

Explainability Guided Multi-Site COVID-19 CT Classification

Ameen Ali, Tal Shaharabany, and Lior Wolf

School of Computer Science, Tel Aviv University, Tel Aviv, Israel

Abstract. Radiologist examination of chest CT is an effective way for screening COVID-19 cases. In this work, we overcome three challenges in the automation of this process: (i) the limited number of supervised positive cases, (ii) the lack of region-based supervision, and (iii) the variability across acquisition sites. These challenges are met by incorporating a recent augmentation solution called SnapMix, by a new patch embedding technique, and by performing a test-time stability analysis. The three techniques are complementary and are all based on utilizing the heatmaps produced by the Class Activation Mapping (CAM) explainability method. Compared to the current state of the art, we obtain an increase of five percent in the F1 score on a site with a relatively high number of cases, and a gap twice as large for a site with much fewer training images.

Keywords: Mixing augmentation · Interpretability scores · COVID-19.

1 Introduction

Deep neural networks are currently the leading image classification method. Their ability to generalize is well-documented. However, in many medical imaging domains, one faces challenges that reduce the effectiveness of the generic solutions. First, due to the cost of acquisition, privacy issues, and the expertise required for labeling, the typical datasets are smaller than those available for many other computer vision tasks. Second, in medical images, the exact capturing apparatus, its setting and its operators all can greatly affect the distribution of the obtained images, causing a sizable domain shift. Third, many diseases are manifested through symptoms that are well localized in images, while the supervision is given at the image-level.

In this work, we demonstrate that explainability methods, which link the classification outcome to specific image regions, can provide an important building block for overcoming the three issues. First, the heatmap obtained from such methods serves as the basis of an augmentation method called SnapMIX [17], which we demonstrate is also effective for the COVID-19 classification task we study in this work. Second, the heatmap can provide a delineation of whether or not local image patches are strongly linked to the obtained classification. By requiring that image patches of similar relevancy have similar embedding, we can improve the classification performance. Third, we can use the heatmap in order

to validate, at test time, the stability of the obtained classification by perturbing the image locations that are the most relevant to the prediction. If the majority of the perturbations do not support the prediction, we flip the predicted label.

We evaluate our method with well-established benchmarks for the classification of Computed Tomography (CT) scans as COVID-19 positive or COVID-19 negative, and present clear evidence for the utility of our method. The gap in performance we obtain is larger than the variance between the state of the art methods. On site A, in which performance (F1 score and accuracy) is over 90%, we improve to over 95%. On site B, in which the performance levels are in the high 70 percentages, we obtain results of almost 90%.

2 Related Work

COVID19 Classification The SARS COV-2 infection (COVID-19) has a devastating impact on the respiratory system and has caused an enormous number of deaths. In the last year, many deep learning methods were developed for classifying COVID-19 in 2D or 3D medical images [8,28,41,36]. Some recent methods use transfer learning from models pretrained on ImageNet [13,1].

Following Wang et al. [37], we study classification in two CT datasets. To overcome the domain shift, their approach adds a contrastive loss that decreases the differences between the latent space distributions. Unlike previous work in the domain of CT diagnosis of COVID-19, our method employs a generic ResNet architecture and our contribution lies solely in the training procedure and in the inference procedure.

Data augmentation Many augmentation approaches were developed over the years as a form of regularization. These include geometric transformations [35] and color space transformations [38], which have been shown to improve many medical applications [20].

Data mixing approaches create virtual samples that combine multiple images from different categories. The generated image has a fuzzy label from the two categories. In MixUp [10], the augmented image is a linear interpolation between two different images. The fuzzy labels are computed using the same weights as the images. Cutmix [40] extracts a box from one image and pastes it to the second. The fuzzy labels are proportional to the area of the box. SnapMix [17] is similar to Cutmix, except that the area of the patch is replaced by the sum of the CAM activations within the extracted and the masked patches. It was shown to be highly effective on fine-grained classification datasets of natural images. Here, it is applied to the binary classification of medical images.

Explainability The task of generating a heatmap that indicates local relevancy from the perspective of a CNN observing an input image has been tackled from many different directions, including gradient-based methods [32,33,31], attribution methods [2,24,25,11], and image manipulation methods [6,7,22].

