

NetAdapt: Platform-Aware Neural Network Adaptation for Mobile Applications Tien-Ju Yang¹, Andrew Howard², Bo Chen², Xiao Zhang², Alec Go², Mark Sandler², Vivienne Sze¹, and Hartwig Adam² (**700**] ¹Massachusetts Institute of Technology, ²Google Inc.

Introduction

- Common problems of network simplification methods:
 - Require **manually** choosing the simplified architecture
 - Guided by **indirect metrics** (e.g., # of MACs), which may not be good approximations to the direct metrics (e.g., latency)
- NetAdapt
 - Automatically and progressively finds the simplified architecture
 - Incorporates **direct metrics** in the network simplification process by using real measurements from the target platform
 - **Improves** the accuracy-latency trade-offs of highly efficient neural networks, MobileNet V1/V2



Fast Resource Consumption Estimation

- Taking measurements can be slow and difficult to parallelize due to the limited number of available devices
- The network latency can be estimated by the sum of the latency of each layer





Real Latency vs. Estimated Latency

NetAdapt Algorithm

Problem formulation:



Experiment Results

NetAdapt is effective on both small and large MobileNet V1 and both

Network	Top-1 A	ccuracy (%)	Latency (ms)		
75% MobileNetV2 (224) [18]	69.8	(+0)	64.5	(100%)	
NetAdapt (Similar Latency)	70.9	(+1.1)	63.6	(99%)	
NetAdapt (Similar Accuracy)	70.2	(+0.4)	55.5	(86%)	

Network	Top-1 A	ccuracy (%)	# of MA	\mathbf{Cs} (×10 ⁶)	Laten	cy (ms)
25% MobileNetV1 (128) [9]	45.1	(+0)	13.6	(100%)	4.65	(100%)
MorphNet $[5]$	46.0	(+0.9)	15.0	(110%)	6.52	(140%)
NetAdapt	46.3	(+1.2)	11.0	(81%)	6.01	(129%)
75% MobileNetV1 (224) [9]	68.8	(+0)	325.4	(100%)	69.3	(100%)
ADC [8]	69.1	(+0.3)	304.2	(93%)	79.2	(114%)
NetAdapt	69.1	(+0.3)	284.3	(87%)	74.9	(108%)