# Aggressive Compression of MobileNets Using Hybrid Ternary Layers
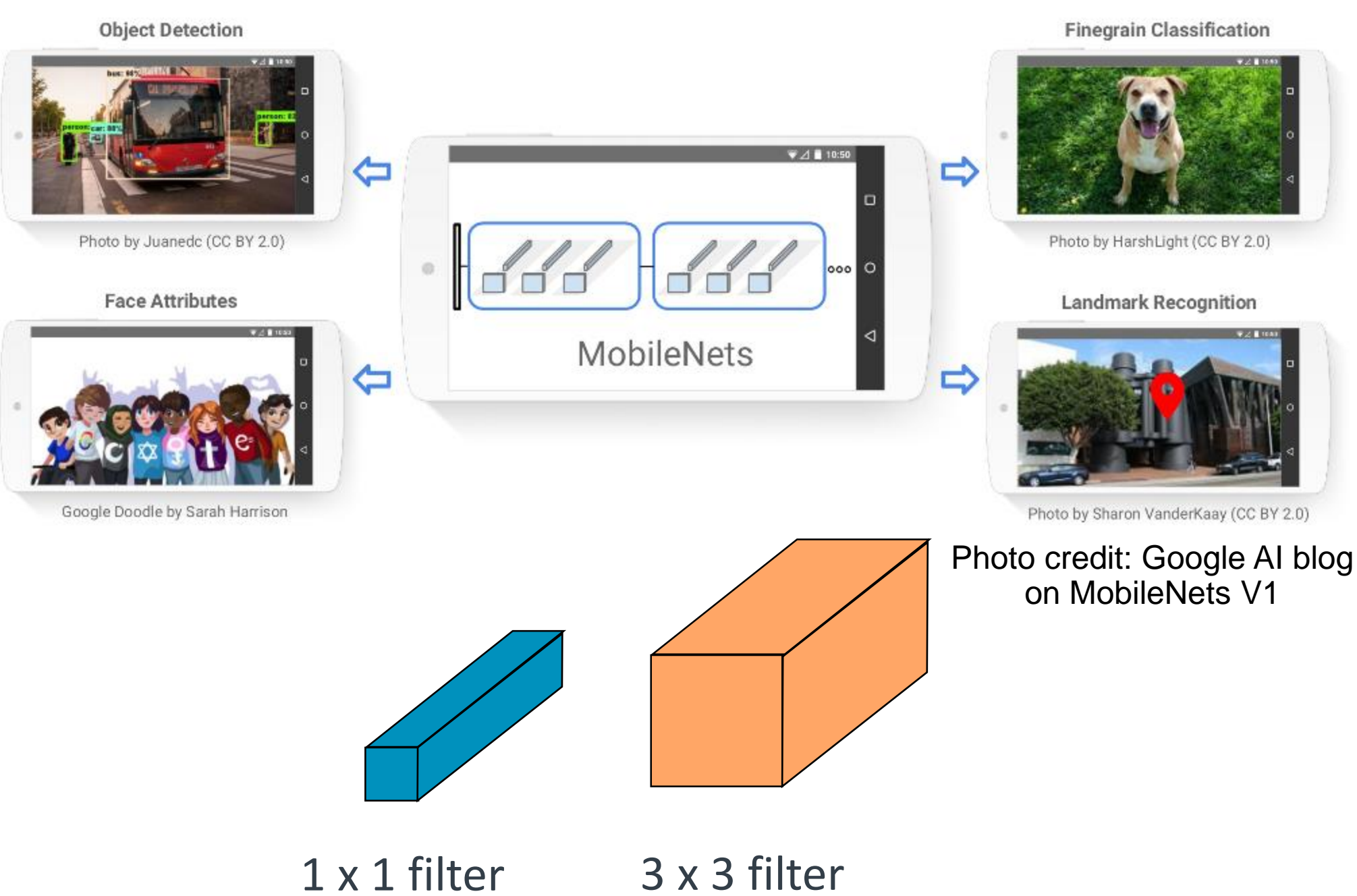
Dibakar Gope, Jesse Beu, Urmish Thakker, and Matthew Mattina

