# New Deep Learning Techniques for Embedded Systems

Tom Michiels, System Architect, Embedded Vision Processors

May 23, 2018
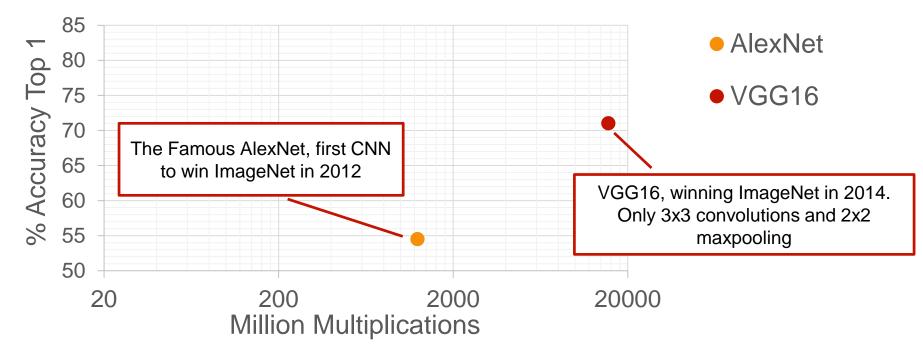
- Evolution towards more efficient and compact CNN graphs

- Memory footprint of weights and feature maps

- External memory bandwidth and impact on power

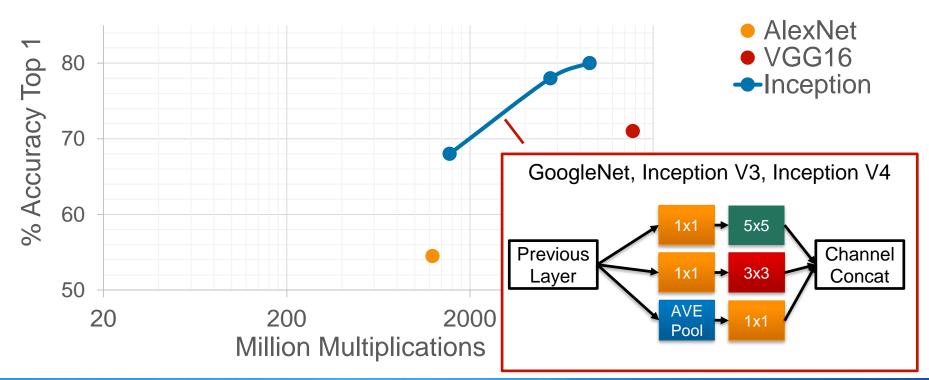- How graph structure can impact bandwidth
  - Tiling and merging layers

# Evolution towards Compact Convolution Networks

*ImageNet ILSVRC-2012 Classification Accuracy vs Compute Requirements*



Legend:
- AlexNet
- VGG16

The Famous AlexNet, first CNN to win ImageNet in 2012

VGG16, winning ImageNet in 2014. Only 3x3 convolutions and 2x2 maxpooling

Y-axis: % Accuracy Top 1 (50, 55, 60, 65, 70, 75, 80, 85)
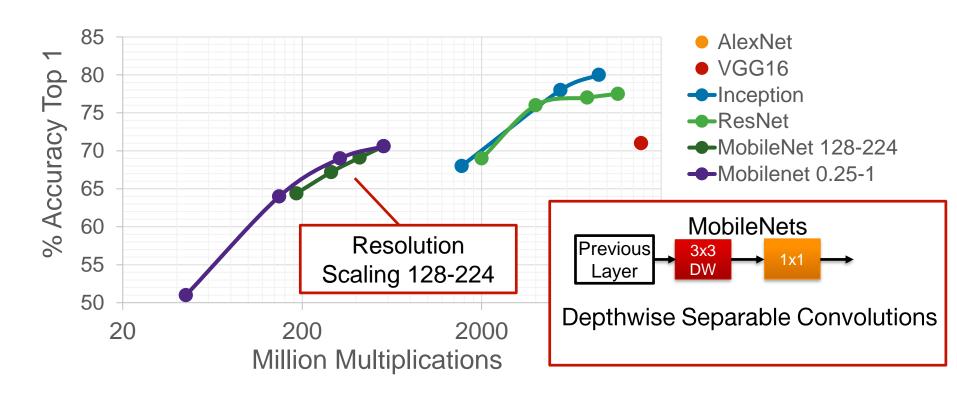X-axis: Million Multiplications (20, 200, 2000, 20000)

