
Self-Supervised Neural Architecture Search for Imbalanced Datasets

Aleksandr Timofeev¹ Grigorios G. Chrysos¹ Volkan Cevher¹

Abstract

Neural Architecture Search (NAS) provides state-of-the-art results when trained on well-curated datasets with annotated labels. However, annotating data or even having balanced number of samples can be a luxury for practitioners from different scientific fields, e.g., in the medical domain. To that end, we propose a NAS-based framework that bears the threefold contributions: (a) we focus on the self-supervised scenario, i.e., where no labels are required to determine the architecture, and (b) we assume the datasets are imbalanced, (c) we design each component to be able to run on a resource constrained setup, i.e., on a single GPU (e.g. Google Colab). Our components build on top of recent developments in self-supervised learning (Zbontar et al., 2021), self-supervised NAS (Kaplan & Giryes, 2020) and extend them for the case of imbalanced datasets. We conduct experiments on an (artificially) imbalanced version of CIFAR-10 and we demonstrate our proposed method outperforms standard neural networks, while using $27\times$ less parameters. To validate our assumption on a naturally imbalanced dataset, we also conduct experiments on ChestM-NIST and COVID-19 X-ray. The results demonstrate how the proposed method can be used in imbalanced datasets, while it can be fully run on a single GPU. Code is available [here](#).

1. Introduction

Deep neural networks (DNNs) have demonstrated success in significant tasks, like image recognition (He et al., 2016) and text processing (Devlin et al., 2019). Their stellar performance can be attributed to the following three pillars: a) well-curated datasets, b) tailored network architectures devised by experienced practitioners, c) specialized hardware, i.e. GPUs and TPUs. The adoption of DNNs by practition-

ers in different fields relies on a critical question: *Are those pillars still holding ground in real-world tasks?*

The first obstacle is that well-curated datasets (e.g., uniform number of samples over the classes) might be hard or even impossible to obtain in different fields, e.g., medical imaging. Similarly, when downloading images from the web, the amount of images of dogs/cats is much larger than the images of ‘Vaquita’ fish. One mitigation of such imbalanced classes is the development of the Neural Architecture Search (NAS) (Zoph & Le, 2017), which enabled researchers to build architectures that can then generalize to similar tasks. Obtaining the annotations required for NAS methods is both a laborious and costly process. To that end, self-supervised learning (SSL) has been proposed for extracting representations (Bengio et al., 2013). One major drawback of both NAS and SSL is that they require substantial computational resources, making their adoption harder.

In this work, we propose a novel framework that combines NAS with self-supervised learning and handling imbalanced datasets. Our method relies on recent progress with self-supervised learning (Zbontar et al., 2021) and self-supervised NAS (Kaplan & Giryes, 2020). Specifically, the proposed method designs a network architecture using only unlabelled data, which are also (naturally) imbalanced, e.g., like data automatically downloaded from the web. We pay particular attention to the resources required, i.e., every component of the proposed framework is designed to run on a single GPU, e.g., on Google Colab. We evaluate our method using both a long-tailed distribution and naturally imbalanced datasets for medical imaging. In both settings, we notice that the proposed method results in accuracy that is similar to or better than well-established handcrafted DNNs with a fraction of the parameters of those networks, i.e., up to $27\times$ less parameters.

2. Preliminaries

2.1. Neural Architecture Search

Neural Architecture Search (NAS) can be roughly separated into three components (Elsken et al., 2019): search space, search strategy, and performance estimation strategy. The first component defines the set of architectures to be

¹EPFL, Switzerland. Correspondence to: Aleksandr Timofeev <aleksandr.timofeev@epfl.ch>.