**Arm ML Research Lab**
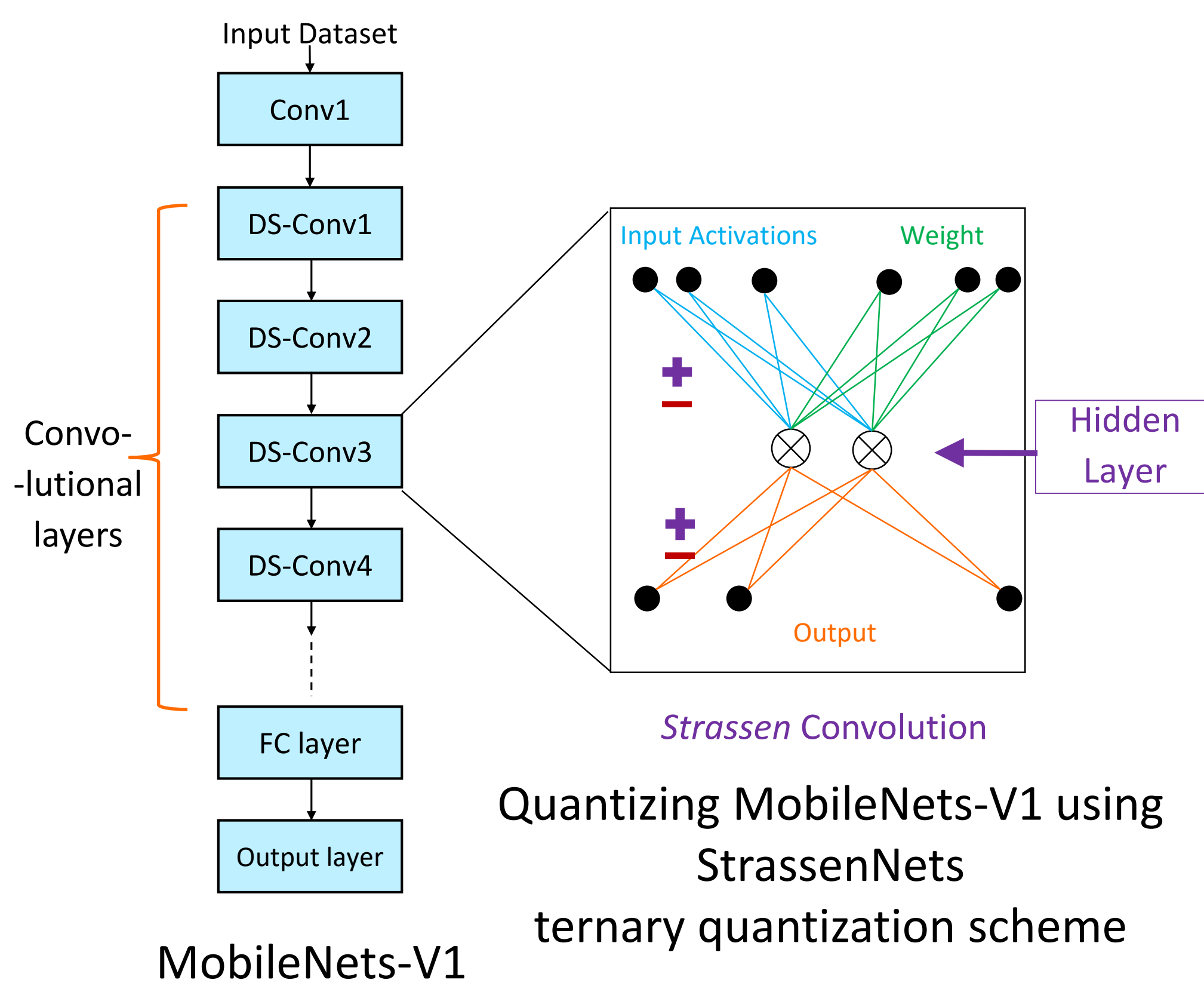
TINY ML

arm Research

## Challenge

- MobileNets [1] family of CV networks are increasingly deployed at mobile/edge devices
- Quantizing MobileNets to ternary weights (2-bit) is necessary to realize siginificant energy savings and runtime speedups
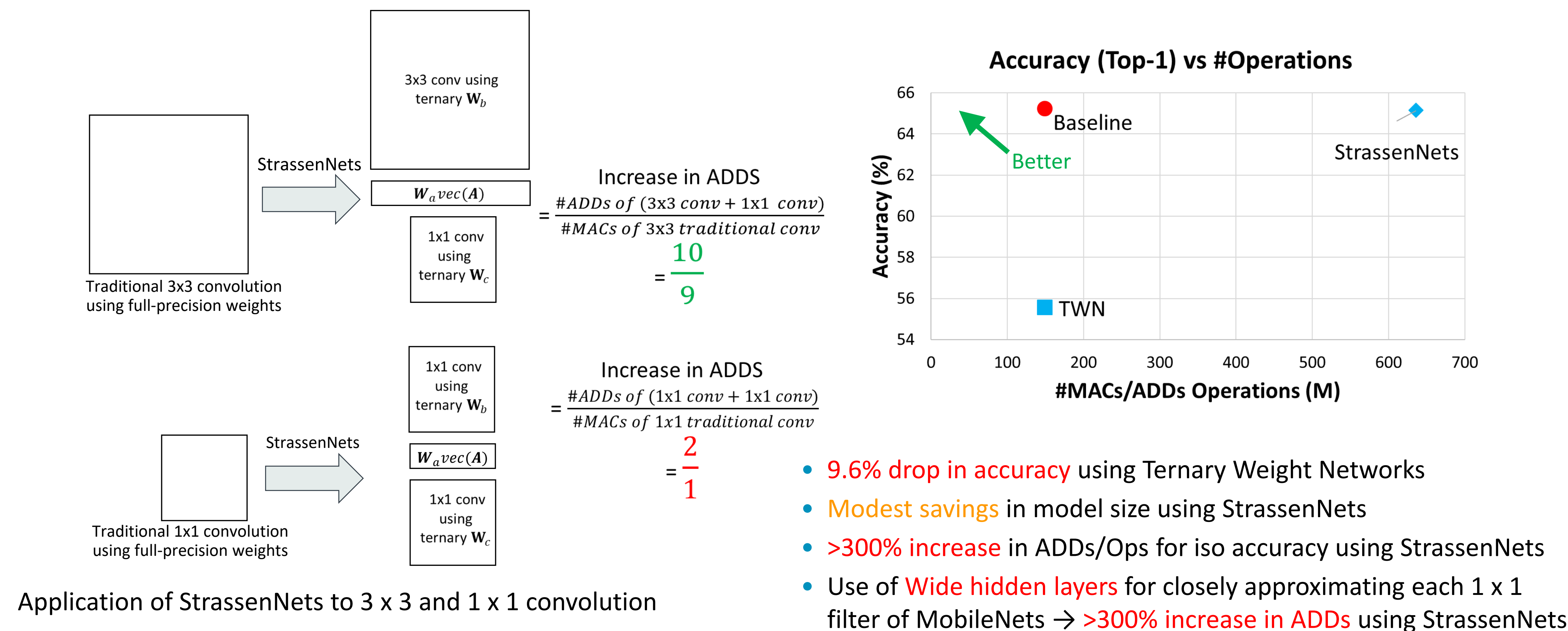
Object Detection | Finegrain Classification
Face Attributes | Landmark Recognition

MobileNets

Photo credit: Google AI blog on MobileNets V1

1 x 1 filter — Very compact
3 x 3 filter — Over parametrized

- MobileNets V1 – 13 depthwise separable (DS) convolutional layers
- Model complexity dominated by compact 1 x 1 filters

