Machine Learning @ Facebook Understanding Inference at the Edge

IEEE International Symposium on High Performance Computer Architecture (HPCA-2019)



Carole-Jean Wu

Research Scientist @ Facebook Al Infra Research







i co

From data centers to the edge

Minimizing network bandwidth

Improving response latency

Exploiting features available only at the edge



Keypoints Segmentation

Augmented Reality with Smart Camera





Applied Machine Learning at Facebook: A Datacenter Infrastructure Perspective. Hazelwood et al. HPCA-2018.





Unique Challenges for Edge Inference

 \sim

Feature-rich edge inference is enabled by the ever increasing mobile performance

Increasing core counts leads to theoretical peak performance increase. But, when looking at the entire ecosystem, the theoretical peak performance is a widespread.



DELIVERING CONSISTENT INFERENCE PERFORMANCE IS CHALLENGING



FRAGMENTED SMARTPHONE ECOSYSTEM POSES UNIQUE CHALLENGES FOR EDGE INFERENCE



Introduction:

& Unique Challenges for

Edge Inference

Machine Learning @ FB



Runs on



Horizontal Integration: Making Inference on Smartphones

Vertical Integration:

Processing Inference for

Oculus VR



Inference in the Wild: Performance Variability

9





i de la companya de l

What is Challenging for Mobile Inference?



Fragmentation

There is no standard mobile SoC to optimize for. Mobile CPUs Show Little Diversity



Performance

The Performance Difference between a Mobile CPU and GPU is Narrow



Programmability

Programmability is a Primary Roadblock for Using Mobile Coprocessors







LESS THAN 15% SMARTPHONES HAVE A GPU THAT IS 3 TIMES AS POWERFUL AS ITS CPU





Quantitative Approach to Mobile Inference Designs

State of the Practice for Mobile Inference is Using CPUs



FRAGMENTATION

 There are more than 2000+ different SoCs but mobile CPUs show little diversity with ARM's Cortex A53 dominating the market



PERFORMANCE

 Performance difference between mobile CPUs and GPUs is narrow



PROGRAMMABILITY

 Programmability is a major road block for co-processors (e.g. Android GPUs)

MOBILE INFERENCE OPTIMIZATION IS TARGETED FOR THE COMMON DENOMINATOR OF THE FRAGMENTED SOC ECOSYSTEM





Horizontal Integration

Backend Neural Network Libraries in Caffe2 Runtime

NNPACK

(32-BIT FLOATING POINT)

- Optimized convolution implementation using Winograd and FFT
- Best for NN with 3x3, 5x5 or larger convolutions

~~~

#### **QNNPACK/QUANTIZED NNPACK**

#### (8-BIT FIXED POINT)

- Optimized direct convolution implementation
- Best for low-intensity convolutions
- Grouped, depth-wise, dilated convolutions
- Eliminate the overhead of im2col and other memory layout transformation

B



## Horizontal Integration

**QNNPACK Performance Evaluation** 





## Vertical Integrated Systems

8 B

• •

. P Processing Inference for Oculus VR





## Vertical Integrated Systems

ŗ

. .

i~

Performance Acceleration with Co-processors

| <b>DNN Features</b> | MACs | Weights |
|---------------------|------|---------|
| Segmentation        | 1X   | 1.5X    |
| Hand Tracking       | 10X  | 1X      |
| Image Model 1       | 10X  | 2X      |
| Image Model 2       | 100X | 1X      |
| Pose Estimation     | 100X | 4X      |



## Vertical Integrated Systems

Making Inference on DSPs Leads to Consistent Performance

CPU thermal throttling causes sudden **FPS drop** 

The primary reason for using co-processors and accelerators are for **lower power** and **more stable performance** 





## Inference in the Wild

Making "Efficient" Inference in the Wild Requires Developers to Deal with Performance Variability



## Inference in the Wild

#### Does the Performance Variability Follow Certain Statistical Distributions?



[3] Improving Smartphone User Experience by Balancing Performance and Energy with Probabilistic Guarantee. Gaudette et al. HPCA-2016.
[4] Optimizing User Satisfaction of Mobile Workloads Subject to Various Sources of Uncertainties. Gaudette et al. TMC-2018.

# What did we learn?

It is important to consider full-picture and system effects for efficient, practical edge inference designs.

## Data-driven approach to summarize the state of the industry practice:

- Lay of the land for mobile SoCs is extremely heterogeneous.
- Majority of mobile inference run on CPUs, that are the current generation.
- Performance difference between a mobile CPU and GPU/DSP is not 100×.
- Inference performance varies much more widely in the field.
- Co-processors and accelerators are used for energy efficiency and stable performance; speedup is often secondary.



# Thank you

facebook