# **On-Device Image Classification with Proxyless Neural Architecture Search and Quantization-Aware Fine-tuning**

Massachusetts Institute of Technology

ProxylessNAS: Direct Neural Architecture Search on Target Task and Hardware, ICLR'19 On-Device Image Classification with Proxyless Neural Architecture Search and Quantization-Aware Fine-tuning, ICCV Workshop'2019



Han Cai, Tianzhe Wang, Zhanghao Wu, Kuan Wang, Ji Lin, Song Han





# **From Manual Design to Automatic Design**



### Use Human Expertise

Manual Architecture Design

VGGNets **Inception Models** ResNets DenseNets

Computational Resources

. . . .

ProxylessNAS: Direct Neural Architecture Search on Target Task and Hardware, ICLR'19 On-Device Image Classification with Proxyless Neural Architecture Search and Quantization-Aware Fine-tuning, ICCV Workshop'2019





Use Machine Learning (NAS)

**Automatic** Architecture Search

Reinforcement Learning **Neuro-evolution Bayesian Optimization** Monte Carlo Tree Search

. . .





# **From General Design to Specialized CNN**

### **Previous Paradigm:** One CNN for all platforms.





ProxylessNAS: Direct Neural Architecture Search on Target Task and Hardware, ICLR'19 On-Device Image Classification with Proxyless Neural Architecture Search and Quantization-Aware Fine-tuning, ICCV Workshop'2019

### Sub-optimal, different in:

- Degree of parallelism
- Cache size
- Memory BW
- . . .





### **Previous Paradigm:** One CNN for all platforms.





ProxylessNAS: Direct Neural Architecture Search on Target Task and Hardware, ICLR'19 On-Device Image Classification with Proxyless Neural Architecture Search and Quantization-Aware Fine-tuning, ICCV Workshop'2019



### **Proxyless NAS:** Customize CNN for each platform.





## **Conventional NAS: Computation Expensive**, thus Proxy-Based



### Current neural architecture search (NAS) is **VERY EXPENSIVE**.

- NASNet: 48,000 GPU hours  $\approx$  5 years on single GPU
- DARTS: 100Gb GPU memory\* ≈ 9 times of modern GPU

### Therefore, previous work have to utilize proxy tasks:

- CIFAR-10 -> ImageNet
- Search a block -> stack to build a full net
- Fewer epochs training -> full training



\*if directly search on ImageNet, like us







## **Conventional NAS: Computation Expensive**, thus Proxy-Based



### **Limitations of Proxy**

- **Suboptimal** for the target task
- Blocks are forced to **share the same structure**.
- Cannot optimize for **specific hardware**.







# **Proxyless, Save GPU Hours by 200x**



**Goal: Directly** learn architectures on the **target task** and **hardware**. We achieved by

2. Cooperating hardware feedback (e.g. latency) into the search process.





1. Reducing the cost of NAS (GPU hours and memory) to the same level of regular training.





# **Save GPU Hours**



### Simplify NAS to be a **single training process** of a over-parameterized network. Build the cumbersome network with all candidate paths. No meta controller.







# **Save GPU Memory**



# Reduce the memory footprint from O(N) to O(1).



- **Binarize** the architecture parameters and allow only **one path of activation** in memory. We propose gradient-based and RL methods to update the architecture parameters.





## **Direct Search on Target Hardware:** Making Latency Differentiable



- Mobile farm infrastructure is expensive and slow.
- Use the latency estimation model as an economical alternative
- Optimize during search stage use **Gradient**.







## Efficiently search a model



ProxylessNAS: Direct Neural Architecture Search on Target Task and Hardware, ICLR'19 On-Device Image Classification with Proxyless Neural Architecture Search and Quantization-Aware Fine-tuning, ICCV Workshop'2019



## Search an efficient model







(1) The history of finding efficient Mobile model





(3) The history of finding efficient GPU model ProxylessNAS: Direct Neural Architecture Search on Target Task and Hardware, ICLR'19 On-Device Image Classification with Proxyless Neural Architecture Search and Quantization-Aware Fine-tuning, ICCV Workshop'2019



# The History of Architectures

(2) The history of finding efficient CPU model

Epoch-00





# **Results for LPIRC**

Model	Setting	Accuracy	Latency
MoblieNetV2	224-0.5	63.7%(65.4%)	28ms
MobileNetV2	192-0.75	67.4%(68.7%)	36ms
MobileNetV2	160-1.0	67.4%(68.8%)	31ms
ProxylessNAS	224-0.5	65.7%(67.0%)	31ms
ProxylessNAS	160-1.0	<b>69.2%</b> (70.3%)	35ms

Table 1. Results of 8-bit model using different preprocessing, the number in the bracket denotes the full-precision model's top-1 accuracy on ImageNet The latency is directly measured on Google Pixel 2. It takes only 200 GPU hours to find the specialized model with ProxylessNAS in the table.







## **Open-source**

### Both search code and models are released on Github:

- # https://github.com/MIT-HAN-LAB/ProxylessNAS
- from proxyless nas import \*
- net = proxyless cpu(pretrained=True)
- net = proxyless gpu(pretrained=True)
- net = proxyless mobile(pretrained=True)









## **Open-source**

### ProxylessNAS is available on PyTorch Hub:

import torch target platform = 'proxyless mobile'





# https://pytorch.org/hub/pytorch\_vision\_proxylessnas

- net = torch.hub.load('mit-han-lab/ProxylessNAS',
  - target platform, pretrained=True)





# Hardware, AI and Neural-nets

songhan@mit.edu

# Thank you!

### 16