The CAM method [43] is based on the gradient of the loss with respect to the input of each layer. CAM and its extension GradCAM [31] have been used by downstream applications, such as weakly-supervised semantic segmentation [19].

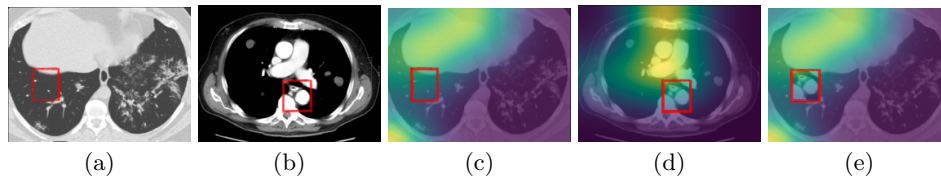


Fig. 1. An illustration of the SnapMix process. (a) A first random image, in this case a positive image from site A, (b) a second random image, shown is a negative image from site B. A random box is marked in each image. (c,d) the CAM maps of (a,b) respectively, with the associated boxes marked. (e) The SnapMix training image obtained by combining images (a) and (b) based on the random boxes. The label of the virtual training image is determined by mixing the labels of the two source images in accordance with the sum of the CAM activations in each box.

Here, we make a novel use of CAM for creating more effective patch embeddings and for test time augmentation.

Contrastive learning The loss we employ between patches of different levels of relevancy is related to contrastive learning methods that have recently made a large impact in the field of self-supervised learning, where it is often used to link an image to its transformed version [14,23,3]. Our work is applied at the patch level. Contrastive learning has emerged in metric learning [4] and subsequently in unsupervised representation learning [12]. The learned embedding brings closer associated samples, while distancing other samples. In our case, association is being determined by the CAM-derived relevancy.

3 Method

Our experiments utilize a Resnet-50 network [15], trained with the conventional binary cross entropy loss L_{BCE} , as the baseline classifier. We then apply (i) SnapMix [17], (ii) a novel optimization term called the Contrastive Patch Embedding loss, and (iii) a novel test time voting procedure. All three techniques utilize the heatmaps produced by the CAM method [43].

3.1 SnapMix [18]

The SnapMix method is illustrated in Fig. 1. It combines two training images, depicted in panels (a) and (b) by considering a random box in each image (marked in red). The importance of each of the two boxes is evaluated by integrating the CAM scores in each box (panels c,d). The virtual sample is generated by pasting the box from the second image onto the selected box of the first image (panel e), and labeling the new image proportionally to the integrated CAM scores. More specifically, a ratio (ρ_a, ρ_b) is computed for each image by considering the sum of all CAM scores in a box over the sum of the CAM scores of the entire image. The labels are then linearly interpolated between the labels of the two images,

using the the complement of the obtained box ratio in the first image $(1 - \rho_a)$ and the ratio in the second image ρ_b .

Unlike the original experiments in [17], which considered datasets with many classes, in our case the problem is binary. It often happens that both images are of the same class. Moreover, since we train using images from two sites, the virtual images created can potentially play a role in overcoming the domain shift.

3.2 Contrastive Patch Embedding

The input images we receive are of size 224×224 , the receptive field of the ResNet-50 architecture is of size 32, and the embedding is of spatial dimensions of 7×7 with a depth of 2,048. For each of the $7 \times 7 = 49$ vectors in \mathbb{R}^{2048} , we compute the sum of the CAM activations in the associated patch of size 32×32 . We then select four vectors out of the 49: two with the highest sum of activations u_1 and u_2 , and two with the lowest sum v_1 and v_2 .

The embedding loss we propose is a contrastive loss [39,14,26] that considers the dot products between the four vectors.

$$L_{\text{CPE}}(u_1, u_2, v_1, v_2) = -\ln \frac{\exp(u_1^\top u_2)}{\exp(u_1^\top u_2) + \sum_{i,j=1}^2 \exp(u_i^\top v_j)} - \ln \frac{\exp(v_1^\top v_2)}{\exp(v_1^\top v_2) + \sum_{i,j=1}^2 \exp(u_i^\top v_j)}. \quad (1)$$

This loss brings together the two most label-supporting embedding-vectors and two most label-opposing embedding-vectors. At the same time, it distances the top label-supporting embedding-vectors from the pair of vectors that support the alternative label.

