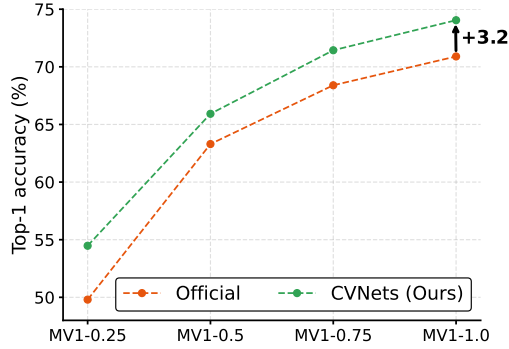# CVNets: High Performance Library for Computer Vision
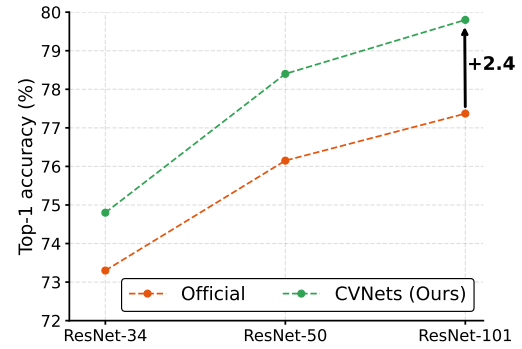
Sachin Mehta*
Apple

Farzad Abdolhosseini
Apple

Mohammad Rastegari
Apple

(a) MobileNetv1 (MV1) at different width factors



(b) ResNet at different depths

**Figure 1: CVNets can be used to improve the performance of different deep neural networks on the ImageNet dataset significantly with simple training recipes (e.g., random resized cropping and horizontal flipping). The official MobileNetv1 [10] and ResNet [7] results are from TensorflowLite [5] and Torchvision [15], respectively.**

## ABSTRACT

We introduce CVNets, a high-performance open-source library for training deep neural networks for visual recognition tasks, including classification, detection, and segmentation. CVNets supports image and video understanding tools, including data loading, data transformations, novel data sampling methods, and implementations of several standard networks with similar or better performance than previous studies. Our source code is available at: https://github.com/apple/ml-cvnets.

## KEYWORDS

Computer vision, deep learning, image and video understanding

## 1 INTRODUCTION

With the rise of deep learning, significant progress has been made in visual understanding tasks, including novel light- and heavy-weight architectures, dedicated hardware and software stacks, advanced data augmentation methods, and better training recipes. There exist several popular libraries that provide implementations for different tasks and input modalities, including Torchvision [15], TensorflowLite [5], timm [21], and PyTorchVideo [4]. Many of these libraries are modular and are designed around a particular task and input modality, and provide implementations and pre-trained weights of different networks with varying performance. However, *reproducibility varies across these libraries*. For example, Torchvision library uses advanced training recipes (e.g., better augmentation) to achieve the same performance for training MobileNetv3 on the ImageNet dataset [3] as TensorflowLite with simple training recipes.

We introduce CVNets, a PyTorch-based deep learning library for training computer vision models with higher performance. With CVNets, we enable researchers and practitioners in academia and industry to train either novel or existing deep learning architectures with high-performance across different tasks and input modalities. CVNets is a modular and flexible framework that aims to train deep neural networks faster with simple or advanced training recipes. Simple recipes are useful for research in resource-constrained environments as they train models for fewer epochs with basic data augmentation (random resized crop and flipping) as compared to advanced training recipes, which trains model for $2-4\times$ longer with advanced augmentation methods (e.g., CutMix and MixUp). With simple recipes (similar to the ones in original publications) and variable batch sampler (Section 3.1), CVNets improves the performance of ResNet-101 significantly (Figure 1) on the ImageNet dataset while for advanced training recipes with the same batch size and number of epochs, it delivers similar performance to previous methods while requiring $1.3\times$ fewer optimization updates.

## 2 CVNETS LIBRARY DESIGN

CVNets follows the design principles below:

***Modularity.*** CVNets provides independent components; allowing users to plug-and-play different components across different visual recognition tasks for both research and production use cases. CVNets implement different components, including datasets and models for different tasks and input modalities, independently. For example, different classification backbones (e.g., ResNet-50) trained in CVNets can be seamlessly integrated with object detection (e.g., SSD) or semantic segmentation (e.g., DeepLabv3) pipelines for studying the generic nature of an architecture.

***Flexibility.*** With CVNets, we would like to enable new use cases in research as well as production. We designed CVNets such that

---

*Project lead and main contributor

new components (e.g., models, datasets, loss functions, data samplers, and optimizers) can be integrated easily. We achieve this by registering each component. As an example, ADE20k dataset for the task of segmentation is registered in `CVNets` as:

`@register_dataset(name="ade20k", task="segmentation")`

To use this dataset for training, one can use `dataset.name` and `dataset.category` as command line arguments.

***Reproducibility.*** `CVNets` provide reproducible implementations of standard models for different computer vision tasks. Each model is benchmarked against the performance reported in original publications as well as the previous best reproduction studies. The pre-trained weights of each model are released online to enable future research.