explored; the second one determines how to explore the search space; the third one designs a way to estimate the performance in each step. The first approaches on NAS used evolutionary algorithms (Real et al., 2017; 2019) or reinforcement learning (Zoph & Le, 2017) and outperform handcrafted neural architectures. One major drawback was the immense computational resources required for running NAS. The first papers that focus on the reduction of the computational cost construct a supernet that covers all possible operations and train exponentially many sub-networks simultaneously (Pham et al., 2018; Liu et al., 2019). This approach indeed reduces the computational cost and is improved by the recent works. Specifically, (Cai et al., 2019) samples architecture paths during the search phase such that only one is trained each step. This allows training architecture of the same depth as the final model which eliminates the depth gap in performance. Similarly, (Xu et al., 2020) samples a small subset of channels and replace the rest with skip-connections. Both methods require less memory and reduce search time by using larger batches. However, it is shown that such approaches are unfair in the choice of operations which leads to deteriorating of sub-network performance (Chu et al., 2019; 2020) Another approach (Liu et al., 2020) is based on growing and trimming candidate architectures which is combined with memory-efficient loss. In our work, we stick to the first approach. We avoid the depth gap using the same depth final networks and induce fair competition between operations leveraging (Chu et al., 2020).

The seminal work of *DARTS* (Liu et al., 2019) constructs a differentiable search space using a cell-based approach. The final architecture is constructed by stacking cells. Each cell is a directed acyclic graph (DAG) with N nodes. Each edge represents a candidate operation o_{ij} with input x_i and output x_j , where $x_j = \sum_{i < j} o_{ij}(x_i)$. This neural network is referred to a supernet or a parent network. Such a search space would be discrete, hence a softmax relaxation between the candidate operations $\mathcal{O} = \{o_{ij}^1, o_{ij}^2, \dots, o_{ij}^M\}$ is used:

$$\bar{o}_{ij} = \sum_{o \in \mathcal{O}} \frac{\exp(\alpha_{o_{ij}})}{\sum_{o' \in \mathcal{O}} \exp(\alpha_{o'_{ij}})} o(x),$$

where $\alpha_{o_{ij}}$ is operation mixing weights. Then, NAS is reduced to learning these weights. The final discrete architecture is obtained by $o_{ij} = \arg \max_{o \in \mathcal{O}} \alpha_{o_{ij}}$. The following bi-level optimization problem describes the objective:

$$\min_{\alpha} \ell_{val}(\omega^*(\alpha), \alpha) \text{ s.t. } \omega^*(\alpha) = \arg \min_{\omega} \ell_{train}(\omega, \alpha).$$

where ω denotes normal neural network weights, \mathcal{L}_{val} and \mathcal{L}_{train} are loss functions computed based on batches from validation and train sets correspondingly. It is hard to solve this task directly. Thus, we approximate $\omega^*(\alpha)$ using only a single training step $\ell_{val}(\omega^*(\alpha), \alpha) \approx \ell_{val}(\omega - \xi \nabla_{\omega} \ell_{train}(\omega, \alpha), \alpha)$.

FairDARTS (Chu et al., 2020) DARTS has the significant overhead of maintaining the supernet during training. To mitigate that, during the search phase a shallow network is assumed, which is then duplicated to obtain the full network for evaluation. However, the weights obtained by a shallow neural network are not appropriate for deep models (Chen et al., 2019). Specifically, skip connections are frequently selected as the operation o_{ij} . Additional drawback is lack of weights significantly outperforming others. The recent work of FairDARTS (Chu et al., 2020) mitigates those issues using the following two modifications: (a) it replaces the softmax operation with a sigmoid function to avoid the competition with skip connections as a candidate operation, (b) it encourages sparsity in the architecture weights by using the following zero-one loss: $\ell_{0-1} = -\frac{1}{N} \sum_{i=1}^N (\sigma(\alpha_i) - 0.5)^2$, where α_i , $1 \leq i \leq N$, are architecture weights. The final loss for the architecture weights is $\ell_{total} = \ell_{val}(\omega^*(\alpha), \alpha) + \omega_{0-1} \ell_{0-1}$, where ω_{0-1} controls the strength of the zero-one loss.

