

Searching for MobileNetV3

Andrew Howard¹ Mark Sandler¹ Grace Chu¹ Liang-Chieh Chen¹ Bo Chen¹ Mingxing Tan²
 Weijun Wang² Yukun Zhu¹ Romain Pang² Vijay Vasudevan² Quoc V. Le² Hartwig Adam¹
¹Google AI²Google Brain
 {howard, sandler, cych, lchen, hechen, tammingsing, weijun, yukun, rpang, vv, qvl, hadam}@google.com

Abstract

We present the next generation of MobileNets based on a combination of complementary search techniques as well as a novel architecture design. MobileNetV3 is tuned to mobile phone CPUs through a combination of hardware-aware network architecture search (NAS) complemented by the NeckAdapt algorithm and then subsequently improved through novel architecture advances. This paper starts the exploration of how automated search algorithms and network design can work together to harness complementary approaches improving the overall state of the art. Through this process we create two new MobileNet models for release: MobileNetV3-Large and MobileNetV3-Small which are targeted for high and low resource use cases. These models are then adapted and applied to the tasks of object detection and semantic segmentation. For the task of semantic segmentation (or any dense pixel prediction), we propose a new efficient segmentation decoder Lite Reduced Atrous Spatial Pyramid Pooling (LR-ASPP). We achieve new state of the art results for mobile classification, detection and segmentation. MobileNetV3-Large is 3.2% more accurate on ImageNet classification while reducing latency 19% compared to MobileNetV2. MobileNetV3-Small is 1.9% more accurate compared to a MobileNetV2 model while reducing latency. MobileNetV3-Large detection faster at roughly the same accuracy as MobileNetV2. MobileNetV3-Large is 1.8% more accurate than MobileNetV2 on COCO.

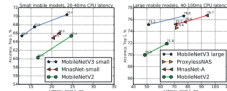


Figure 1. The trade-off between Pixel 1 latency and top-1 ImageNet accuracy. All models use the input resolution 224. V3 large and V3 small use multipliers 0.75, 1 and 1.25 to show optimal frontier. All latencies were measured on a single large core of the same device using TensorFlow Lite. MobileNetV3-Small and Large are our proposed next-generation mobile models.

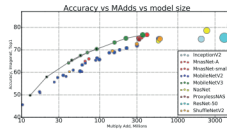


Figure 2. The trade-off between mAP and top-1 accuracy. This allows to compare models that were targeted different hardware

Everything you need to know about MobileNetV3



Vardit Jain · Follow

Published in Towards Data Science · 8 min read · Nov 22, 2019

180 2

When MobileNet V1 came in 2017, it essentially started a new section of deep learning research in computer vision, i.e. coming up with models that can run in embedded systems. This led to several important works including but not limited to ShuffleNet(V1 and V2), MNasNet, CondenseNet, EffNet, among others. Somewhere in between came the second version of MobileNet as well last year. Now, this year's iteration gives us the third version of MobileNet called MobileNetV3. This story is a review of MobileNetV3 from Google that was presented at ICCV in Seoul, South Korea this year.

Contents:

1. Efficient Mobile Building Blocks
2. Neural Architecture Search for Block-Wise Search
3. NetAdapt for Layer wise search
- Network Improvements — Layer removal and H-swish
- Structure

AI Label

AI LABEL
 MobileNetV3Small
 infer ImageNet (ILSVRC2012)
 Scan for further information

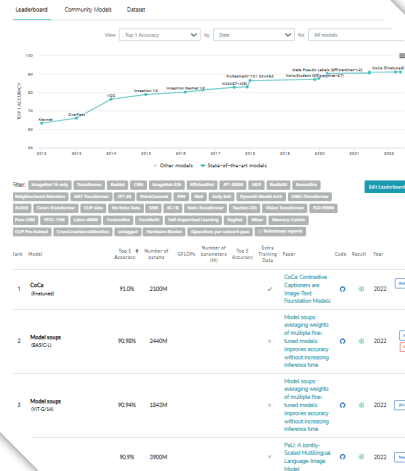


A100 x8 - TensorFlow 2.8.0

608.827 [mW]
Power Draw per Inference59.2788 [%]
Corrupted Robustness63.2031 [%]
Top1 Accuracy1.28563 [s]
Running Time per Inference

Benchmark (PWC)

Image Classification on ImageNet



Platform (PWC)

MobileNetV3

Introduced by Howard et al. in Searching for MobileNetV3

MobileNetV3 is a convolutional neural network that is tuned to mobile phone CPUs through a combination of hardware-aware network architecture search (NAS) complemented by the NeckAdapt algorithm, and then subsequently improved through novel architecture advances. Advances include (1) complementary search techniques, (2) new efficient versions of nonlinearities practical for the mobile setting, (3) new efficient network designs.

The network designs include the use of a hard switch activation and squeeze-and-excitation modules in the MBConv blocks.

Source: Searching for MobileNetV3

Read Paper See Code

Papers

Search for a paper or author

Paper	Code	Results	Date	Stars
Searching for MobileNetV3 Howard, Andrew Howard, Weijun Wang, Mingxing Tan, Quoc V. Le, Mark Sandler, Yukun Zhu, Vijay Vasudevan, Romain Pang, Bo Chen, Grace Chu, Liang-Chieh Chen			6 Mar 2019	76,787
Hardware-Efficient Ghost Module via Re-parameterization Wang, Jian Dong, Chengpeng Chen, Huan Zeng			11 Nov 2022	30,617

This model is an implementation of MobileNetV3-Small found [here](#). This repository provides to run MobileNetV3-Small on Qualcomm® devices. More details on model performance across various devices, can be found [here](#).

Model Details

- Model Type: Image classification
- Model Stats:
 - Model checkpoint: Imagenet
 - Input resolution: 224x224
 - Number of parameters: 2.54M
 - Model size: 9.72 MB

Device	Chipset	Target Runtime	Inference Time (ms)	Peak Memory Range (MB)	Precision	Primary Compute Unit	Target Model
Samsung Galaxy S23 Ultra (Android 13)	Snapdragon® 8 Gen 2	TfLite	0.844 ms	0 - 2 MB	FP16	NPU	MobileNetV3-Small.tflite
Samsung Galaxy S23 Ultra (Android 13)	Snapdragon® 8 Gen 2	QNN Library	0.879 ms	1 - 5 MB	FP16	NPU	MobileNetV3-Small.qnn

Installation

Model can be installed as a Python package via pip.

```
! ! pip install oai-hub-models
```

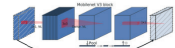


Figure 4. MobileNetV3 is a squeeze-and-excite [35]. In contrast with [35] we apply the squeeze-and-excite in the residual layers. We use different nonlinearity depending on the layers, see section 3.2 for details.