***Compatibility.*** `CVNets` is compatible with hardware accelerated frameworks (e.g., CoreML) and domain-specific libraries (e.g., PyTorchVideo). The models from domain-specific libraries can be easily consumed in the `CVNets`, as shown in Listing 1; reducing researchers overhead in implementing new components or submodules in `CVNets`.

```
1  from pytorchvideo.models import resnet
2
3  @register_video_cls_models("resnet_3d")
4  class ResNet3d(BaseVideoEncoder):
5      def __init__(self, opts):
6          super().__init__(opts=opts)
7          self.model = resnet.create_resnet(
8              input_channel=3,
9              model_depth=50,
10             model_num_class=400,
11             norm=nn.BatchNorm3d,
12             activation=nn.ReLU,
13         )
```

**Listing 1: An example of registering a video classification model from PyTorchVideo inside CVNets on the Kinetics-400 dataset**

***Beyond ImageNet.*** Many standard models are benchmarked on the ImageNet dataset. With `CVNets`, we would like to enable researchers to build generic computer vision architectures that can be easily scaled to down-stream tasks such as segmentation and detection. Any classification backbone in `CVNets` (either existing or new) can seamlessly be integrated with down-stream networks (e.g., PSPNet and SSD) and enables researchers to study the generic nature of different classification models.

## 3 CVNETS LIBRARY COMPONENTS

`CVNets` include efficient data sampling (Section 3.1) and training methods (Section 3.2), in addition to standard components (e.g., optimizers; Section 3.3), which are discussed below.

### 3.1 Data Samplers

`CVNets` offer data samplers with three sampling strategies: (1) single-scale with fixed batch size, (2) multi-scale with fixed batch size, and (3) multi-scale with variable batch size. These sampling strategies are visualized in Figure 2a and discussed below:

- **Single-scale with fixed batch size (SSc-FBS):** This method is the default sampling strategy in most deep learning frameworks (e.g., PyTorch, Tensorflow, and MixNet) and libraries built on top

of them (e.g., the `timm` library [21]). At the $t$-th training iteration, this method samples a batch of $b$ images per GPU with a pre-defined spatial resolution of height $H$ and width $W$.
- **Multi-scale with fixed batch size (MSc-FBS):** The SSc-FBS method allows a network to learn representations at a single scale (or resolution). However, objects in the real-world are composed at different scales. To allow a network to learn representations at multiple scales, MSc-FBS extends SSc-FBS to multiple scales [16]. Unlike the SSc-FBS method that takes a pre-defined spatial resolution as an input, this method takes a sorted set of $n$ spatial resolutions $\mathcal{S} = \{(H_1, W_1), (H_2, W_2), \cdots, (H_n, W_n)\}$ as an input. At the $t$-th iteration, this method randomly samples $b$ images per GPU of spatial resolution $(H_t, W_t) \in \mathcal{S}$.
- **Multi-scale with variable batch size (MSc-VBS):** Networks trained using the MSc-FBS methods are more robust to scale changes as compared to SSc-FBS [13]. However, depending on the maximum spatial resolution in $\mathcal{S}$, MSc-FBS methods may have a higher peak GPU memory utilization (see Figure 2c) as compared to SSc-FBS; causing out-of-memory errors on GPUs with limited memory. For example, MSc-FBS with $\mathcal{S} = \{(128, 128), (192, 192), (224, 224), (320, 320)\}$ and $b = 256$ would need about 2× more GPU memory (for images only) than SSc-FBS with a spatial resolution of $(224, 224)$ and $b = 256$. To address this memory issue, we extend MSc-FBS to variably-batch sizes in our previous work [13]. For a given sorted set of spatial resolutions $\mathcal{S} = \{(H_1, W_1), (H_2, W_2), \cdots, (H_n, W_n)\}$ and a batch size $b$ for a maximum spatial resolution of $(H_n, W_n)$, a spatial resolution $(H_t, W_t) \in \mathcal{S}$ with a batch size of $b_t = H_n W_n b / H_t W_t$ is sampled randomly at $t$-th training iteration on each GPU.

These samplers offer different training costs and performance for different models, as shown in Figure 2c. Compared to SSc-FBS and MSc-FBS, MSc-VBS is a memory-efficient sampler that speeds-up the training significantly while maintaining performance.
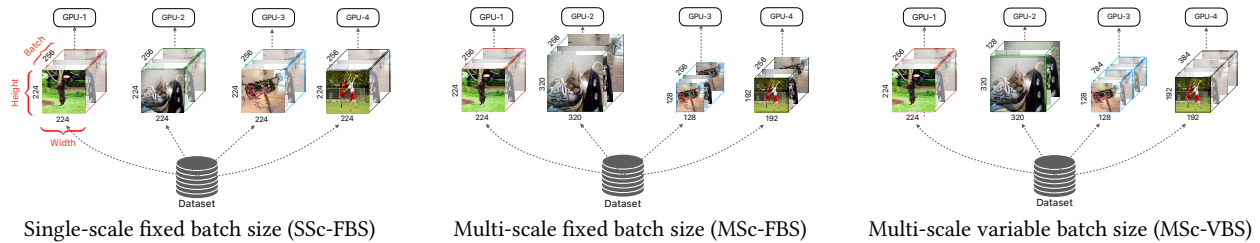
***Variably-sized video sampler.*** Samplers discussed above can be easily extended for videos. `CVNets` provide variably-sized sampler for videos, wherein different video-related input variables (e.g., number of frames, number of clips per video, and spatial size) can be controlled for learning space- and time-invariant representations.