2.2. Barlow Twins

Supervised learning has demonstrated success in a number of domains, but requires a massive amount of annotations, while it ignores the enormous amount of unlabelled data that can provide complementary information. The effort to utilize unsupervised learning has been decades old process (Radford et al., 2015; Doersch et al., 2015; Bengio et al., 2013), with the concept of self-supervised learning as a popular method of learning. The idea is to devise one task that the "target label" is known, and use losses developed for supervised learning. For instance, predicting the next word in a sentence enables utilizing the virtually unlimited text on the internet for unsupervised training; this is precisely the method used in the recent successful GPT (Radford et al.) and BERT models (Devlin et al., 2019). Similarly, in visual computing a host of tasks has been used for self-supervised learning (Noroozi & Favaro, 2016; Komodakis & Gidaris, 2018; Chen et al., 2020).

By analogy to other successful self-supervised methods (Chen et al., 2020; Chen & He, 2020; Caron et al., 2020) Barlow Twins (Zbontar et al., 2021) creates a pair of images for every original image. The pair is created by applying two randomly sampled transformations (e.g., random crop, horizontal flip, color distortion in pixels, etc). Similar pairs of images are created for every sample in the mini-batch. The model extracts the representations z^A and z^B of the two corresponding distorted versions of the original mini-batch. The idea is then to make the cross-correlation between z^A and z^B close to the identity. Specifically, the objective function is $\ell_{BT} = \sum_i (1 - \mathcal{C}_{ii})^2 + \lambda \sum_i \sum_{j \neq i} \mathcal{C}_{ij}^2$, where λ is a positive coefficient, \mathcal{C} is a cross-correlation matrix of

the outputs size computed between two outputs:

$$C_{ij} = \frac{\sum_b z_{b,i}^A z_{b,j}^B}{\sqrt{\sum_b (z_{b,i}^A)^2} \sqrt{\sum_b (z_{b,j}^B)^2}},$$

where b indexes batch samples and i, j index the vector dimension of the outputs. In other words, the model is encouraged to differentiate the two distinct images in each pair. The advantages of Barlow Twins is that this decorrelation removes redundant information about samples in the output units. Unlike other recent self-supervised methods, Barlow Twins does not require large batches which is important when there are constrained resources (e.g., a single GPU).

3. Methodology

We propose a new NAS-based approach for real-world datasets, which might not have available labels or they might be imbalanced. The approach consists of two steps: architecture search and subsequent fine-tuning.

Our method is build on top of FairDARTS. This allows eliminating the shortcomings of DARTS as mentioned in Section 2.1. We also replace the supervised loss with a self-supervised one. As (Yang & Xu, 2020) shows, it is also beneficial for learning imbalanced datasets. Namely, the recent Barlow Twins loss which does not require labels. Additionally, we use the supernet with only 3 cells for all steps. This is beneficial for three reasons. Firstly, the training process is efficient and affordable even for slow GPUs, while it produces small but powerful architectures. Secondly, the designed architecture is appropriate for the final model as its depth is unchanged. Thirdly, the learned weights in the first step are fully utilized in the second step (i.e., unlike other NAS methods that the weight values are typically discarded). To fine-tune the designed model, we add on the top another layer which projects the output matrix into the output classes and train it with the focal loss (see Appendix A) in a supervised manner. We do not freeze weights of the rest of the network. Furthermore, we apply the logit adjustment (see Appendix A) to improve learning of rare classes. The ablation study on imbalance handling techniques is in Appendix B.

Our work shares some similarities with the recent Self-Supervised Neural Architecture Search (SSNAS) (Kaplan & Giryes, 2020), which we describe next. SSNAS consists of three steps. In the first step, DARTS is used to determine a cell architecture by a shallow neural network. SSNAS assumes the unlabelled data and uses SimCLR (Chen et al., 2020) for determining the architecture. In the next step, SSNAS stacks the constructed cells from the previous step to obtain the architecture, which is sequentially trained with the same self-supervised loss. In the last step, the architecture is fine-tuned using an annotated dataset. Despite the

similarities, our method differs from SSNAS in four critical ways: (a) Using DARTS has several drawbacks as aforementioned, while DARTS is not robust to initialization and it requires several runs, (b) SimCLR should be used with large batches, while Barlow Twins exhibits better performance and can be executed with a smaller batch size, (c) we use a smaller supernet which improves training time and produces more efficient architectures, (d) we skip the self-supervised pretraining step without occurring a loss in performance (this step required a $\approx 30\%$ overhead in the training time). Lastly, our method is specifically developed for imbalanced datasets by leveraging the logit adjustment and the focal loss.