## Prior Solutions

- Ternary weight networks (TWN) [2]
  (-) Drops accuracy
- StrassenNets [3]
  (+) 99% reduction in MULs for 3 x 3 filters
  (+) mostly ternary weights, preserve accuracy
  (-) Never looked into DS (1 x 1) layers
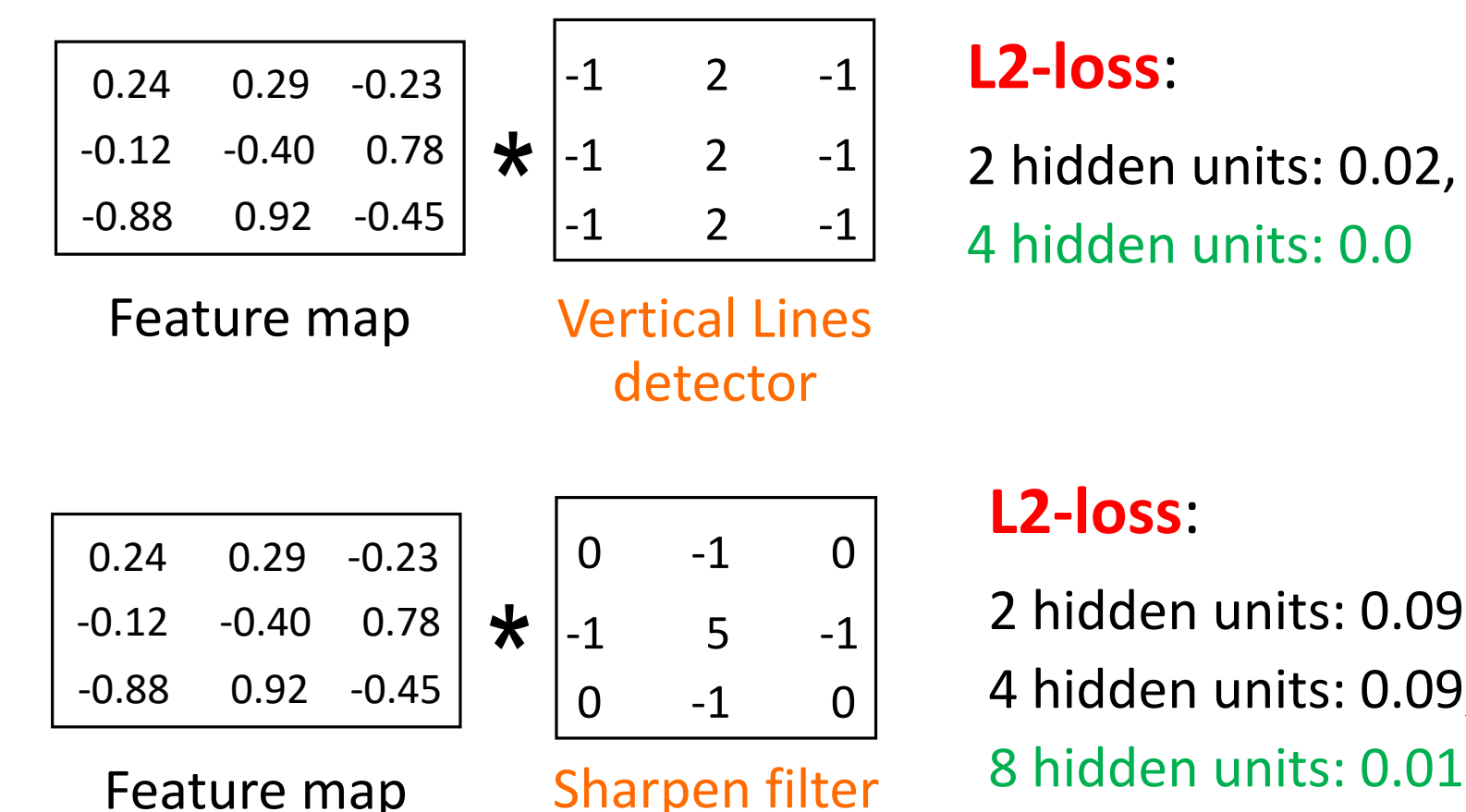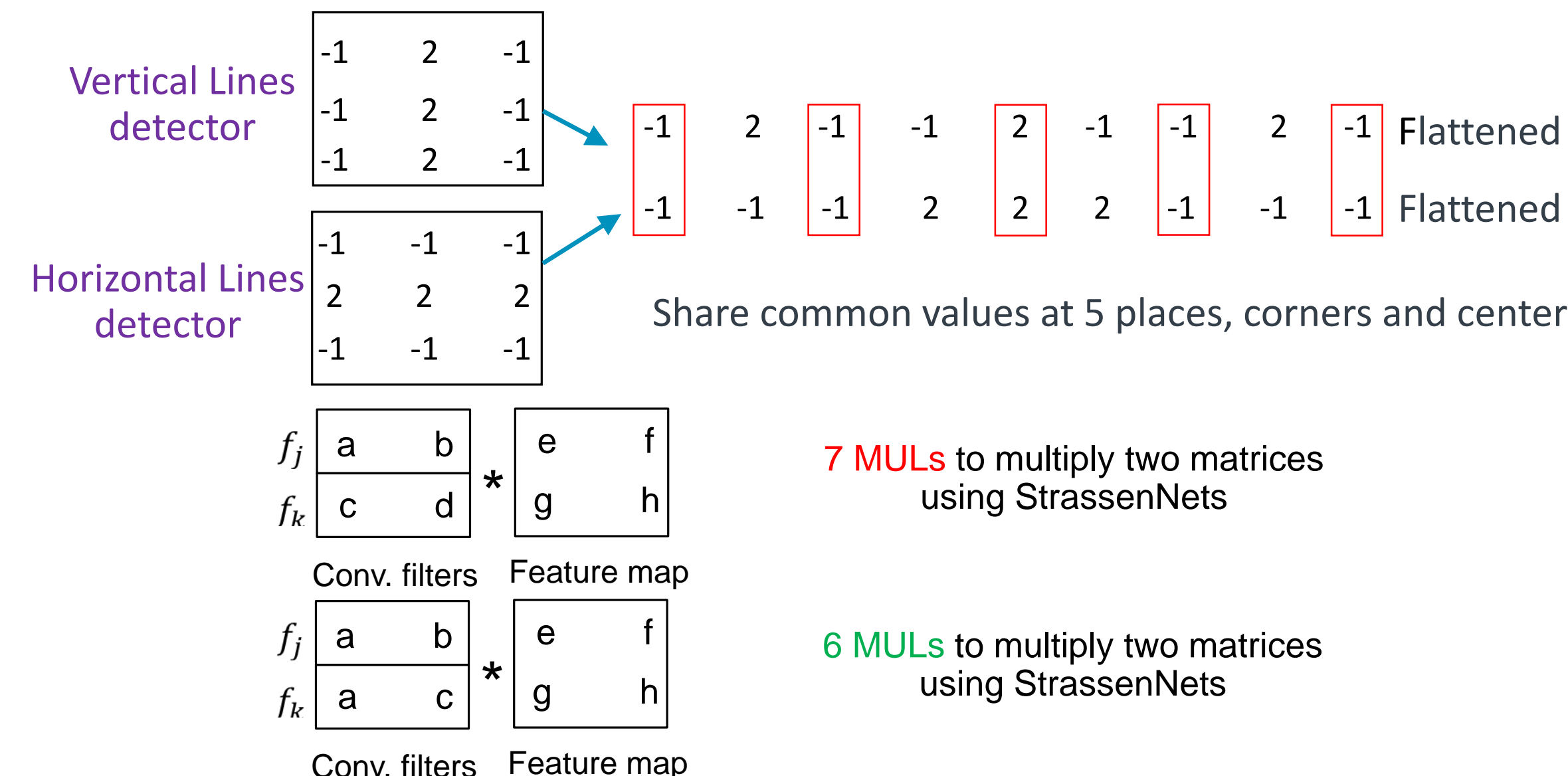- Prior solutions come with their own advantages and limitations