3.3 CAM-Directed Test Time Augmentation

It may be the case that the decision for a certain label is based on local artifacts that bias the network into giving the wrong prediction. In order to avoid such cases, we classify each image $k + 1$ times: using the entire image, and when masking one out of k different patches.

For this purpose, we divide the image into small non-overlapping patches of size 8×8 , obtaining a grid of size 28×28 . For each cell in the grid, we compute the sum of the CAM activations. We then create $k = 31$ alternative images, by masking out sequentially the k patches with the highest sum of activations. In the first alternative image, we mask out the patch with the highest CAM scores; in the second, we mask out the two patches with the highest CAM scores; and so on. See Fig. 2 for an illustration.

The label that we report is obtained by performing voting among the classifier output of the k images. A supporting vote occurs when the pseudo-probability obtained from the network classifier is at least $\theta = 0.2$ if the original image has a positive label (i.e., a pseudo-probability larger than 0.5), or lower than $1 - \theta$ for

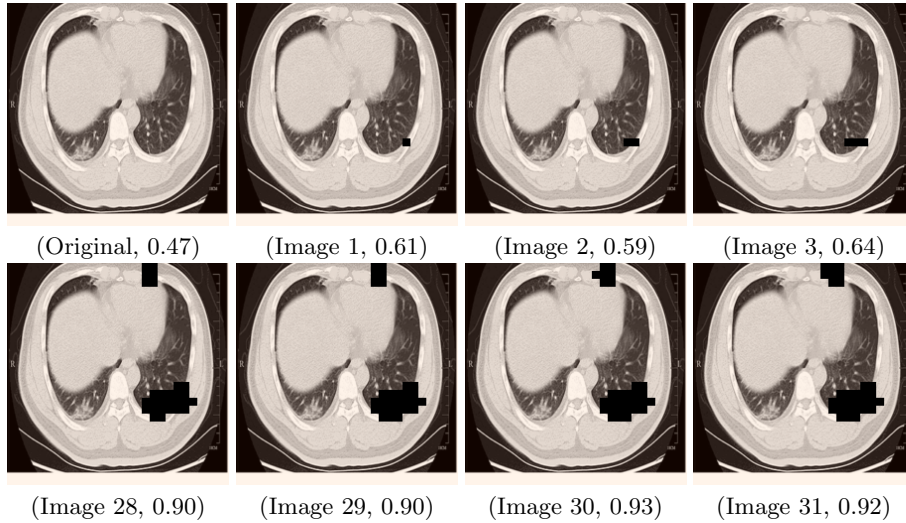


Fig. 2. The original test time image (top left) got a negative classification score, with a probability of 0.47 of being positive. Removing even a small number of patches (Images 1-3) increased this probability to be above 0.5. Subsequently, as more and more patches were removed, the probability of a positive case increased further, and became higher than $1 - \theta$, see the last derived images (out of $k = 31$ images), Images 28-31.

images with negative samples. If more than half the k votes are not supporting, we flip the label. In other words, if the inferred labels of more than half out of the k alternative images contradict, with a high certainty, the label the classifier assigns to the entire image, we flip the predicted label of the image.

4 Experiments

Data We evaluate the proposed method on three datasets, the first two contains CT images of patients who are COVID-19 positive or negative. The SARS-CoV-2 dataset (denoted as site A) consists of 2,482 CT images from 120 patients, in which 1252 are positive with COVID-19. The 1,230 negative samples are inflicted with other lung disease. The resolution of these images vary between 119×104 and 416×512 . The COVID-CT dataset [42] (denoted as site B) is much smaller and includes 349 CT images from 216 COVID-19 positive patients and 397 CT images from 171 control patients. The resolution of the images of site B ranges from 102×137 to 1853×1485 . Following [37], the images of both datasets are resized to a fixed resolution of 224×224 and are intensity normalized to zero mean and unit variance. The classification accuracy, F1 score, Sensitivity, and Precision are reported in percents using the train/test splits of the different datasets.