4. Experiments

Below, we conduct an empirical validation of the proposed method in both a long-tailed distribution and a naturally imbalanced medical dataset. Furthermore, we validate whether the learned architecture can be used for transfer learning in the crucial domain of medical imaging, where we utilize a recent COVID-19 X-ray dataset. Our empirical evidence confirms that the proposed method can achieve similar performance with well-established networks using a fraction of their parameters.

Training details. We apply a first-order version of FairDARTS to accelerate search. A supernet has 3 cells with 4 nodes each. The search space is the same as in (Liu et al., 2019; Kaplan & Giryes, 2020). We use SGD with learning rate 0.025, momentum 0.9, and weight decay 3×10^{-4} with cosine annealing learning rate scheduler (Loshchilov & Hutter, 2017). A batch size of 32 is used, while we train the architecture for 100 epochs. The experiments are performed on NVIDIA Tesla K40c. In fine-tuning step, the focal loss is applied with $\gamma = 2$ and $\alpha_t = 1$. We use light data augmentation: random image cropping and horizontal flipping. The training is run for 600 epochs or until convergence. The rest parameters are the same.

4.1. Evaluation on a long-tailed distribution

We evaluate the proposed method on the long-tailed version of CIFAR-10 dataset (Krizhevsky et al., 2009). To this end, we reduce the number of samples for each class according to an exponential function $n_c^{lt} = \beta^{c+1} n_c$, $0 \leq c < |C|$, where $\beta \in (0, 1)$, n_c and n_c^{lt} are numbers of samples in each class before and after transformation correspondingly, C is a set of classes. The test set stays unchanged. We define the imbalance factor ρ of a dataset as the number of training samples in the largest class divided by the smallest one. The imbalanced dataset is used for both architecture search and subsequent fine-tuning. No label is used for the architecture search.

Table 1: Image classification on CIFAR-10 LT ($\rho = 10$). The number of parameters is reported in millions ($\times 10^6$).

Method	# Params	Error (\downarrow)
ResNet-32 + Focal	21.80	13.34
ResNet-32 + SGM	21.80	12.97
ResNet-32 + BSGM	21.80	12.51
LDAM-DRW	21.80	11.84
smDragon	21.80+	12.17
VE2 + smDragon	21.80+	11.84
SSNAS (Kaplan & Giryes, 2020)	0.83	18.84
Our method	0.81	10.91

In Table 1, we summarize the results of our experiments and compare them against the previous representative works on long-tailed distributions: ResNet-32 + Focal loss (Lin et al., 2017), ResNet-32 + Sigmoid (SGM) and Balanced Sigmoid (BSGM) Cross Entropy losses (Cui et al., 2019), LDAM-DRW (Cao et al., 2019), smDragon and VE2 + smDragon (Samuel et al., 2021), SSNAS (Kaplan & Giryes, 2020). To provide a fair comparison to SSNAS method, we implement it in a common framework with common hyper-parameters. Notably, NAS-based methods require orders of magnitude less parameters. Though, SSNAS result in an increased error for the reduced parameters. Our method significantly improves the results of other methods achieving it with much fewer parameters.

4.2. Evaluation on naturally imbalanced datasets

We assess the performance of the method on *ChestMNIST* (Wang et al., 2017) which is a naturally imbalanced dataset. The dataset contains 78,468 images of chest X-ray scans. There are 14 non-exclusive pathologies. The results are presented in Table 2. The accuracy of models for comparison (ResNet (He et al., 2016), auto-sklearn (Feurer et al., 2019), AutoKeras (Jin et al., 2019), and Google Auto ML) are reported from (Yang et al., 2021). The proposed method is able to achieve the best seen before result wherein it keeps a tiny number of parameters. The smaller number of parameters for SSNAS is caused by a lot of skip-connections.