Input Dataset
Conv1
DS-Conv1
DS-Conv2
DS-Conv3
DS-Conv4
FC layer
Output layer

MobileNets-V1

Convolutional layers

Input Activations | Weight
+ / − | Hidden Layer
Output
*Strassen* Convolution

Quantizing MobileNets-V1 using StrassenNets ternary quantization scheme

## Observations with Prior Solutions

3x3 conv using ternary $W_b$

Traditional 3x3 convolution using full-precision weights

StrassenNets → $W_a vec(A)$ → 1x1 conv using ternary $W_c$

Increase in ADDS
$$\frac{\#ADDs \ of \ (3x3 \ conv + 1x1 \ conv)}{\#MACs \ of \ 3x3 \ traditional \ conv} = \frac{10}{9}$$

Traditional 1x1 convolution using full-precision weights

StrassenNets → $W_a vec(A)$ → 1x1 conv using ternary $W_c$

1x1 conv using ternary $W_b$

Increase in ADDS
$$\frac{\#ADDs \ of \ (1x1 \ conv + 1x1 \ conv)}{\#MACs \ of \ 1x1 \ traditional \ conv} = \frac{2}{1}$$

Application of StrassenNets to 3 x 3 and 1 x 1 convolution

### Accuracy (Top-1) vs #Operations

Baseline
Better
StrassenNets
TWN

#MACs/ADDs Operations (M)

- 9.6% drop in accuracy using Ternary Weight Networks
- Modest savings in model size using StrassenNets
- >300% increase in ADDs/Ops for iso accuracy using StrassenNets
- Use of Wide hidden layers for closely approximating each 1 x 1 filter of MobileNets → >300% increase in ADDs using StrassenNets

## Different filters respond differently to ternary quantization

### Different sensitivity of individual filters to StrassenNets

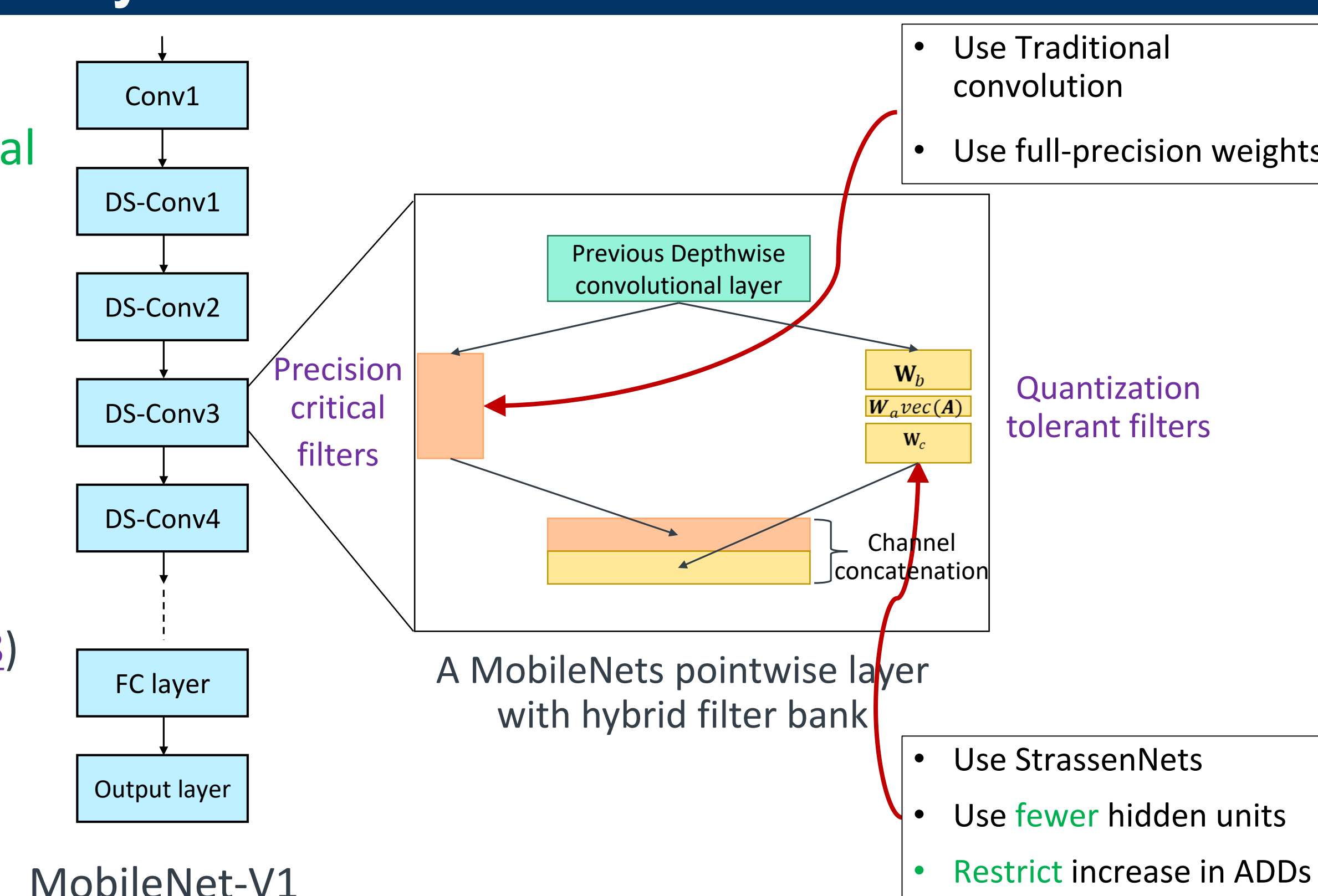| 0.24 | 0.29 | -0.23 |
| -0.12 | -0.40 | 0.78 |
| -0.88 | 0.92 | -0.45 |

