# ComputerVision

## Week6

# From ResNet to Lightweight CNNs: MobileNet and EfficientNet

- **Why Do We Need Lightweight Models?**
  - ResNet and deeper models improve accuracy but are **too large** for real-time/mobile use.
  - **Mobile apps, drones, AR, and robotics** require fast and lightweight models.
  - Need a **trade-off between accuracy, speed, and size**.



ResNet50 Model Architecture

| Model | 2D-CNN Params | 3D-CNN Params |
|---|---|---|
| VGG-16 | 134.7 M | 179.1 M |
| ResNet-18 | 11.4 M | 33.3 M |
| ResNet-34 | 21.5 M | 63.6 M |
| ResNet-50 | 23.9 M | 46.4 M |
| ResNet-101 | 42.8 M | 85.5 M |
| ResNet-152 | 58.5 M | 117.6 M |

**"Low memory, real-time requirement"**

# From ResNet to Lightweight CNNs: MobileNet and EfficientNet

- **Three Directions in CNN Efficiency**
  - We can reduce the computational cost of CNNs by adjusting three key design axes.

| Axis | What it Means | Effect on Computation |
|---|---|---|
| **Depth** | Number of layers | Linear increase in FLOPs |
| **Width** | Number of channels per layer | Quadratic increase in FLOPs |
| **Resolution** | Size of input images | Quadratic increase in FLOPs |

  - Efficient design must **balance all three** rather than scaling just one.

  - This is the idea behind **EfficientNet's compound scaling**.

# From ResNet to Lightweight CNNs: MobileNet and EfficientNet

- ▪ **What is a FLOP?**

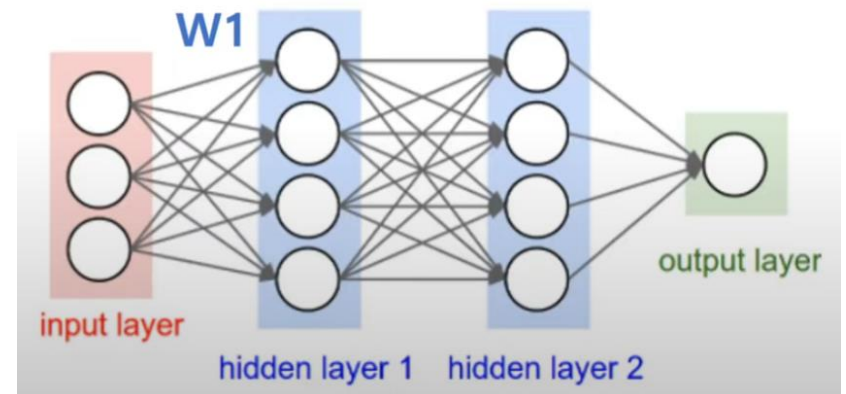  - • **Understanding FLOPs – What Do We Actually Count?**

    - ○ **FLOP = Floating Point Operation**

    - ○ A single **multiplication** or **addition** is counted as 1 FLOP

    - ○ In deep learning, FLOPs are used to estimate the **computational cost** of models

  - • **Example 1: Dot Product**

    - ○ $Y = w[0] \cdot x[0] + w[1] \cdot x[1] + w[2] \cdot x[2]$

      - ✓ **3** multiplications and **2** additions → **Total: 5 FLOPs** → **In general: 2n−1 FLOPs for n elements**