As typically done in NAS, we evaluate the optimized architecture on transfer learning using *COVID-19 X-ray* dataset (Ozturk et al., 2020). This dataset consists of chest images and naturally imbalanced with $\rho = 4$. This dataset represents several difficulties that arise in real-world settings: (a) there is an imbalance factor $\rho = 4$, (b) the images are slightly different since they are collected from different sources, (c) there is noise in some images since there are some overlaid labels which are not related to the task in hand (see Appendix C). Unfortunately, because such datasets are recent, we have the only model to compare which is DarkCovidNet (Ozturk et al., 2020). Table 3 shows

Table 2: Image classification on ChestMNIST. The number of parameters is reported in millions ($\times 10^6$). The resolution of input images are indicated in the parenthesis.

Method	# Params	Accuracy (\uparrow)
ResNet-18 (28)	11.68	94.7
ResNet-18 (224)	11.68	94.8
ResNet-50 (28)	25.56	94.7
ResNet-50 (224)	25.56	94.7
auto-sklearn (28)	-	64.7
AutoKeras (28)	-	93.9
Google Auto ML (28)	-	94.7
SSNAS (28)	0.57	94.7
Our method (28)	0.82	94.8

Table 3: Image classification on COVID-19 X-ray. The number of parameters is reported in millions ($\times 10^6$). The resolution of input images are indicated in the parenthesis.

Method	# Params	Accuracy (\uparrow)
DarkCovidNet (224)	1.16	98.08
SSNAS (224)	0.57	98.40
SSNAS (28)	0.57	98.08
Our method (224)	0.82	98.40
Our method (28)	0.82	98.40

that our architecture is successfully transferred to another task achieving slightly better results DarkCovidNet but with smaller resolution and number of parameters.

5. Conclusion

In this paper, we propose a NAS framework which is well-suited for scenarios with real-world tasks, where the data are naturally imbalanced and do not have label annotations. Our framework designs an architecture based on the provided unlabelled data using self-supervised learning. To evaluate our method, we conduct experiments on a long-tailed version of CIFAR-10 as well as ChestMNIST and COVID-19 X-ray which are medical datasets that are naturally imbalanced. For all the experiments, we show that the proposed approach provides more compact architecture while maintaining an accuracy on par with strong performing baselines. We expect our method to provide a reasonable framework for practitioners from different fields that want to capitalize on the success of deep neural networks but do not necessarily have well-curated datasets. In addition, our method is suitable for researchers on a constrained budget (e.g., using only the publicly-available Google Colab).

6. Acknowledgements

This project was sponsored by the Department of the Navy, Office of Naval Research(ONR) under a grant number N62909-17-1-2111, by Hasler Foundation Program: Cyber Human Systems (project number 16066). This project has received funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme (grant agreement number 725594). The project was also supported by 2019 Google Faculty Research Award.