* 

| -1 | 2 | -1 |
| -1 | 2 | -1 |
| -1 | 2 | -1 |

Feature map | Vertical Lines detector

**L2-loss**:
2 hidden units: 0.02,
4 hidden units: 0.0

| 0.24 | 0.29 | -0.23 |
| -0.12 | -0.40 | 0.78 |
| -0.88 | 0.92 | -0.45 |

* 

| 0 | -1 | 0 |
| -1 | 5 | -1 |
| 0 | -1 | 0 |

Feature map | Sharpen filter

**L2-loss**:
2 hidden units: 0.09,
4 hidden units: 0.09,
8 hidden units: 0.01

Not all filters do require wide hidden layers to be approximated well using StrassenNets

### Different sensitivity of group of filters to StrassenNets

Vertical Lines detector
| -1 | 2 | -1 |
| -1 | 2 | -1 |
| -1 | 2 | -1 |

Horizontal Lines detector
| -1 | -1 | -1 |
| 2 | 2 | 2 |
| -1 | -1 | -1 |

-1  2  -1  -1  2  -1  2  -1  -1  2  -1  Flattened
-1  -1  -1  2  2  2  -1  -1  -1  Flattened

Share common values at 5 places, corners and center

$f_j$ | a | b    $f_j$ | e | f
$f_k$ | c | d    $f_k$ | g | h

7 MULs to multiply two matrices using StrassenNets

$f_j$ | a | b    $f_j$ | e | f
$f_k$ | a | c    $f_k$ | g | h

6 MULs to multiply two matrices using StrassenNets

Conv. filters | Feature map
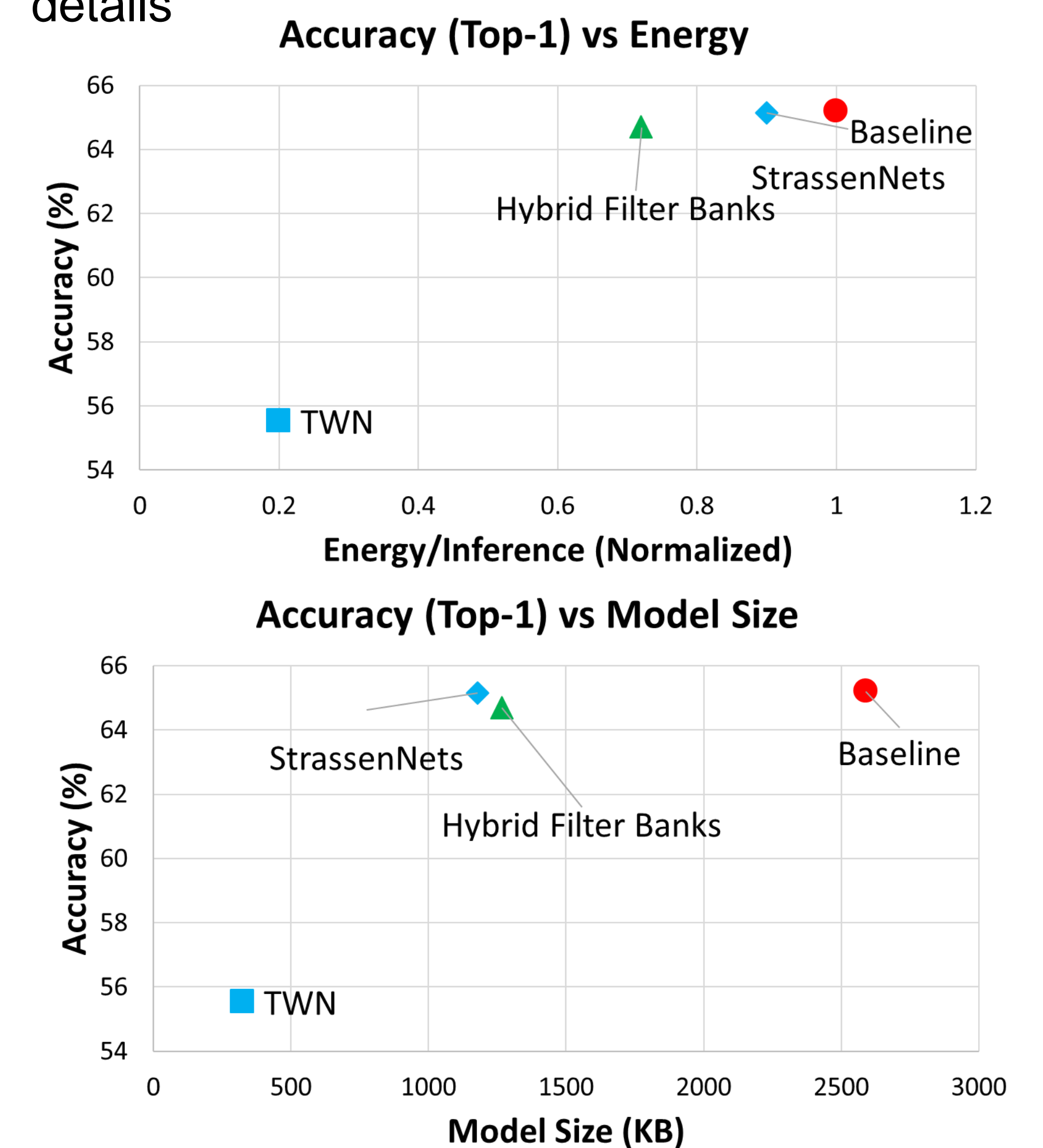
## Per-Layer Hybrid Filter Banks