### 3.2 Sample efficient training

Previous works [11, 12, 14] remove and re-weight data samples to reduce optimization updates (or number of forward passes) at the cost of performance degradation. Moreover, these methods are compute- and memory-intensive, and do not scale well to large models and datasets [12]. This work aims to reduce optimization updates with minimal or no performance degradation.

Models learn quickly during the initial phase of training (see Figure 2b). In other words, models can accurately classify many training data samples during earlier epochs and such samples do not contribute much to learning. Therefore, a natural question arises: *Can we remove such samples to reduce total optimization updates?*

`CVNets` answer this question with a simple heuristic method, which we call *Sample Efficient Training (SET)*. At each epoch, SET categorizes each sample as either *hard* or *easy*. To do so, SET uses a simple heuristic: if model predicts the training data sample correctly with a confidence greater than a pre-defined threshold $\tau$, then it is

Single-scale fixed batch size (SSc-FBS)    Multi-scale fixed batch size (MSc-FBS)    Multi-scale variable batch size (MSc-VBS)

**(a) Different samplers**



MobileNetv1-1.0      ResNet-50

**(b) Training and validation loss curves for models trained with different samplers**

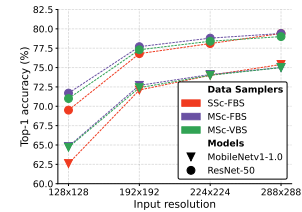| Sampler | # optim. updates | Max. GPU Memory | Train time (in sec) | Top-1 @ 224x224 (%) |
|---|---|---|---|---|
| **MobileNetv1-1.0** | | | | |
| SSc-FBS | 751k (1.00×) | **15.9 GB** (1.00×) | 122k (1.00×) | 73.95 ( 0.00) |
| MSc-FBS | 751k (1.00×) | 26.0 GB (1.64×) | 137k (1.12×) | **74.11** (+0.16) |
| MSc-VBS | **574k** (0.76×) | 18.1 GB (1.14×) | 87k (0.71×) | 74.05 (+0.10) |
| **ResNet-50** | | | | |
| SSc-FBS | 188k (1.00×) | **24.1 GB** (1.00×) | 43k (1.00×) | 78.12 ( 0.00) |
| MSc-FBS | 188k (1.00×) | 43.5 GB (1.81×) | 48k (1.12×) | **78.81** (+0.69) |
| MSc-VBS | **143k** (0.79×) | 28.0 GB (1.16×) | 36k (0.84×) | 78.44 (+0.32) |

**(c) Effect of different samplers on the training performance of different models.**



**(d) Effect on the performance of different models when evaluated at different input resolutions.**

**Figure 2: Effect of training deep learning models with different sampling methods on the ImageNet dataset. Models trained with MSc-VBS delivers similar performance, trains faster with fewer optimization updates, and generalizes better (higher train loss; similar validation loss) as compared to the ones trained with SSc-FBS and MSc-FBS. Notably, models trained with MSc-VBS require similar computational resources as SSc-FBS. Here, all models are trained with simple training recipes. For more results, see Appendix B.**



**(a) Optim. updates vs. confidence threshold $\tau$**    **(b) Top-1 accuracy vs. confidence threshold $\tau$**    **(c) Training samples vs. epoch at $\tau = 0.5$**

**Figure 3: Sample efficient training for the ResNet-50 model. SET reduces the optimizer updates (a) while maintaining performance (b). Fluctuations in (c) represents that easy samples were classified as hard and added back to training data.**

an easy sample and we remove it from the training data. At each epoch, model only trains using hard samples. Because of randomness in training due to data augmentation (e.g., random cropping), it is possible that region of interest corresponding to the object category may be partially (or not) present in model's input and

model may classify such inputs correctly. To make SET robust to such randomness, a training sample is classified easy only when it is classified correctly with a confidence greater than $\tau$ for a moving window of $w$ epochs. Because SET uses a window of $w$ epochs to determine easy samples, it is possible that some samples classified

| Model | Params | FLOPs | Top-1 (in %) | |
| --- | --- | --- | --- | --- |
| | | | CVNets (Ours) | Prev. |
| MobileNetv1-0.25 | 0.5 M | 46.3 M | **54.5** | 49.8 |
| MobileNetv1-1.0 | 4.2 M | 568.7 M | **74.1** | 70.9 |
| MobileNetv2-0.25 | 1.5 M | 50.5 M | **53.6** | – |
| MobileNetv2-1.0 | 3.5 M | 300.7 M | **72.9** | 72.0 |
| MobileNetv3-Large | 5.4 M | 210.7 M | 75.1 | 75.2 (TPU) 74.6 (GPU) |
| ResNet-101 (simple recipe) | 44.5 M | 7.7 G | 79.8 | 77.4[†] |
| ResNet-101 (adv. recipe) | 44.5 M | 7.7 G | 81.8 | 81.8[†] |
| ViT-Tiny | 5.7 M | 1.3 G | 72.9 | 72.2[★] |

**Table 1: Classification on the ImageNet dataset.**
[†]**Torchvision [15] requires** $1.3\times$ **more optimization updates (forward passes) as compared to CVNets; see Appendix A for details.** [★] **Results are from [20].**

easy earlier may be classified harder during the later training stages. SET adds such samples back to the training data (see Figure 3c).

