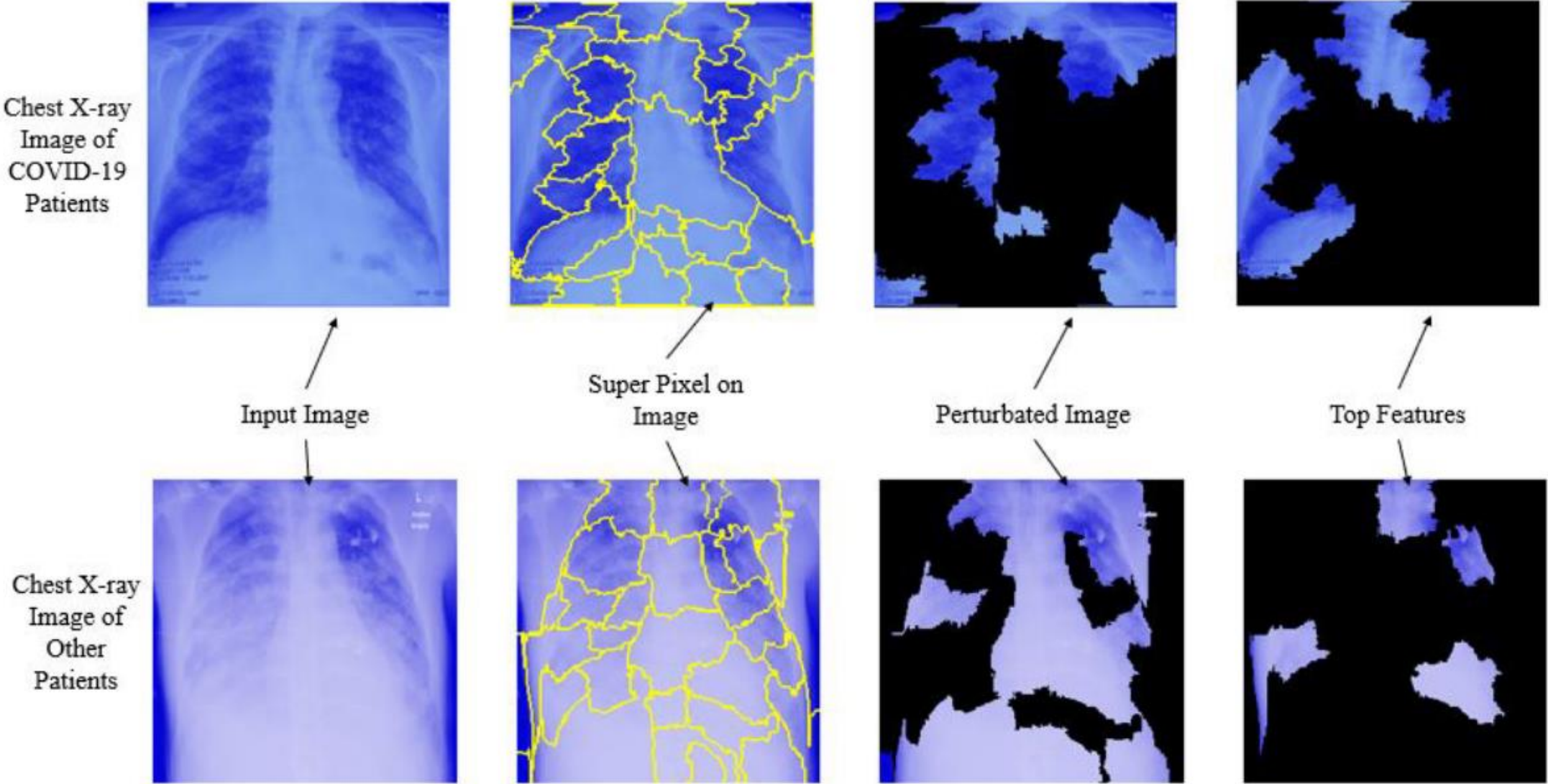


Figure 15: Overall prediction analysis using LIME

Jaiswal, A. K., Tiwari, P., Rath, V. K., Qian, J., Pandey, H. M., & Albuquerque, V. H. C. (2020). COVIDPEN: A Novel COVID-19 Detection Model using Chest X-Rays and CT Scans. *MedRxiv*, 2020.07.08.20149161. <https://doi.org/10.1101/2020.07.08.20149161>