The third dataset we employ is COVIDx-CT [9] this dataset is considered one of the largest in terms of the number of annotated samples provided, containing

Table 1. COVID-19 classification results for site A

| Method | Accuracy | Precision | Recall | F1 |
|--------------------------------|-------------------|-------------------|-------------------|-------------------|
| Series Adapter [29] | 85.73±0.71 | 90.98±0.79 | 81.91±2.61 | 86.19±1.65 |
| Parallel Adapter [30] | 82.13±1.91 | 83.51±1.87 | 80.02±2.47 | 82.39±1.78 |
| MS-Net [21] | 87.98±1.31 | 93.78±2.76 | 84.91±2.83 | 88.73±1.20 |
| Single (Covidnet) [36] | 77.12±0.98 | 80.04±2.87 | 70.97±2.37 | 76.03±1.13 |
| Single (Redesign) [37] | 89.09±1.08 | 94.58±2.07 | 83.78±0.62 | 88.97±0.91 |
| Joint (Covidnet) [36] | 68.72±1.94 | 68.27±1.21 | 69.41±3.91 | 69.17±1.93 |
| Joint (Redesign) [37] | 78.42±2.19 | 80.82±1.05 | 74.07±3.16 | 77.86±2.01 |
| SepNorm [37] | 88.76±0.78 | 95.46±0.74 | 82.97±1.66 | 87.88±0.81 |
| SepNorm + Contrastive [37] | 90.83±0.93 | 95.75±0.43 | 85.89±1.05 | 90.87±1.29 |
| Baseline architecture | 89.68±0.46 | 95.02±0.40 | 83.99±0.51 | 89.13±0.47 |
| Baseline + CPE loss (ablation) | 91.71±1.21 | 97.02±1.65 | 85.13±1.34 | 90.03±0.61 |
| SnapMix [17] | 92.38±0.32 | 98.33±1.81 | 86.42±1.50 | 91.92±0.37 |
| SnapMix + Contrastive | 91.99±0.13 | 99.02±0.52 | 84.22±1.22 | 91.03±0.83 |
| SnapMix + CPE (ablation) | 95.73±0.07 | 98.97±0.33 | 92.49±0.47 | 95.59±0.12 |
| Our full method | 95.90±0.24 | 98.64±0.12 | 92.93±0.40 | 95.87±0.25 |

35996 training images of negative samples and 82286 of positive samples, the test split contains 12245 and 6018 samples for positive and negative patients respectively. For the quantitative analysis we report accuracy as well as sensitivity and PPV (positive predictive value) for each infection type at the image level.

Implementation Details The architecture of our model is based on ResNet50, followed by an MLP classifier, the ResNet model is initialized with pretrained ImageNet weights. We train the model for 200 epochs. The cross entropy loss is used unweighted on the original samples or on virtual SnapMix samples, as dictated by a beta distribution with a parameter of $\alpha = 1$, which is the default parameter in [17]. The L_{CPE} loss is applied to all samples and is summed, unweighted with the cross entropy loss.

Baseline methods The first two baseline methods used for sites A and B are methods that address domain shift in medical images. Series Adapter [29] and Parallel Adapter [30] include a domain adapter model that is based on a filter bank, in order to learn a joint representation from multiple datasets. MS-Net [21] was originally developed for a multi-site prostate segmentation task. It uses domain-specific auxiliary decoders. For classification tasks, each site is associated with an auxiliary classification head. The results for these three methods are obtained from [37].

The single and joint methods from [36], employ an architecture called Covidnet. The difference is whether the method is trained on each dataset separately or not. It was also rerun in [37], using a modified architecture (redesign). The SepNorm method of [37] uses features that are normalization for each site separately. It is further augmented with a contrastive loss that minimizes the domain shift (“SepNorm + Contrastive”).