References

- Bengio, Y., Courville, A., and Vincent, P. Representation learning: A review and new perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence (T-PAMI)*, 35(8):1798–1828, 2013.
- Cai, H., Zhu, L., and Han, S. ProxylessNAS: Direct neural architecture search on target task and hardware. In *International Conference on Learning Representations (ICLR)*, 2019. URL <https://openreview.net/forum?id=HylVB3AqYm>.
- Cao, K., Wei, C., Gaidon, A., Arechiga, N., and Ma, T. Learning imbalanced datasets with label-distribution-aware margin loss. In Wallach, H., Larochelle, H., Beygelzimer, A., d’Alché-Buc, F., Fox, E., and Garnett, R. (eds.), *Advances in neural information processing systems (NeurIPS)*, volume 32. Curran Associates, Inc., 2019. URL <https://proceedings.neurips.cc/paper/2019/file/621461af90cadfdaf0e8d4cc25129f91-Paper.pdf>.
- Caron, M., Misra, I., Mairal, J., Goyal, P., Bojanowski, P., and Joulin, A. Unsupervised learning of visual features by contrasting cluster assignments. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M. F., and Lin, H. (eds.), *Advances in neural information processing systems (NeurIPS)*, volume 33, pp. 9912–9924. Curran Associates, Inc., 2020. URL <https://proceedings.neurips.cc/paper/2020/file/70feb62b69f16e0238f741fab228fec2-Paper.pdf>.
- Chen, T., Kornblith, S., Norouzi, M., and Hinton, G. A simple framework for contrastive learning of visual representations. In *International Conference on Machine Learning (ICML)*, pp. 1597–1607. PMLR, 2020.
- Chen, X. and He, K. Exploring simple siamese representation learning. *arXiv preprint arXiv:2011.10566*, 2020.
- Chen, X., Xie, L., Wu, J., and Tian, Q. Progressive differentiable architecture search: Bridging the depth gap between search and evaluation. In *International Conference on Computer Vision (ICCV)*, pp. 1294–1303, 2019. doi: 10.1109/ICCV.2019.00138.
- Chu, X., Zhang, B., Xu, R., and Li, J. Fairnas: Rethinking evaluation fairness of weight sharing neural architecture search. *arXiv preprint arXiv:1907.01845*, 2019.
- Chu, X., Zhou, T., Zhang, B., and Li, J. Fair darts: Eliminating unfair advantages in differentiable architecture search. In Vedaldi, A., Bischof, H., Brox, T., and Frahm, J.-M. (eds.), *European Conference on Computer Vision (ECCV)*, pp. 465–480, Cham, 2020. Springer International Publishing. ISBN 978-3-030-58555-6.
- Cui, Y., Jia, M., Lin, T.-Y., Song, Y., and Belongie, S. Class-balanced loss based on effective number of samples. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 9268–9277, 2019.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL <https://www.aclweb.org/anthology/N19-1423>.
- Doersch, C., Gupta, A., and Efros, A. A. Unsupervised visual representation learning by context prediction. In *International Conference on Computer Vision (ICCV)*, December 2015.
- Dong, Y., Shen, X., Jiang, Z., and Wang, H. Recognition of imbalanced underwater acoustic datasets with exponentially weighted cross-entropy loss. *Applied Acoustics*, 174:107740, 2021.
- Elsken, T., Metzen, J. H., and Hutter, F. Neural architecture search: A survey. *Journal of Machine Learning Research*, 20(55):1–21, 2019. URL <http://jmlr.org/papers/v20/18-598.html>.
- Feurer, M., Klein, A., Eggenberger, K., Springenberg, J. T., Blum, M., and Hutter, F. Auto-sklearn: efficient and robust automated machine learning. In *Automated Machine Learning*, pp. 113–134. Springer, Cham, 2019.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- Jin, H., Song, Q., and Hu, X. Auto-keras: An efficient neural architecture search system. In *International Conference on Knowledge Discovery & Data Mining (SIGKDD)*, pp. 1946–1956, 2019.