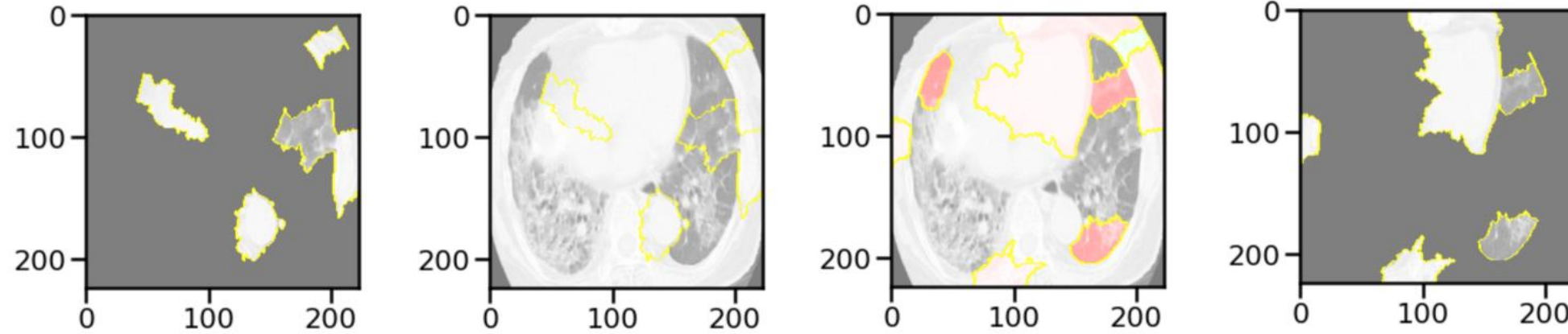


Fig. 6: Model Interpretability - COVID-19 CT Scan Dataset

The regions shaded in pink and green detect superpixels that contributed to and against prediction of chest radiographs and CT scans.