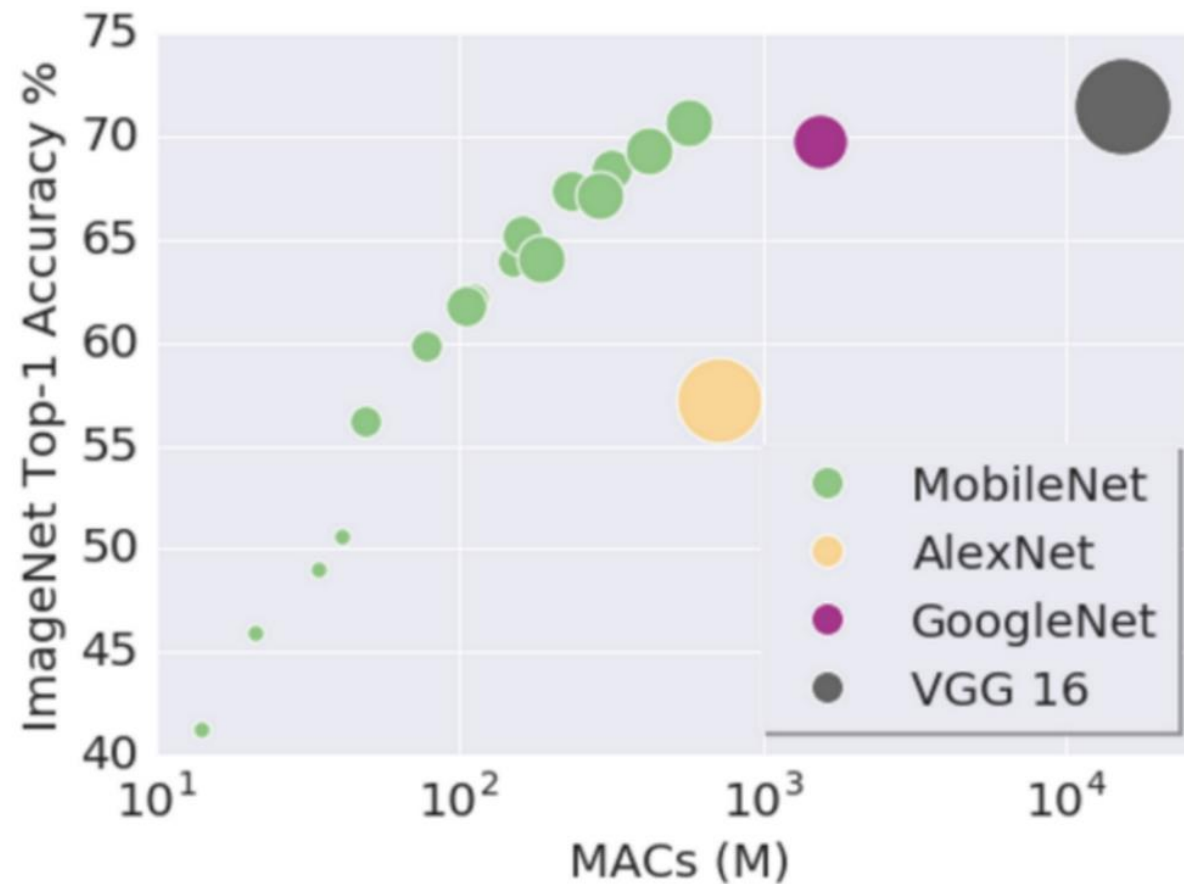  - • **Example 2: Dense (Fully Connected) Layer**

    - ○ **Input vector:** 3-dimensional / **Output vector:** 4-dimensional
      → **Weight matrix:** $W1 \in \mathbb{R}^{4 \times 3}$

    - ○ **4 dot products performed** → 5×4=20 FLOPs

      - ✓ **12 multiplications** (4 rows × 3 elements)

      - ✓ **8 additions** (each dot product needs n-1 additions)



W1

input layer

hidden layer 1    hidden layer 2

output layer

# What is MobileNet?

- ## Introducing MobileNet – Efficient CNNs for Mobile Vision

  - Proposed by Google in 2017

  - Designed for **mobile and embedded vision applications**

  - **Key ideas**
    - **Depthwise Separable Convolution**
    - **Two hyperparameters** for flexible scaling

  - Used in real-world tasks
    - object detection, face recognition, geo-localization