We present results for the ResNet-50 based architecture that our method utilizes (“Baseline architecture”), and also study the effect of our CPE loss

Table 2. COVID-19 classification results for site B

| Method | Accuracy | Precision | Recall | F1 |
|--------------------------------|-------------------|-------------------|-------------------|-------------------|
| Series Adapter [29] | 70.01±3.82 | 63.04±4.87 | 74.91±1.89 | 67.08±3.09 |
| Parallel Adapter [30] | 74.93±1.83 | 79.84±1.75 | 71.81±2.47 | 73.46±1.68 |
| MS-Net [21] | 76.23±1.81 | 79.29±1.48 | 74.07±1.29 | 76.54±1.73 |
| Single (Covidnet) [36] | 63.12±2.09 | 64.03±3.91 | 57.73±2.94 | 61.09±1.28 |
| Single (Redesign) [37] | 77.07±1.92 | 79.48±0.96 | 74.69±3.91 | 77.04±2.17 |
| Joint (Covidnet) [36] | 63.27±2.82 | 64.27±3.81 | 54.19±4.17 | 59.78±3.12 |
| Joint (Redesign) [37] | 69.67±0.92 | 64.98±3.17 | 66.94±5.86 | 66.89±4.91 |
| SepNorm [37] | 76.89±0.65 | 80.74±2.98 | 70.34±3.76 | 75.02±1.14 |
| SepNorm + Contrastive [37] | 78.69±1.54 | 78.02±1.34 | 79.71±1.42 | 78.83±1.43 |
| Baseline architecture | 85.23±0.41 | 86.54±0.84 | 83.58±0.81 | 84.51±0.62 |
| Baseline + CPE loss (ablation) | 85.96±1.22 | 87.03±1.22 | 84.71±1.02 | 85.22±0.79 |
| SnapMix [17] | 87.56±0.41 | 88.76±0.53 | 85.19±1.02 | 86.85±0.48 |
| SnapMix + Contrastive | 87.03±0.35 | 88.33±0.85 | 84.22±0.79 | 85.72±0.65 |
| SnapMix + CPE (ablation) | 87.02±0.49 | 88.32±0.69 | 84.69±1.02 | 86.95±0.76 |
| Our full method | 88.76±0.26 | 87.44±0.42 | 88.48±0.19 | 88.25±0.22 |

Table 3. Classification results on COVIDx-CT dataset, which is considered the biggest dataset

| Method | Acc | Sensitivity | | PPV | |
|----------------------|--------------|--------------|--------------|--------------|--------------|
| | | Non-Covid-19 | Covid-19 | Non-Covid-19 | Covid-19 |
| ResNet-50 [16] | 98.7% | 98.7% | 96.2% | 97.8% | 99.1% |
| NASNet-A-Mobile[44] | 98.6% | 97.9% | 96.8% | 99.6% | 97.1% |
| EfficientNet-B0 [34] | 98.3% | 97.8% | 95.8% | 98.7% | 98.6% |
| COVIDNet-CT [9] | 99.1% | 99.0% | 97.3% | 98.4% | 99.7% |
| Ours | 99.5% | 99.7% | 99.7% | 99.8% | 99.8% |

(Eq. 1) on it (“Baseline+CPE loss”). Results are also presented when augmenting this architecture with the SnapMix method. As additional ablations, we present result for SnapMIX combined with either the contrastive loss of [37] (“SnapMix+Contrastive loss”) or with our CPE loss (“SnapMix + CPE”). Finally, we present our full method, which includes SnapMix augmentation, the CPE loss, as well as the CAM-driven test time augmentation and voting. For the COVIDx-CT dataset we compare our method with the reported baselines in [9], the COVIDNet-CT baseline [9] was pre-trained on ImageNet [5] and later finetuned on COVIDx-CT [9] dataset using stochastic gradient descent with momentum [27]. We also compare our model with existing models for image recognition (ResNet50 , EfficientNet-B0 , NASNet-A-Mobile [16,44,34]) for image recognition finetuned on COVIDx-CT dataset.

The results are reported in Tab. 1 for site A, and Tab. 2 for site B. Evidently, for both sites, the baseline architecture is already competitive with the best method from the literature, which is SepNorm with the Contrastive loss. In site A the baseline is slightly inferior, and for site B it is considerably preferable.