Figure 3 shows results for ResNet-50 trained with and without SET using MSc-VBS. ResNet-50 without SET requires 22% more optimization updates while delivering similar performance; demonstrating the effectiveness of SET on top of MSc-VBS. Note that SET has an overhead. Therefore, the reduction in optimization updates do not translate to reduction in training time. We believe SET can serve as a baseline in this direction and inspire future research to improve training speed while maintaining performance.

## 3.3 Standard components

CVNets support different tasks (e.g., image classification, detection, segmentation), data augmentation methods (e.g., flipping, random resized crop, RandAug, and CutMix), datasets (e.g., ImageNet-1k/21k for image classification, Kinetics-400 for video classification, MS-COCO for object detection, and ADE20k for segmentation), optimizers (e.g., SGD, Adam, and AdamW), and learning rate annealing methods (e.g., fixed, cosine, and polynomial).

## 4 BENCHMARKS

CVNets support different visual recognition tasks, including classification, detection, and segmentation. We provide comprehensive benchmarks for standard methods along with pre-trained weights.

***Classification on ImageNet dataset.*** CVNets implement popular light- and heavy-weight image classification models. The performance of some of these models on the ImageNet dataset is shown in Table 1. With CVNets, we are able to achieve better performance (e.g., MobileNetv1/v2) or similar performance (ResNet-50/101) with fewer optimization updates (faster training). See Appendix A, including training recipe comparisons.

***Detection and segmentation.*** Similar to image classification, CVNets can be used to train standard detection and segmentation models with better performance. For example, SSD with ResNet-101 backbone trained with CVNets at a resolution of $384 \times 384$ delivers a 1.6% better mAP than the same model trained at a resolution of $512 \times 512$ as reported in [6]. Similarly, on the task of semantic segmentation on the ADE20k dataset using DeepLabv3 with MobileNetv2 as the backbone, CVNets delivers 1.1% better performance

than MMSegmentation library [1] with 2× fewer epochs and optimization updates. For more details, please see our benchmarking results at https://github.com/apple/ml-cvnets.

## 5 CONCLUSION

This paper introduces CVNets, a modular deep learning library for visual recognition tasks with high performance. In future, we plan to continue enhancing CVNets with novel and reproducible methods. We welcome contributions from the research and open-source community to support further innovation.

## REFERENCES

[1] MMSegmentation Contributors. 2020. MMSegmentation: OpenMMLab Semantic Segmentation Toolbox and Benchmark.

[2] Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le. 2020. Randaugment: Practical automated data augmentation with a reduced search space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*. 702–703.

[3] Jia Deng et al. 2009. ImageNet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*.

[4] Haoqi Fan et al. 2021. PyTorchVideo: A Deep Learning Library for Video Understanding. In *Proceedings of the 29th ACM International Conference on Multimedia*.

[5] Martín Abadi et al. 2015. TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems.

[6] Cheng-Yang Fu et al. 2017. Dssd: Deconvolutional single shot detector. *arXiv preprint arXiv:1701.06659* (2017).

[7] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.

[8] Elad Hoffer et al. 2019. Augment your batch: better training with larger batches. *arXiv preprint arXiv:1901.09335* (2019).

[9] Andrew Howard et al. 2019. Searching for mobilenetv3. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 1314–1324.

[10] Andrew G Howard et al. 2017. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861* (2017).

[11] Krishnateja Killamsetty et al. 2021. GLISTER: Generalization based Data Subset Selection for Efficient and Robust Learning. *Proceedings of the AAAI Conference on Artificial Intelligence* 35 (May 2021), 8110–8118.

[12] Krishnateja Killamsetty et al. 2021. GRAD-MATCH: Gradient Matching based Data Subset Selection for Efficient Deep Model Training. In *Proceedings of the 38th International Conference on Machine Learning*.

[13] Sachin Mehta and Mohammad Rastegari. 2022. MobileViT: Light-weight, General-purpose, and Mobile-friendly Vision Transformer. In *International Conference on Learning Representations*.

[14] Baharan Mirzasoleiman, Jeff Bilmes, and Jure Leskovec. 2020. Coresets for Data-efficient Training of Machine Learning Models. In *Proceedings of the 37th International Conference on Machine Learning*. PMLR.

[15] Adam Paszke et al. 2019. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In *Advances in Neural Information Processing Systems 32*. 8024–8035.

[16] Joseph Redmon and Ali Farhadi. 2017. YOLO9000: better, faster, stronger. In *Proceedings of the IEEE conference on computer vision and pattern recognition*.