# Standard Convolution vs. Depthwise Separable Convolution

- **Rethinking Convolution: Lighter and Faster**

  - **Standard convolution**

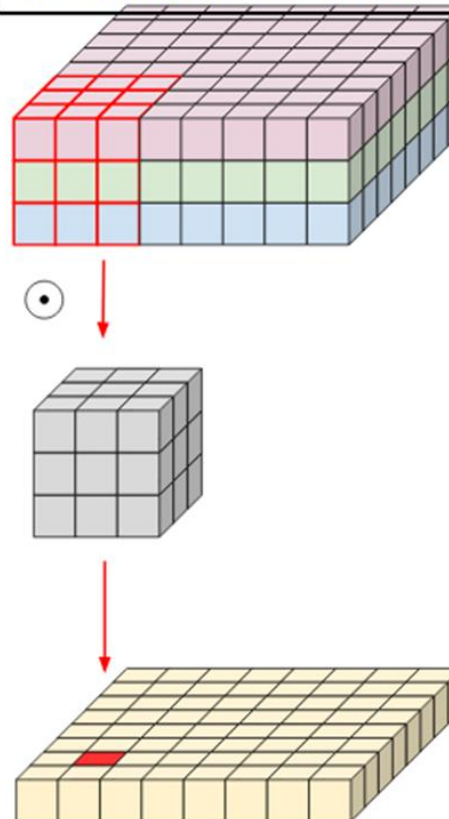    o Applies multiple filters across all channels simultaneously

  - **Depthwise separable convolution**

    o Depthwise

    ✓ Applies one filter per input channel
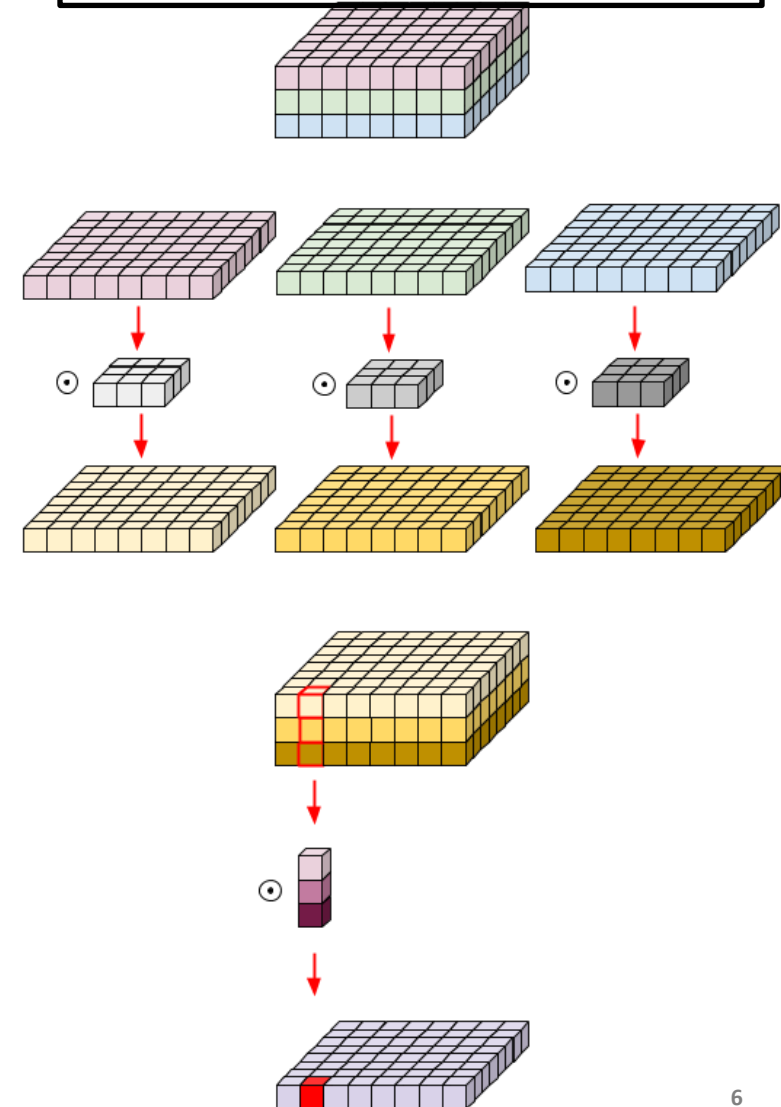
    o Pointwise