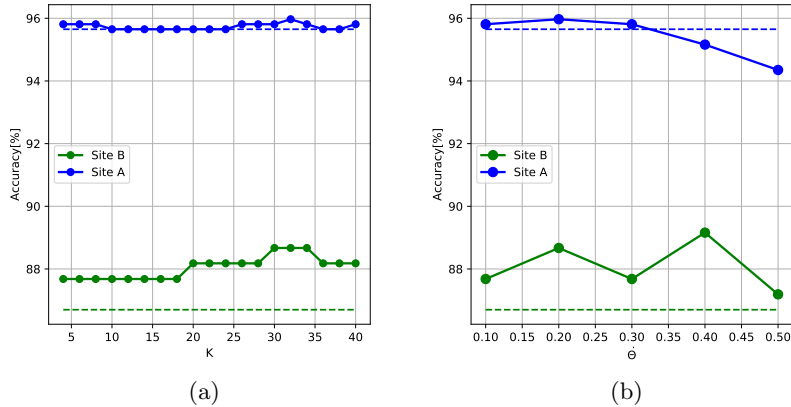


Fig. 3. Parameter sensitivity analysis of the test time augmentation method. (a) The effect of varying the number of alternative images k , (b) The effect of varying the the certainty probability threshold θ .

Adding the CPE loss (Sec. 3.2) improves the results on both sites. So does the SnapMix augmentation, by a larger amount. The two contributions are complementary and adding both the CPE loss and SnapMix provides considerably better results than both separately in site A. In site B, the combination of both provides a slightly higher F1 score than either contributions alone. However, SnapMix by itself is slightly better in terms of the three other scores. The ablation of using the contrastive loss of [37] combined with the SnapMix technique hurts the F1 performance relative to SnapMix in site A, by increasing the precision on the expense of the recall. On site B, it hurts all four scores.

Our complete method, which adds the test time augmentation of Sec. 3.3 on top of SnapMix and the CPE loss, obtains the best accuracy, recall, and F1 score among all methods. Its precision is slightly less than the best ablation method. However, the gap in performance in the F1 score (which combines both precision and recall) is substantial in comparison to the ablation method with the highest precision (almost 5% in site A, and 1.5% in site B). In Tab 3 we show the results for the COVIDx-CT dataset, as shown our full model achieves superior performance over all of the reported baselines.

Parameter Sensitivity SnapMix employs the default augmentation parameters prescribed by [17]. The CPE loss is defined without a temperature parameter that is commonly used in other contrastive learning methods and it employs the minimal number of patches. It is, therefore, virtually parameter-free.

The parameter sensitivity of the CAM-driven test-time voting is explored in Fig. 3, in which performance without this voting (“SnapMix+CPE”) is depicted as a dashed horizontal line. When varying the number of augmented images k (panel a), we observe that for any value of k , there is a performance boost for site B, and this is maximized between $k = 30$ and $k = 35$. The performance

boost for site A is smaller for all k , and peaks at the value of $k = 31$. However, no value of k hurts the performance in site A .

Varying the value of the probability threshold depicted in Fig. 3(b), shows that there is a positive benefit for all tested values $\theta \in \{0.1, 0.2, 0.3, 0.4, 0.5\}$ considering site B. The largest contribution is for the value of 0.4. For site A , however, the contribution is positive only for conservative values (smaller than 0.3, when flipping the label of the test image becomes less frequent). The value of $\theta = 0.2$ provides a small boost to site A and is also the 2nd highest for site B.

5 Conclusions

We present a method of COVID-19 detection in CT scans. The method tackles many of the challenges faced by medical imaging classification systems: distribution shifts across sites, limited training data, and the lack of region based tagging. We propose to combine three different techniques, which have in common the reliance on the heatmap produced by the CAM explainability method. The first method is a powerful regularizer called SnapMix, which was previously used for fine-grained classification. The second is a novel patch embedding method that considers the two patches that show the strongest CAM activations in a given image and the two that present the lowest activations. Finally, we propose a voting method that constructs multiple masked images based on the CAM score. Taken together, our method obtains, despite using a generic network architecture, state of the art results on the two publicly available COVID-19 CT datasets. The gap in performance is extremely sizable, and we demonstrate the individual contribution of each component to it.

Acknowledgment

This project has received funding from the European Research Council (ERC) under the European Unions Horizon 2020 research and innovation programme (grant ERC CoG 725974).