*ImageNet ILSVRC-2012 Classification Accuracy vs Compute Requirements*



GoogleNet, Inception V3, Inception V4

*ImageNet ILSVRC-2012 Classification Accuracy vs Compute Requirements*



Legend:
- AlexNet
- VGG16
- Inception
- ResNet

ResNet 18, 50, 101, 152

Previous Layer → 1x1 → 3x3 → 1x1 → Channel Addition

X-axis: Million Multiplications (20, 200, 2000, 20000)

Y-axis: % Accuracy Top 1 (50–85)

Resolution
Scaling 128-224

MobileNets

Previous Layer → 3x3 DW → 1x1 →

Depthwise Separable Convolutions

Legend:
- AlexNet
- VGG16
- Inception
- ResNet
- MobileNet 128-224
- Mobilenet 0.25-1

# Evolution towards Compact Convolution Networks

MultiScale DenseNet

*ImageNet ILSVRC-2012 Classification*



MobileNet V2

Legend:
- AlexNet
- VGG16
- Inception
- ResNet
- MobileNet 128-224
- Mobilenet 0.25-1
- DenseNet
- MultiScale DensNet
- Mobilenet V2

Axes:
- Y-axis: % Accuracy Top 1 (50 to 85)
- X-axis: Million Multiplications (20 to 2000)

# Model Sizes in Million Weights

Power Cost of Computation and local storage

~1 pJoule per MAC

Power of fetching weights and spilling intermediate feature maps off-chip
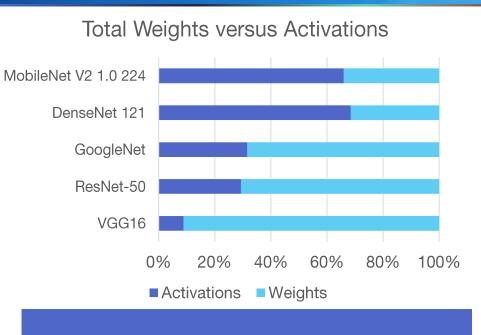
~80 pJoules per byte

Power of computational cost and local storage can relatively easy be computed from the CNN graph structure.

**How is external memory bandwidth related to the graph size and structure?**

# Size of Feature Maps & Coefficients

## Total Weights versus Activations



- MobileNet V2 1.0 224
- DenseNet 121
- GoogleNet
- ResNet-50
- VGG16

0%  20%  40%  60%  80%  100%

■ Activations  ■ Weights

**How is DRAM bandwidth related to the size of feature maps and the number of weights?**

- For more compact graphs, the relative size of the feature maps increases.

- This is even more true when graphs run on large frame resolutions
  - Detection
  - Scene Segmentation
  - Style Transfer

- Weights can be compressed more than feature maps
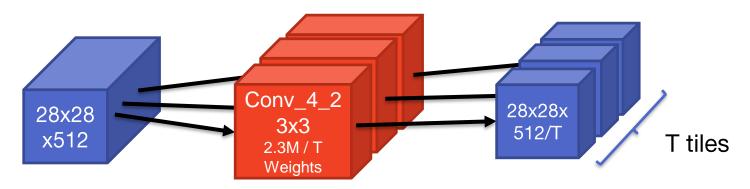  - Pruning
  - Offline Compression

*Example Conv4_2 of VGG16*

```
DataType F_in[512][28][28];
DataType F_out[512][28][28];
DataType W[512][512][3][3];

For (m=0; m<512; m++) {
  for (x=0; x<W; x++) {
    for (y=0; y<H; y++) {
      F_out[m][x][y] = Bias[m];
      for (i=0; i<I; i++) {
        for (j=0; j<J; j++){
          for (k=0; k<512; k++) {
            F_out[m][x][y] += F_in[k][x+i][y+j] * W[m][k][i][j];
          }
        }
      }
    }
  }
}
```

400K Values

400K Values

2.4M Values

- On a memory constrained (e.g. 256KB) architecture, none of these will completely fit in the local memory

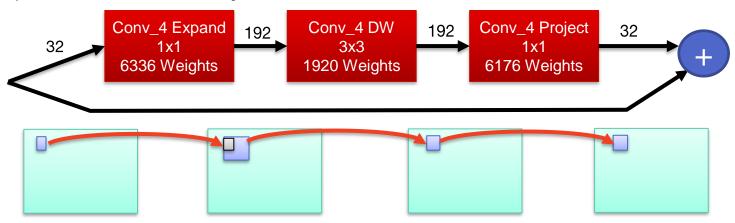- How would we tile it to fit in local memory?

*Example Conv4_2 of VGG16*



We can partition the channels of the activation of the convolution in T parts and re-read the feature maps of the previous layer T times.