    ✓ 1×1 convolution to combine outputs

  - Result

    o **~8–9x less computation**
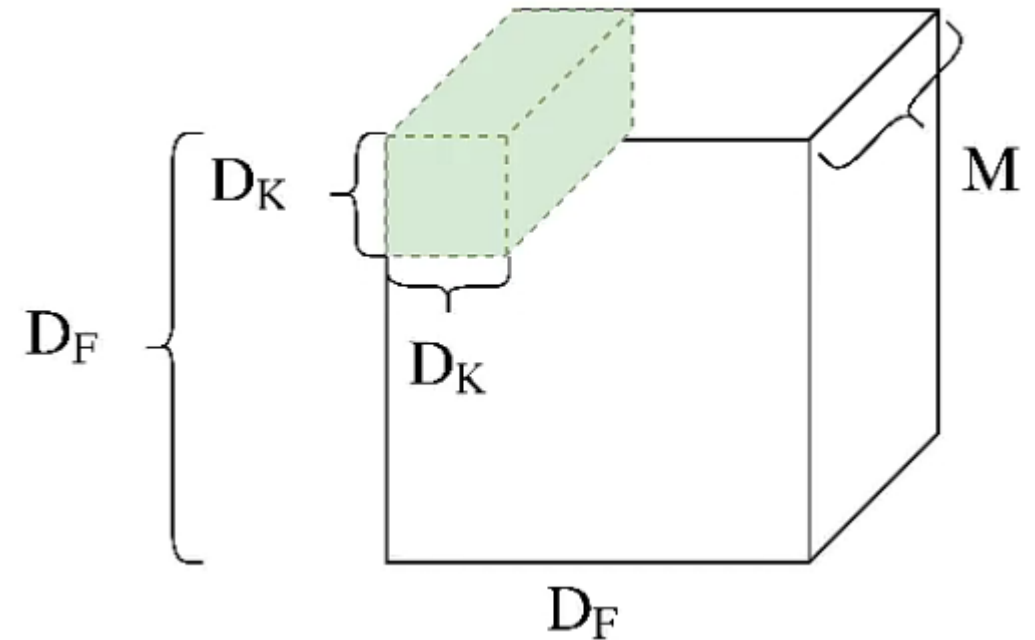


Standard Convolution

Depthwise Convolution

# How Efficient is Depthwise Separable Conv?

- **FLOPs Comparison: Standard vs. Depthwise Separable**
  - **How to Compute FLOPs in Standard Convolution**
    - **Key Definitions**
      - ✓ $D_K$: Kernel size (e.g., 3 → 3×3)
      - ✓ $M$: Number of input channels
      - ✓ $N$: Number of output channels (i.e., number of filters)
      - ✓ $D_F$: Input feature map spatial dimension
      - ✓ $D_G$: Output feature map spatial dimension
        - ➤ Usually, $D_G \approx D_F$ (same padding, stride 1)

# How Efficient is Depthwise Separable Conv?

- **FLOPs Comparison: Standard vs. Depthwise Separable**
  - **How to Compute FLOPs in Standard Convolution**
    - **Step-by-Step FLOP Count**
      - ✓ **Step 1. Single Location, Single Filter**
        - ➢ Each filter performs: $D_K^2 \times M$ multiplications