References

1. Apostolopoulos, I.D., Aznaouridis, S.I., Tzani, M.A.: Extracting possibly representative covid-19 biomarkers from x-ray images with deep learning approach and image data related to pulmonary diseases. *Journal of Medical and Biological Engineering* **40**, 462–469 (2020)
2. Bach, S., Binder, A., Montavon, G., Klauschen, F., Müller, K.R., Samek, W.: On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PloS one* **10**(7), e0130140 (2015)
3. Chen, T., Kornblith, S., Norouzi, M., Hinton, G.: A simple framework for contrastive learning of visual representations. In: *International Conference on Machine Learning (ICML)* (2020)

4. Chopra, S., Hadsell, R., LeCun, Y.: Learning a similarity metric discriminatively, with application to face verification. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2005)
5. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: 2009 IEEE conference on computer vision and pattern recognition. pp. 248–255. Ieee (2009)
6. Fong, R., Patrick, M., Vedaldi, A.: Understanding deep networks via extremal perturbations and smooth masks. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 2950–2958 (2019)
7. Fong, R.C., Vedaldi, A.: Interpretable explanations of black boxes by meaningful perturbation. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 3429–3437 (2017)
8. Gozes, O., Frid-Adar, M., Greenspan, H., Browning, P.D., Zhang, H., Ji, W., Bernheim, A., Siegel, E.: Rapid ai development cycle for the coronavirus (covid-19) pandemic: Initial results for automated detection & patient monitoring using deep learning ct image analysis. arXiv preprint arXiv:2003.05037 (2020)
9. Gunraj, H., Wang, L., Wong, A.: Covidnet-ct: A tailored deep convolutional neural network design for detection of covid-19 cases from chest ct images. *Frontiers in medicine* **7** (2020)
10. Guo, H., Mao, Y., Zhang, R.: Mixup as locally linear out-of-manifold regularization. In: AAAI Conference on Artificial Intelligence. pp. 3714–3722 (2019)
11. Gur, S., Ali, A., Wolf, L.: Visualization of supervised and self-supervised neural networks via attribution guided factorization. AAAI (2020)
12. Hadsell, R., Chopra, S., LeCun, Y.: Dimensionality reduction by learning an invariant mapping. In: 2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR’06). vol. 2, pp. 1735–1742. IEEE (2006)
13. Hall, L.O., Paul, R., Goldgof, D.B., Goldgof, G.M.: Finding covid-19 from chest x-rays using deep learning on a small dataset. arXiv preprint arXiv:2004.02060 (2020)
14. He, K., Fan, H., Wu, Y., Xie, S., Girshick, R.: Momentum contrast for unsupervised visual representation learning. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2020)
15. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 770–778 (2016)
16. He, K., Zhang, X., Ren, S., Sun, J.: Identity mappings in deep residual networks. In: European conference on computer vision. pp. 630–645. Springer (2016)
17. Huang, S., Wang, X., Tao, D.: Snapmix: Semantically proportional mixing for augmenting fine-grained data. In: AAAI Conference on Artificial Intelligence (2021)
18. Huang, Z., Wang, X., Wang, J., Liu, W., Wang, J.: Weakly-supervised semantic segmentation network with deep seeded region growing. In: CVPR (2018)
19. Li, K., Wu, Z., Peng, K.C., Ernst, J., Fu, Y.: Tell me where to look: Guided attention inference network. In: CVPR (2018)
20. Litjens, G., Kooi, T., Bejnordi, B.E., Setio, A.A.A., Ciompi, F., Ghafoorian, M., Van Der Laak, J.A., Van Ginneken, B., Sánchez, C.I.: A survey on deep learning in medical image analysis. *Medical image analysis* **42**, 60–88 (2017)
21. Liu, Q., Dou, Q., Yu, L., Heng, P.A.: Ms-net: Multi-site network for improving prostate segmentation with heterogeneous mri data. *IEEE transactions on medical imaging* **39**(9), 2713–2724 (2020)
22. Lundberg, S.M., Lee, S.I.: A unified approach to interpreting model predictions. In: *Advances in Neural Information Processing Systems*. pp. 4765–4774 (2017)