Exploit the difference in sensitivity of individual and groups of filters to ternary quantization

- Bank similar value structure filters together
- Share hidden units of StrassenNets
- Use fewer hidden units → fewer ADDs/Ops to approximate a major portion of filters at each layer
- See our paper (https://arxiv.org/abs/1911.01028) for Mathematical proof, details

Conv1
DS-Conv1
DS-Conv2
DS-Conv3
DS-Conv4
FC layer
Output layer

MobileNet-V1

Precision critical filters

Previous Depthwise convolutional layer
$W_b$
$W_a vec(A)$
$W_c$
Quantization tolerant filters
Channel concatenation

- Use Traditional convolution
- Use full-precision weights

- Use StrassenNets
- Use fewer hidden units
- Restrict increase in ADDs

A MobileNets pointwise layer with hybrid filter bank

## Evaluation Results

- Dataset: ImageNet, Network: MobileNet-V1 (width multiplier of 0.5)
- 47% reduction in MULs, only 48% reduction in ADDs, when compared to >300%
- 51% reduction in MobileNets-V1 model size,
- 28% reduction in energy/inference
- No degradation in inference throughput on an area-equivalent ML accelerator comprising both MAC and adder units
- 0.27% loss in top-1 accuracy
- Hybrid filter banks is effective in compressing ResNet architecture comprising 3x3 convolutional filters also; see our paper for details

### Accuracy (Top-1) vs Energy

Baseline
StrassenNets
Hybrid Filter Banks
TWN

Energy/Inference (Normalized)

### Accuracy (Top-1) vs Model Size

StrassenNets
Hybrid Filter Banks
Baseline
TWN

Model Size (KB)

Top-1 accuracy, energy/inference, and model size of hybrid filter banks and improvement over state-of-the-art ternary quantization schemes

## Read Our Paper for Details

Gope et al., "Ternary MobileNets via Per-Layer Hybrid Filter Banks", 2019

arXiv link: https://arxiv.org/abs/1911.01028

## References

[1] Howard et al., "MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications", 2017
[2] Li et al., "Ternary weight networks," NeurIPS 2016
[3] Tschannen et al., "StrassenNets: Deep Learning with a Multiplication Budget", ICML 2018