[17] Olga Russakovsky et al. 2015. Imagenet large scale visual recognition challenge. *International journal of computer vision* 115, 3 (2015), 211–252.

[18] Mark Sandler et al. 2018. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 4510–4520.

[19] Hugo Touvron, Andrea Vedaldi, Matthijs Douze, and Hervé Jégou. 2019. Fixing the train-test resolution discrepancy. *Advances in neural information processing systems* 32 (2019).

[20] Hugo Touvron et al. 2021. Training data-efficient image transformers & distillation through attention. In *International Conference on Machine Learning*. PMLR, 10347–10357.

[21] Ross Wightman, Hugo Touvron, and Hervé Jégou. 2021. Resnet strikes back: An improved training procedure in timm. *arXiv preprint arXiv:2110.00476* (2021).

[22] Sangdoo Yun et al. and Han. 2019. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *Proceedings of the IEEE/CVF international conference on computer vision*. 6023–6032.

[23] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. 2017. mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412* (2017).

[24] Zhun Zhong et al. 2020. Random erasing data augmentation. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 34. 13001–13008.

## A  COMPARISON OF TRAINING PROCEDURES

Training ResNet-50 [7] on the ImageNet dataset [17] is widely used for comparing training procedures [15, 21]. Table 2 compares the performance of CVNets' training recipe with previous methods. We can make following observations:

- **Simple training recipes** are useful for research in resource-constrained environments, where number of GPUs available for training are limited. With simple recipes, ResNet-50 trained with CVNets achieves better performance as compared to previous works (e.g., ResNet, Torchvision, and FixRes in Table 2).
- **Advanced training recipes** are useful for achieving state-of-the-art performance, as they train models for longer with better augmentation methods. With advanced recipes, ResNet-50 trained with CVNets achieve similar performance to previous works, but with 1.3× fewer optimization updates (Torchvision-Adv vs. CVNets' recipe; see highlighted text in Table 2).

*Comparison with light-weight model training recipes.* Table 3 and Table 4 compares the performance of MobileNetv2 and MobileNetv3 models trained with different training recipes. With simple training recipes, models trained with CVNets achieve similar or better performance as compared to other libraries (e.g., TensorflowLite or Torchvision). For example, both CVNets and Torchvision are able to train MobileNetv3 with reported performance of 75.2% [9]. However, CVNets achieve the reported performance with basic data augmentation while Torchvision achieves this with more data augmentation (e.g., RandAugment, Cutmix, MixUp, and RandomErase).

## B  EFFECT OF DIFFERENT SAMPLERS

CVNets implement three different samplers: (1) single-scale with fixed batch size (SSc-FBS), (2) multi-scale with fixed batch size (MSc-FBS), and (3) multi-scale with variable batch size (MSc-VBS). For different models, these samplers may offer different training costs and performance. Table 5 compares the effect of different samplers on training efficiency and performance on ImageNet. Clearly, MSc-VBS is a memory-efficient sampler that speed-up the training significantly while delivering similar performance to other two samplers.

| | ResNet [7] | Torchvision [15] | FixRes [19] | DeiT [20] | timm-A3 [21] | timm-A1 [21] | Torchvision-Adv [15] | CVNets (Ours) | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Train res | 224 | 224 | 224 | 224 | 160 | 224 | 176 | 224 | {128, 192, 224, 288, 320} | | |
| Test res | 224 | 224 | 224 | 224 | 224 | 224 | 224 | 224 | 224 | 224 | 224 |
| Test crop ratio | 0.875 | 0.875 | 0.875 | 0.875 | 0.95 | 0.95 | 0.95 | 0.875 | 0.875 | 0.875 | 0.95 |
| Data sampler | SSc-FBS | SSc-FBS | SSc-FBS | SSc-FBS | SSc-FBS | SSc-FBS | SSc-FBS | SSc-FBS | MSc-FBS | MSc-VBS | MSc-VBS |
| Epochs | 90 | 90 | 120 | 300 | 100 | 600 | 600 | 150 | 150 | 150 | 600 |
| Batch size | 256 | 256 | 512 | 1024 | 2048 | 2048 | 1024 | 1024 | 1024 | 1024 | 1024 |
| # optim. updates | 450k | 450k | 300k | 375k | 63k | 375k | 750k | 188k | 188k | 143k | 571k |
| Max. LR | 0.1 | 0.1 | 0.2 | 0.001 | 0.008 | 0.005 | 0.5 | 0.4 | 0.4 | 0.4 | 0.4 |
| LR Annealing | step | step | step | cosine | cosine | cosine | cosine | cosine | cosine | cosine | cosine |
| Weight decay | $10^{-4}$ | $10^{-4}$ | $10^{-4}$ | 0.05 | 0.02 | 0.01 | $2^{-5}$ | $10^{-4}$ | $10^{-4}$ | $10^{-4}$ | $10^{-4}$ |
| Warmup epochs | 0 | 0 | 0 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 |
| Optimizer | SGD-M | SGD-M | SGD-M | AdamW | LAMB | LAMB | SGD-M | SGD-M | SGD-M | SGD-M | SGD-M |
| Loss fn. | CE | CE | CE | CE | BCE | BCE | CE | CE | CE | CE | CE |
| EMA | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Label smoothing $\epsilon$ | ✗ | ✗ | ✗ | 0.1 | ✗ | 0.1 | 0.1 | 0.1 | 0.1 | 0.1 | 0.1 |
| Stoch. Depth | ✗ | ✗ | ✗ | 0.1 | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ |
| Mixed Precision | ✗ | ✗ | ✗ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| H. flip | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| RRC | ✗ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Repeated Aug [8] | ✗ | ✗ | ✗ | ✓ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ |
| Rand Augment [2] | ✗ | ✗ | ✗ | ✓ | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ | ✓ |
| Mixup [23] | ✗ | ✗ | ✗ | ✓ | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ | ✓ |
| Cutmix [22] | ✗ | ✗ | ✗ | ✓ | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ | ✓ |
| Random Erase [24] | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ |
| ColorJitter | ✗ | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| PCA lighting | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| Top-1 acc. | 75.3% | 76.1% | 77.0% | 78.4% | 78.1% | **80.4%** | **80.4%** | 78.1% | 78.8% | 78.4% | **80.4%** |