      - ✓ **Step 2. All Spatial Locations in One Filter**
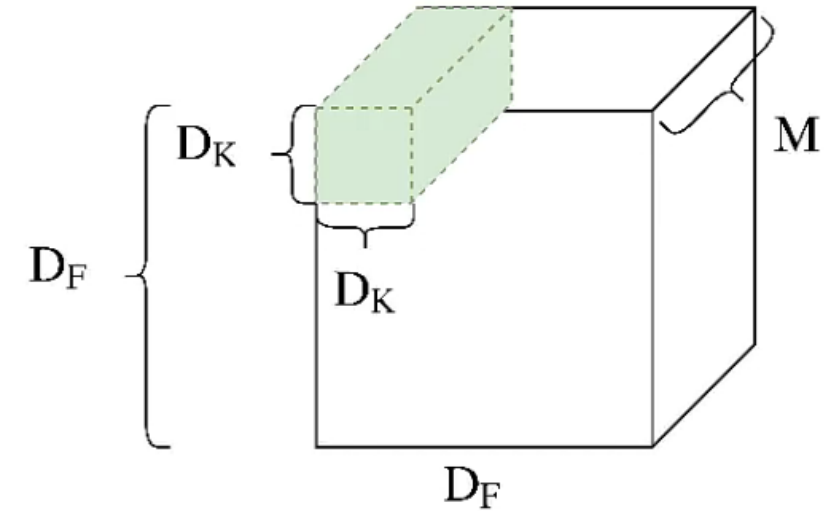        - ➢ Spatial positions: $D_G^2$ (assuming square output)
          → Typically $D_G = D_F$ for stride 1 & padding
        - ➢ FLOPs per filter: $D_G^2 \times D_K^2 \times M$

      - ✓ **Step 3. All N Filters (Output Channels)**
        - ➢ Multiply by $N$

      - ✓ Total FLOPs = $D_G^2 \times D_K^2 \times M \times N$



$$\text{Mults once} = D_K^2 \times M$$

$$\text{Mults per Kernel} = D_G^2 \times D_K^2 \times M$$

$$\text{Mults N Kernels} = N \times D_G^2 \times D_K^2 \times M$$

# How Efficient is Depthwise Separable Conv?

- **FLOPs Comparison: Standard vs. Depthwise Separable**
  - What is Depthwise Separable Convolution?
    - **Key Concepts**
      - ✓ Standard convolution combines **spatial filtering + channel mixing** in one operation.

      - ✓ Depthwise separable convolution splits this into two steps
        - ➤ **Step1. Depthwise Convolution**: applies spatial filtering **per input channel**
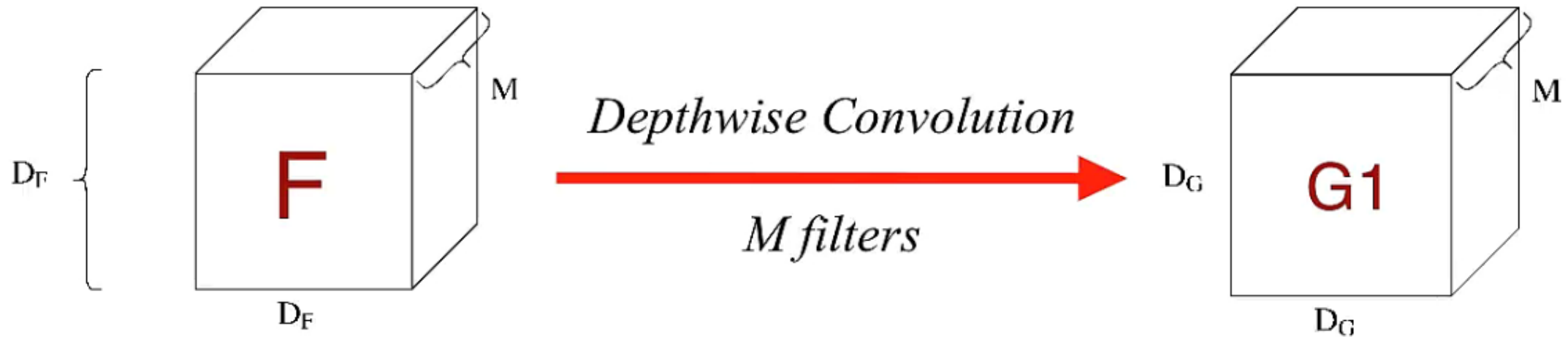        - ➤ **Step 2. Pointwise Convolution**: mixes the output channels using **1×1 conv**

      - → This separation significantly **reduces computation and parameters**.