23. Misra, I., van der Maaten, L.: Self-supervised learning of pretext-invariant representations. arXiv preprint arXiv:1912.01991 (2019)
24. Montavon, G., Lapuschkin, S., Binder, A., Samek, W., Müller, K.R.: Explaining nonlinear classification decisions with deep taylor decomposition. *Pattern Recognition* **65**, 211–222 (2017)
25. Nam, W.J., Gur, S., Choi, J., Wolf, L., Lee, S.W.: Relative attributing propagation: Interpreting the comparative contributions of individual units in deep neural networks. arXiv preprint arXiv:1904.00605 (2019)
26. Oord, A.v.d., Li, Y., Vinyals, O.: Representation learning with contrastive predictive coding. arXiv preprint arXiv:1807.03748 (2018)
27. Qian, N.: On the momentum term in gradient descent learning algorithms. *Neural networks* **12**(1), 145–151 (1999)
28. Rahimzadeh, M., Attar, A.: A modified deep convolutional neural network for detecting covid-19 and pneumonia from chest x-ray images based on the concatenation of xception and resnet50v2. *Informatics in Medicine Unlocked* (2020)
29. Rebuffi, S.A., Bilen, H., Vedaldi, A.: Learning multiple visual domains with residual adapters. arXiv preprint arXiv:1705.08045 (2017)
30. Rebuffi, S.A., Bilen, H., Vedaldi, A.: Efficient parametrization of multi-domain deep neural networks. In: *CVPR* (2018)
31. Selvaraju, R.R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., Batra, D.: Grad-cam: Visual explanations from deep networks via gradient-based localization. In: *Proceedings of the IEEE international conference on computer vision* (2017)
32. Shrikumar, A., Greenside, P., Kundaje, A.: Learning important features through propagating activation differences. In: *Proceedings of the 34th International Conference on Machine Learning-Volume 70*. pp. 3145–3153. JMLR. org (2017)
33. Srinivas, S., Fleuret, F.: Full-gradient representation for neural network visualization. In: *Advances in Neural Information Processing Systems* (2019)
34. Tan, M., Le, Q.: Efficientnet: Rethinking model scaling for convolutional neural networks. In: *International Conference on Machine Learning*. pp. 6105–6114. PMLR (2019)
35. Taylor, L., Nitschke, G.: Improving deep learning using generic data augmentation. arXiv preprint arXiv:1708.06020 (2017)
36. Wang, L., Lin, Z.Q., Wong, A.: Covid-net: A tailored deep convolutional neural network design for detection of covid-19 cases from chest x-ray images. *Scientific Reports* **10**(1), 1–12 (2020)
37. Wang, Z., Liu, Q., Dou, Q.: Contrastive cross-site learning with redesigned net for covid-19 ct classification. *IEEE Journal of Biomedical and Health Informatics* **24**(10), 2806–2813 (Oct 2020)
38. Wu, R., Yan, S., Shan, Y., Dang, Q., Sun, G.: Deep image: Scaling up image recognition. arXiv preprint arXiv:1501.02876 **7**(8) (2015)
39. Wu, Z., Xiong, Y., Yu, S.X., Lin, D.: Unsupervised feature learning via non-parametric instance discrimination. In: *CVPR* (2018)
40. Yun, S., Han, D., Oh, S.J., Chun, S., Choe, J., Yoo, Y.: Cutmix: Regularization strategy to train strong classifiers with localizable features. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. pp. 6023–6032 (2019)
41. Zhang, K., Liu, X., Shen, J., Li, Z., Sang, Y., Wu, X., Zha, Y., Liang, W., Wang, C., Wang, K., et al.: Clinically applicable ai system for accurate diagnosis, quantitative measurements, and prognosis of covid-19 pneumonia using computed tomography. *Cell* **181**(6), 1423–1433 (2020)
42. Zhao, J., Zhang, Y., He, X., Xie, P.: Covid-ct-dataset: a ct scan dataset about covid-19. arXiv preprint arXiv:2003.13865 (2020)

43. Zhou, B., Khosla, A., Lapedriza, A., Oliva, A., Torralba, A.: Learning deep features for discriminative localization. In: CVPR (2016)
44. Zoph, B., Vasudevan, V., Shlens, J., Le, Q.V.: Learning transferable architectures for scalable image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 8697–8710 (2018)