**Table 2: Comparison of ResNet-50 training recipes and performance in different papers with CVNets (ours).**
Link to Torchvision-Adv recipe: https://pytorch.org/blog/how-to-train-state-of-the-art-models-using-torchvision-latest-primitives/. Last accessed on May 9, 2022.

| | Torchvision [15] | TensorflowLite [18] | CVNets (Ours) | | |
|---|---|---|---|---|---|
| Train res | 224 | 224 | 224 | {128, 192, 224, 288, 320} | {128, 192, 224, 288, 320} |
| Test res | 224 | 224 | 224 | 224 | 224 |
| Test crop ratio | 0.875 | 0.875 | 0.875 | 0.875 | 0.875 |
| Data sampler | SSc-FBS | SSc-FBS | SSc-FBS | MSc-FBS | MSc-VBS |
| Epochs | 300 | 450 | 300 | 300 | 300 |
| Batch size | 256 | 768 | 1024 | 1024 | 1024 |
| # optim. updates | 1500k | 700k | 375k | 375k | 288k |
| Max. LR | 0.36 | 0.045 | 0.4 | 0.4 | 0.4 |
| LR Annealing | poly | poly | cosine | cosine | cosine |
| Weight decay | $4^{-5}$ | $4^{-5}$ | $4^{-5}$ | $4^{-5}$ | $4^{-5}$ |
| Warmup epochs | 0 | 0 | 5 | 5 | 5 |
| Optimizer | SGD-M | RMSProp | SGD-M | SGD-M | SGD-M |
| Loss fn. | CE | CE | CE | CE | CE |
| EMA | ✗ | ✓ | ✓ | ✓ | ✓ |
| Label smoothing $\epsilon$ | 0.0 | 0.1 | 0.1 | 0.1 | 0.1 |
| H. flip | ✓ | ✓ | ✓ | ✓ | ✓ |
| RRC | ✓ | ✓ | ✓ | ✓ | ✓ |
| Top-1 acc. | 71.8% | 72.0% | 73.3% | **73.4%** | 72.9% |

**Table 3: Comparison of MobileNetv2-1.0 training recipes and performance in different libraries with CVNets (ours).**

| | Torchvision-adv [15] | TensorflowLite [9] | | CVNets (Ours) | |
|---|---|---|---|---|---|
| Train res | 224 | 224 | 224 | {128, 192, 224, 288, 320} | {128, 192, 224, 288, 320} |
| Test res | 224 | 224 | 224 | 224 | 224 |
| Test crop ratio | 0.875 | 0.875 | 0.875 | 0.875 | 0.875 |
| Data sampler | SSc-FBS | SSc-FBS | SSc-FBS | MSc-FBS | MSc-VBS |
| Epochs | 600 | 840$^{\dagger}$ | 300 | 300 | 300 |
| Batch size | 1024 | 1536 | 2048 | 2048 | 2048 |
| # optim. updates | 750k | 700k | 188k | 188k | 144k |
| Max. LR | 0.064 | 0.16 | 0.4 | 0.4 | 0.4 |
| LR Annealing | poly | poly | cosine | cosine | cosine |
| Weight decay | $1^{-5}$ | $1^{-5}$ | $4^{-5}$ | $4^{-5}$ | $4^{-5}$ |
| Warmup epochs | 0 | 5 | 5 | 5 | 5 |
| Optimizer | RMSProp | RMSProp | SGD-M | SGD-M | SGD-M |
| Loss fn. | CE | CE | CE | CE | CE |
| EMA | ✗ | ✓ | ✓ | ✓ | ✓ |
| Label smoothing $\epsilon$ | 0.0 | 0.1 | 0.1 | 0.1 | 0.1 |
| H. flip | ✓ | ✓ | ✓ | ✓ | ✓ |
| RRC | ✓ | ✓ | ✓ | ✓ | ✓ |
| Rand Augment | ✓ | ✗ | ✗ | ✗ | ✗ |
| Random Erase | ✓ | ✗ | ✗ | ✗ | ✗ |
| Cutmix | ✓ | ✗ | ✗ | ✗ | ✗ |
| MixUp | ✓ | ✗ | ✗ | ✗ | ✗ |
| Top-1 acc. | 75.3% | 74.6% | 75.3% | **75.3%** | 75.1% |