- Kaplan, S. and Giryes, R. Self-supervised neural architecture search. *arXiv preprint arXiv:2007.01500*, 2020.
- Komodakis, N. and Gidaris, S. Unsupervised representation learning by predicting image rotations. In *International Conference on Learning Representations (ICLR)*, 2018.
- Krizhevsky, A. et al. Learning multiple layers of features from tiny images. 2009.
- Lin, T.-Y., Goyal, P., Girshick, R., He, K., and Dollár, P. Focal loss for dense object detection. In *International Conference on Computer Vision (ICCV)*, pp. 2980–2988, 2017.
- Liu, H., Simonyan, K., and Yang, Y. DARTS: Differentiable architecture search. In *International Conference on Learning Representations (ICLR)*, 2019. URL <https://openreview.net/forum?id=S1eYHoC5FX>.
- Liu, P., Wu, B., Ma, H., and Seok, M. Memnas: Memory-efficient neural architecture search with grow-trim learning. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- Loshchilov, I. and Hutter, F. Sgdr: Stochastic gradient descent with warm restarts. In *International Conference on Learning Representations (ICLR)*, 08 2017.
- Menon, A. K., Jayasumana, S., Rawat, A. S., Jain, H., Veit, A., and Kumar, S. Long-tail learning via logit adjustment. In *International Conference on Learning Representations (ICLR)*, 2021. URL <https://openreview.net/forum?id=37nvvqkCo5>.
- Noroozi, M. and Favaro, P. Unsupervised learning of visual representations by solving jigsaw puzzles. In *European Conference on Computer Vision (ECCV)*, pp. 69–84. Springer, 2016.
- Ozturk, T., Talo, M., Yildirim, E. A., Baloglu, U. B., Yildirim, O., and Rajendra Acharya, U. Automated detection of covid-19 cases using deep neural networks with x-ray images. *Computers in Biology and Medicine*, 121:103792, 2020. ISSN 0010-4825. doi: <https://doi.org/10.1016/j.compbiomed.2020.103792>. URL <https://www.sciencedirect.com/science/article/pii/S0010482520301621>.
- Pham, H., Guan, M., Zoph, B., Le, Q., and Dean, J. Efficient neural architecture search via parameters sharing. In *International Conference on Machine Learning (ICML)*, pp. 4095–4104. PMLR, 2018.
- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., and Sutskever, I. Language models are unsupervised multitask learners.
- Radford, A., Metz, L., and Chintala, S. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.
- Real, E., Moore, S., Selle, A., Saxena, S., Suematsu, Y. L., Tan, J., Le, Q. V., and Kurakin, A. Large-scale evolution of image classifiers. In *International Conference on Machine Learning (ICML)*, pp. 2902–2911. PMLR, 2017.
- Real, E., Aggarwal, A., Huang, Y., and Le, Q. V. Regularized evolution for image classifier architecture search. In *AAAI Conference on Artificial Intelligence*, volume 33, pp. 4780–4789, 2019.
- Sambasivam, G. and Opiyo, G. D. A predictive machine learning application in agriculture: Cassava disease detection and classification with imbalanced dataset using convolutional neural networks. *Egyptian Informatics Journal*, 22(1):27–34, 2021.
- Samuel, D., Atzmon, Y., and Chechik, G. From generalized zero-shot learning to long-tail with class descriptors. In *Winter Conference on Applications of Computer Vision (WACV)*, pp. 286–295, 2021.
- Wang, X., Peng, Y., Lu, L., Lu, Z., Bagheri, M., and Summers, R. M. Chestx-ray8: Hospital-scale chest x-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 3462–3471, 2017. doi: 10.1109/CVPR.2017.369.
- Xu, Y., Xie, L., Zhang, X., Chen, X., Qi, G.-J., Tian, Q., and Xiong, H. Pc-darts: Partial channel connections for memory-efficient architecture search. In *International Conference on Learning Representations (ICLR)*, 2020. URL <https://openreview.net/forum?id=BJ1S634tPr>.
- Yang, J., Shi, R., and Ni, B. Medmnist classification decathlon: A lightweight automl benchmark for medical image analysis. In *International Symposium on Biomedical Imaging (ISBI)*, pp. 191–195, 2021. doi: 10.1109/ISBI48211.2021.9434062.
- Yang, Y. and Xu, Z. Rethinking the value of labels for improving class-imbalanced learning. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M. F., and Lin, H. (eds.), *Advances in neural information processing systems (NeurIPS)*, volume 33, pp. 19290–19301. Curran Associates, Inc., 2020. URL <https://proceedings.neurips.cc/paper/2020/file/e025b6279c1b88d3ec0eca6fcb6e6280-Paper.pdf>.

Zbontar, J., Jing, L., Misra, I., LeCun, Y., and Deny, S. Barlow twins: Self-supervised learning via redundancy reduction. *arXiv preprint arXiv:2103.03230*, 2021.

Zoph, B. and Le, Q. V. Neural architecture search with reinforcement learning. In *International Conference on Learning Representations (ICLR)*, 2017.

A. Handling imbalanced datasets

Logit adjustment. In many real-world applications, gathering balanced datasets is difficult or even impossible. For instance, in medical analysis, a small group of the patients has a specific pathology that the majority of the population does not have. In such settings, doing predictions on imbalanced datasets is crucial.