**For some layers/graphs the DRAM bandwidth can be much higher than the sum of the weights and the feature maps**

*Example MobileNet V2 Layers*



If the weights of multiple consecutive layers fit the local memory, we can merge the convolution layers, by in the X/Y domain and avoid spilling intermediate feature maps

**For some layers/graphs the DRAM bandwidth can be much lower than the sum of weights and the feature maps**

# How Much Bandwidth for Processing a Layer?

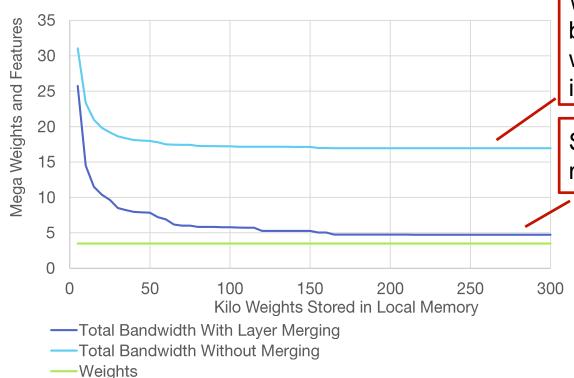Small           Size of Weights of Layers           Large

If the weights of multiple layers can be stored in local memory, spilling of intermediate feature maps can be avoided by merging layers

If the weights of a single layer do not fit, feature maps need to be read multiple times from DRAM

# Bandwidth Example, MobileNet

*Example MobileNet V2 1.0 224*



Without merging of layers, bandwidth decreases until the weight of every single layer fits in local memory
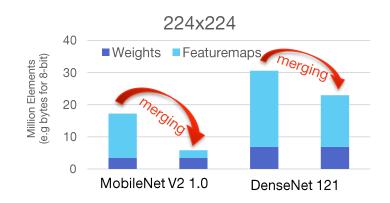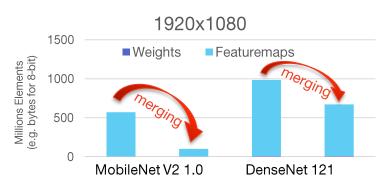
Significant bandwidth gain from merging layers already

Weight compression does not only reduce bandwidth of loading weights, but also that of spilling feature maps

224x224

MobileNet V2 1.0 — DenseNet 121

1920x1080
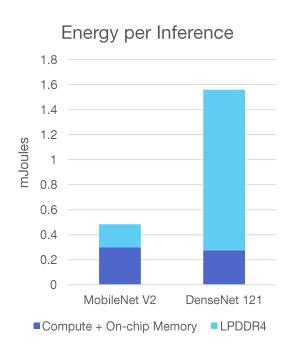
MobileNet V2 1.0 — DenseNet 121

- The opportunities for layer merging depend on the network structure.
  - Densely connected layers are harder to merge because the same feature map is consumed by multiple layers

- Real applications will run on larger frames than 224x224
  - Scene Segmentation,
  - Detection  (Yolo, SSD)

- Feature maps can be compressed, but not as much as coefficients.

## Energy per Inference



- Differences in bandwidth can have big impact on the total energy per inference.

- The graph on the left assumes:
  - 1 pJoules per MAC for the compute and on-chip memory
  - 10 pJoules per bit to read from DRAM
  - 8-bit accuracy
  - Only counting feature maps

- Other architectural choices have impact on the bandwidth.  Example: Feature map compression

# Takeaways

- DRAM bandwidth can be a dominant part in the power consumption of a CNN.

- The bandwidth is determined by weights and feature maps
  - Graphs are often optimized to have small amount of weights, but the bandwidth of feature maps often dominates.

- Bandwidth of feature maps can be reduced significantly by merging convolution layers.
  - This requires an architecture that allows merged execution of convolutions.
  - Storing weights compressed in local memory not only reduces the bandwidth for fetching the weights, but also for fetching and storing feature maps

- Today's networks are optimized for weight-size and macs.
  Optimizing for bandwidth, will likely lead to new network structures.

- Visit **www.synopsys.com/EV**

- Visit the Synopsys booth to check out demos:

  - *Accelerating Android Neural Network Performance with DesignWare EV6x*

  - *Real-Time Object Classification & Tracking with DesignWare EV6x*

  - *AI, 3D Imaging & SLAM on-a-Single Chip for Embedded Markets (by Inuitive)*

  - *Face Recognition for Driver Monitoring System (by PathPartner)*

- Embedded Vision Alliance Website: <u>Software Frameworks and Toolsets for Deep Learning-based Vision Processing</u>

# Thank You

Tom Michiels, System Architect, Embedded Vision Processors