**Table 4: Comparison of MobileNetv3-Large training recipes and performance in different libraries with CVNets (ours).** $^{\dagger}$ **Epochs are computed based on the training recipe in the official repository.**

**Link to official training recipe: https://github.com/tensorflow/models/tree/master/research/slim/nets/mobilenet. Last accessed on May 9, 2022.**

| Model | Sampler | | | Peak GPU Memory (in GB) | Train time (in sec) | # optim. updates (in thousands) | Top-1 (in %) |
|---|---|---|---|---|---|---|---|
| | SSc-FBS | MSc-FBS | MSc-VBS | | | | |
| **ResNet-34** | ✓ | ✗ | ✗ | 15.7 $_{(1.00\times)}$ | 35k $_{(1.00\times)}$ | 188k $_{(1.00\times)}$ | 75.1 $_{(\ 0.0)}$ |
| | ✗ | ✓ | ✗ | 24.2 $_{(1.54\times)}$ | 40k $_{(1.14\times)}$ | 188k $_{(1.00\times)}$ | 75.2 $_{(+0.1)}$ |
| | ✗ | ✗ | ✓ | 17.8 $_{(1.13\times)}$ | 31k $_{(0.89\times)}$ | 143k $_{(0.76\times)}$ | 74.8 $_{(-0.3)}$ |
| **ResNet-50** | ✓ | ✗ | ✗ | 24.1 $_{(1.00\times)}$ | 43k $_{(1.00\times)}$ | 188k $_{(1.00\times)}$ | 78.1 $_{(\ 0.0)}$ |
| | ✗ | ✓ | ✗ | 43.5 $_{(1.80\times)}$ | 48k $_{(1.12\times)}$ | 188k $_{(1.00\times)}$ | 78.8 $_{(+0.7)}$ |
| | ✗ | ✗ | ✓ | 28.0 $_{(1.16\times)}$ | 36k $_{(0.84\times)}$ | 143k $_{(0.76\times)}$ | 78.4 $_{(+0.3)}$ |
| **ResNet-101** | ✓ | ✗ | ✗ | 31.1 $_{(1.00\times)}$ | 57k $_{(1.00\times)}$ | 188k $_{(1.00\times)}$ | 79.6 $_{(\ 0.0)}$ |
| | ✗ | ✓ | ✗ | 61.7 $_{(1.98\times)}$ | 64k $_{(1.12\times)}$ | 188k $_{(1.00\times)}$ | 80.0 $_{(+0.4)}$ |
| | ✗ | ✗ | ✓ | 35.3 $_{(1.14\times)}$ | 45k $_{(0.79\times)}$ | 143k $_{(0.76\times)}$ | 79.8 $_{(+0.2)}$ |
| **MobileNetv1-0.25** | ✓ | ✗ | ✗ | 10.0 $_{(1.00\times)}$ | 122k $_{(1.00\times)}$ | 751k $_{(1.00\times)}$ | 55.5 $_{(\ 0.0)}$ |
| | ✗ | ✓ | ✗ | 13.1 $_{(1.31\times)}$ | 125k $_{(1.02\times)}$ | 751k $_{(1.00\times)}$ | 54.4 $_{(-1.1)}$ |
| | ✗ | ✗ | ✓ | 10.5 $_{(1.05\times)}$ | 85k $_{(0.69\times)}$ | 574k $_{(0.76\times)}$ | 54.5 $_{(-1.0)}$ |
| **MobileNetv1-0.50** | ✓ | ✗ | ✗ | 12.0 $_{(1.00\times)}$ | 118k $_{(1.00\times)}$ | 751 k $_{(1.00\times)}$ | 66.6 $_{(\ 0.0)}$ |
| | ✗ | ✓ | ✗ | 16.2 $_{(1.35\times)}$ | 134k $_{(1.14\times)}$ | 751k $_{(1.00\times)}$ | 66.5 $_{(-0.1)}$ |
| | ✗ | ✗ | ✓ | 12.3 $_{(1.03\times)}$ | 86k $_{(0.73\times)}$ | 574k $_{(0.76\times)}$ | 65.9 $_{(-0.7)}$ |
| **MobileNetv1-0.75** | ✓ | ✗ | ✗ | 13.9 $_{(1.00\times)}$ | 126k $_{(1.00\times)}$ | 751k $_{(1.00\times)}$ | 71.5 $_{(\ 0.0)}$ |
| | ✗ | ✓ | ✗ | 21.8 $_{(1.57\times)}$ | 129k $_{(1.02\times)}$ | 751k $_{(1.00\times)}$ | 71.7 $_{(+0.2)}$ |