In (Menon et al., 2021), the authors propose the logit adjustment for the loss function to handle imbalance which corrects the output of the model before softmax operation. Specifically, they introduce the logit adjusted softmax cross-entropy loss:

$$\ell(y, f(x)) = -\log \frac{e^{f_y(x) + \tau \log \pi_y}}{\sum_{y' \in [L]} e^{f_{y'}(x) + \tau \log \pi_{y'}}},$$

where L is a number of classes, $f_y(x)$ is a logit of the given class, π_y is empirical frequencies of classes. Therefore, we induce the label-dependent prior offset which requires a larger margin for rare classes.

Focal loss. The focal loss (Lin et al., 2017) is frequently used in imbalanced datasets (Sambasivam & Opiyo, 2021; Dong et al., 2021). The idea behind the focal loss is to give a lower weight to easily classified samples. In a binary case, we introduce $p_t = \mathbb{1}[y = 1]p + \mathbb{1}[y = -1](1 - p)$, where $y \in \{-1, 1\}$ are labels, p is model’s estimated probability, and $\mathbb{1}[\cdot]$ is an indicator function. The cross-entropy loss is then $\text{CE}(p_t) = -\log p_t$. To tackle the imbalance problem, the focal loss adds a modulating factor to the weighted cross-entropy $\text{FL}(p_t) = -\alpha_t(1 - p_t)^\gamma \log p_t$, where α_t are loss weights which can be inverse class frequencies, γ is a tunable parameter. If a sample is misclassified and p_t is small, the loss is unaffected. However, when $p_t \rightarrow 1$ then $(1 - p_t)^\gamma \rightarrow 0$ which gives less weight to this sample.

B. Ablation study

To show effectiveness of components responsible for handling imbalance, we analyse the performance of all combinations of the Focal loss and the logit adjustment as well as their absence. In the latter case, we use simply the cross-entropy loss. The results are summarized in Table B1. The best result is achieved by combination of the focal loss and

Table B1: Ablation study on imbalance handling techniques for image classification on CIFAR-10 LT ($\rho = 10$). CE = Cross-Entropy loss, FL = Focal loss.

Method	Error (\downarrow)
CE	13.21
CE + Logit adj.	13.05
FL	11.78
FL + Logit adj.	10.91

logit adjustment. Removing the latter slightly deteriorates the performance while absence of the focal loss is significant.

C. COVID-19 X-ray dataset

In Figure C1, we show four representative samples of the COVID-19 X-ray dataset. All images are collected from different sources, while some images contain unrelated content overlaid on image. Likewise, the light intensity, the resolutions, and the image formats might differ from image to image which makes learning harder.

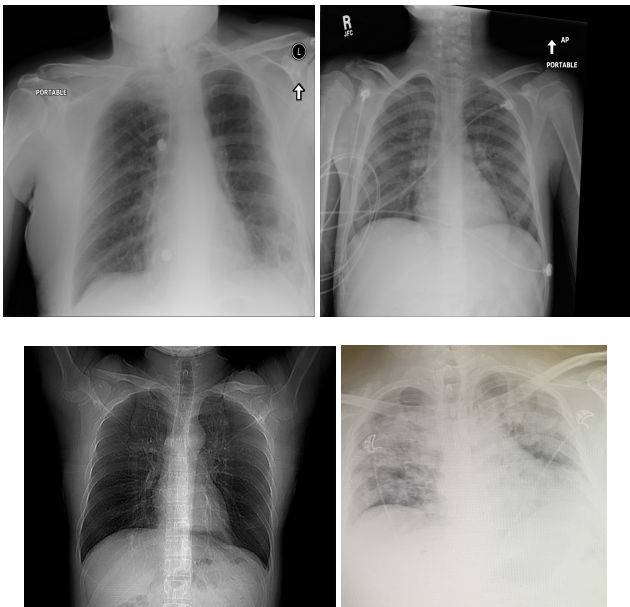


Figure C1: The samples of COVID-19 X-ray dataset. *Top row*: no findings. *Bottom row*: COVID-19.