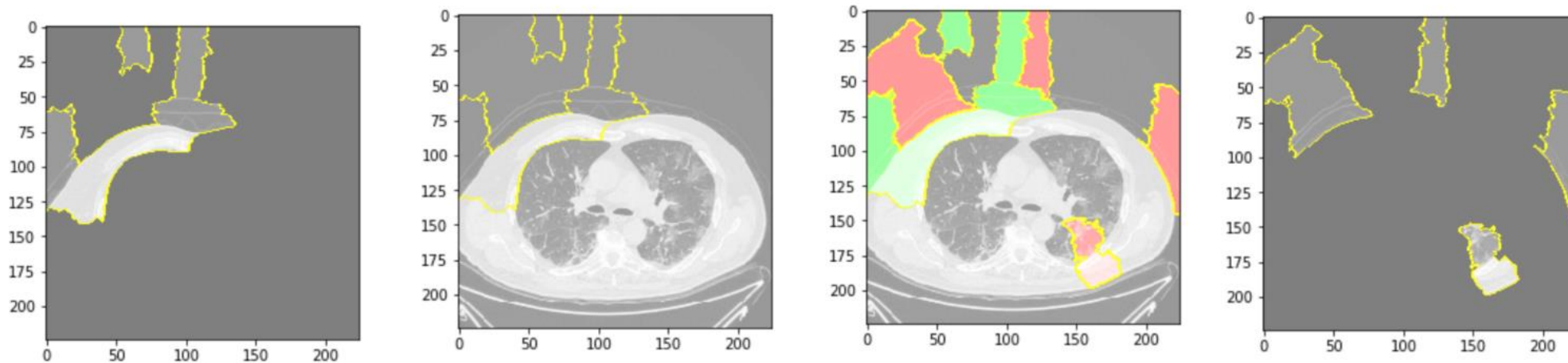


Fig. 7: Model Interpretability - COVID-19 Chest Radiographs Dataset

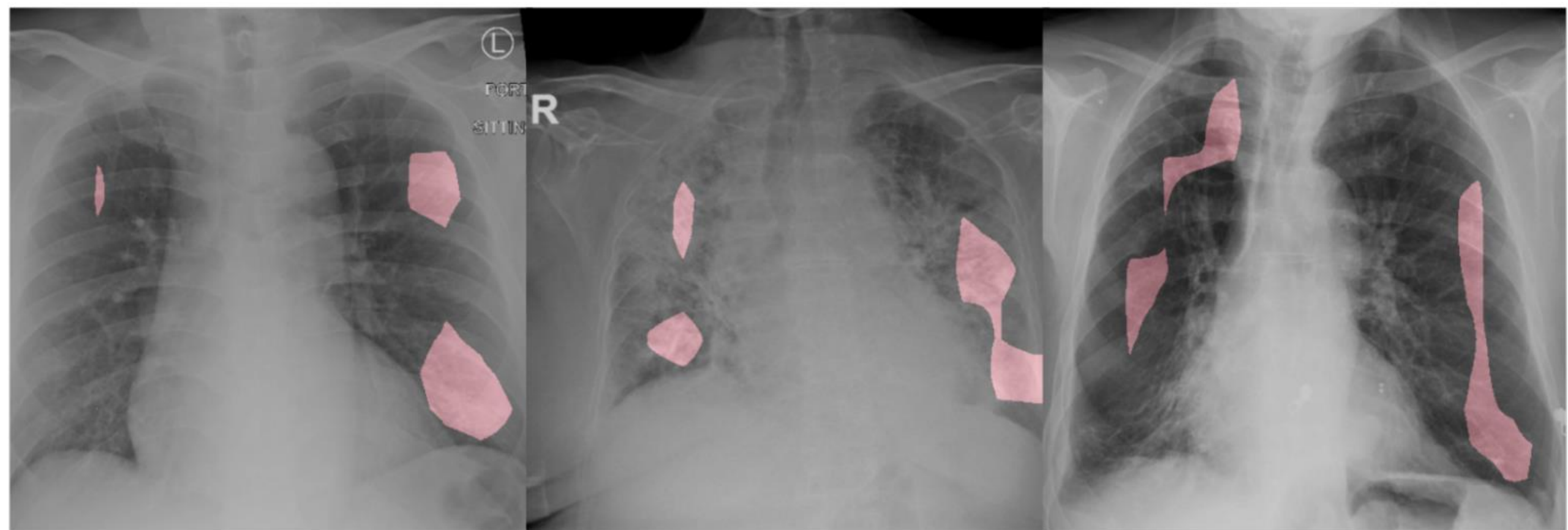


Figure 7. Example CXR images of COVID-19 cases from several different patients and their associated critical factors (highlighted in red) as identified by GSInquire³⁷.

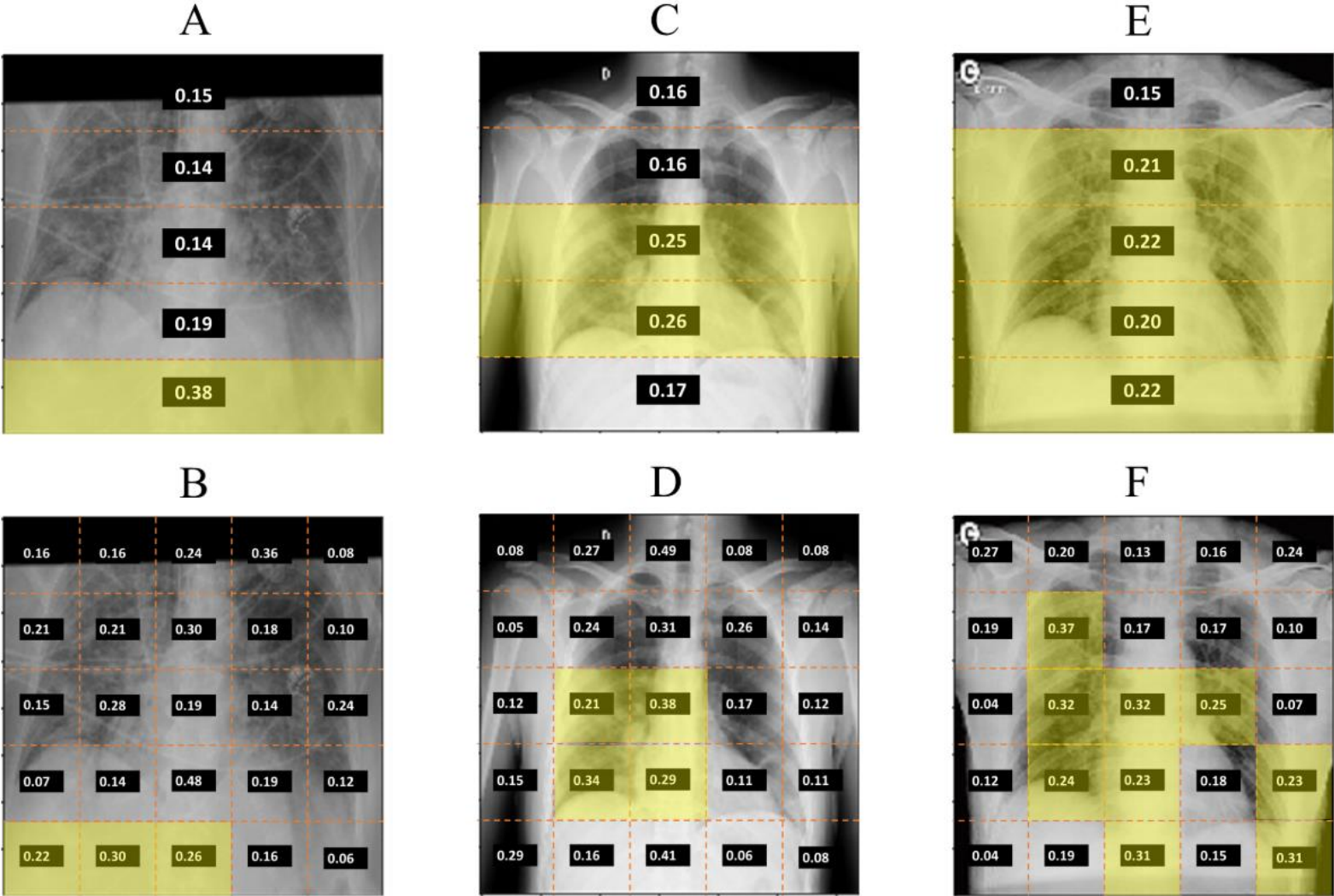


Figure 6: Attention scores for different zones of the lung (horizontal level) and different blocks of the image for 3 patients with COVID-19. Signs of COVID-19 were detected in the lower zone for Patient 1 (A-B), middle zone for Patient 2 (C-D), and lower and middle zones for Patient 3 (E-F).