| | ✗ | ✗ | ✓ | 15.3 $_{(1.10\times)}$ | 114k $_{(0.90\times)}$ | 574k $_{(0.76\times)}$ | 71.4 $_{(-0.1)}$ |
| **MobileNetv1-1.0** | ✓ | ✗ | ✗ | 15.9 $_{(1.00\times)}$ | 122k $_{(1.00\times)}$ | 751k $_{(1.00\times)}$ | 74.0 $_{(\ 0.0)}$ |
| | ✗ | ✓ | ✗ | 26.0 $_{(1.64\times)}$ | 137k $_{(1.12\times)}$ | 751k $_{(1.00\times)}$ | 74.1 $_{(+0.1)}$ |
| | ✗ | ✗ | ✓ | 18.1 $_{(1.14\times)}$ | 87k $_{(0.71\times)}$ | 574k $_{(0.76\times)}$ | 74.1 $_{(+0.1)}$ |
| **MobileNetv2-0.25** | ✓ | ✗ | ✗ | 21.3 $_{(1.00\times)}$ | 119k $_{(1.00\times)}$ | 375k $_{(1.00\times)}$ | 54.9 $_{(\ 0.0)}$ |
| | ✗ | ✓ | ✗ | 33.5 $_{(1.57\times)}$ | 135k $_{(1.13\times)}$ | 375k $_{(1.00\times)}$ | 54.0 $_{(-0.9)}$ |
| | ✗ | ✗ | ✓ | 25.5 $_{(1.20\times)}$ | 75k $_{(0.63\times)}$ | 288k $_{(0.77\times)}$ | 53.6 $_{(-1.3)}$ |
| **MobileNetv2-0.50** | ✓ | ✗ | ✗ | 26.4 $_{(1.00\times)}$ | 120k $_{(1.00\times)}$ | 375k $_{(1.00\times)}$ | 65.8 $_{(\ 0.0)}$ |
| | ✗ | ✓ | ✗ | 45.2 $_{(1.71\times)}$ | 138k $_{(1.15\times)}$ | 375k $_{(1.00\times)}$ | 65.6 $_{(-0.2)}$ |
| | ✗ | ✗ | ✓ | 30.9 $_{(1.17\times)}$ | 99k $_{(0.82\times)}$ | 288k $_{(0.77\times)}$ | 65.3 $_{(-0.5)}$ |
| **MobileNetv2-0.75** | ✓ | ✗ | ✗ | 34.3 $_{(1.00\times)}$ | 127k $_{(1.00\times)}$ | 375k $_{(1.00\times)}$ | 70.7 $_{(\ 0.0)}$ |
| | ✗ | ✓ | ✗ | 65.3 $_{(1.90\times)}$ | 139k $_{(1.09\times)}$ | 375k $_{(1.00\times)}$ | 70.8 $_{(+0.1)}$ |
| | ✗ | ✗ | ✓ | 41.0 $_{(1.20\times)}$ | 85k $_{(0.67\times)}$ | 288k $_{(0.77\times)}$ | 70.4 $_{(-0.3)}$ |
| **MobileNetv2-1.0** | ✓ | ✗ | ✗ | 36.8 $_{(1.00\times)}$ | 128k $_{(1.00\times)}$ | 375k $_{(1.00\times)}$ | 73.3 $_{(\ 0.0)}$ |
| | ✗ | ✓ | ✗ | 70.7 $_{(1.92\times)}$ | 142k $_{(1.11\times)}$ | 375k $_{(1.00\times)}$ | 73.4 $_{(+0.1)}$ |
| | ✗ | ✗ | ✓ | 43.9 $_{(1.19\times)}$ | 86k $_{(0.67\times)}$ | 288k $_{(0.77\times)}$ | 72.9 $_{(-0.4)}$ |
| **MobileNetv3-Small** | ✓ | ✗ | ✗ | 14.5 $_{(1.00\times)}$ | 63k $_{(1.00\times)}$ | 188k $_{(1.00\times)}$ | 67.1 $_{(\ 0.0)}$ |
| | ✗ | ✓ | ✗ | 24.7 $_{(1.70\times)}$ | 71k $_{(1.13\times)}$ | 188k $_{(1.00\times)}$ | 66.8 $_{(-0.3)}$ |
| | ✗ | ✗ | ✓ | 15.4 $_{(1.06\times)}$ | 57k $_{(0.90\times)}$ | 144k $_{(0.77\times)}$ | 66.7 $_{(-0.4)}$ |
| **MobileNetv3-Large** | ✓ | ✗ | ✗ | 25.7 $_{(1.00\times)}$ | 71k $_{(1.00\times)}$ | 188k $_{(1.00\times)}$ | 75.3 $_{(\ 0.0)}$ |
| | ✗ | ✓ | ✗ | 47.3 $_{(1.84\times)}$ | 77k $_{(1.08\times)}$ | 188k $_{(1.00\times)}$ | 75.4 $_{(+0.1)}$ |
| | ✗ | ✗ | ✓ | 30.5 $_{(1.19\times)}$ | 60k $_{(0.85\times)}$ | 144k $_{(0.77\times)}$ | 75.1 $_{(-0.2)}$ |

Table 5: Effect of different samplers on the performance of different models.