# How Efficient is Depthwise Separable Conv?

- **FLOPs Comparison: Standard vs. Depthwise Separable**
  - FLOPs Breakdown of Depthwise Separable Convolution
    - **Key Definition**



- $D_K$: kernel size (e.g., 3 for 3×3)
- $D_F$: input spatial resolution
- $M$: number of input channels
- $N$: number of output channels
- Assume $D_G = D_F$ for simplicity

# How Efficient is Depthwise Separable Conv?

## ▪ FLOPs Comparison: Standard vs. Depthwise Separable

- FLOPs Breakdown of Depthwise Separable Convolution
  - Step-by-step FLOP Calculation
    - ✓ **Step 1. Depthwise Convolution (per channel)**



➢ Each of the **M** channels gets a **filter**

➢ For each channel: $D_F^2 \times D_K^2$ operations

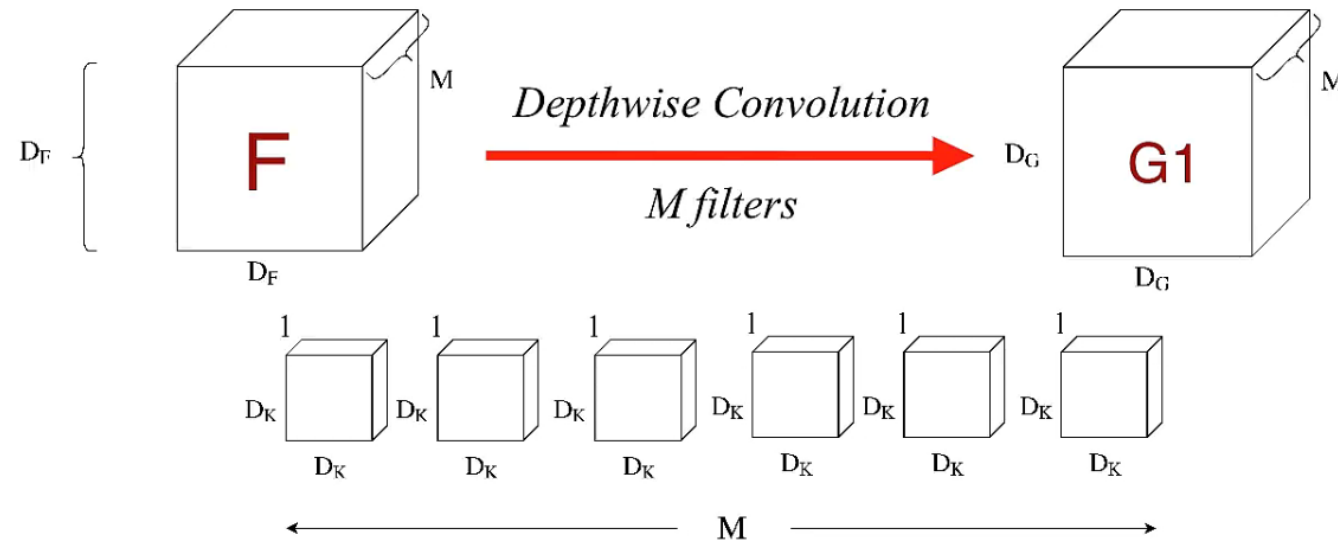→ Total: $D_F^2 \times D_K^2 \times M$
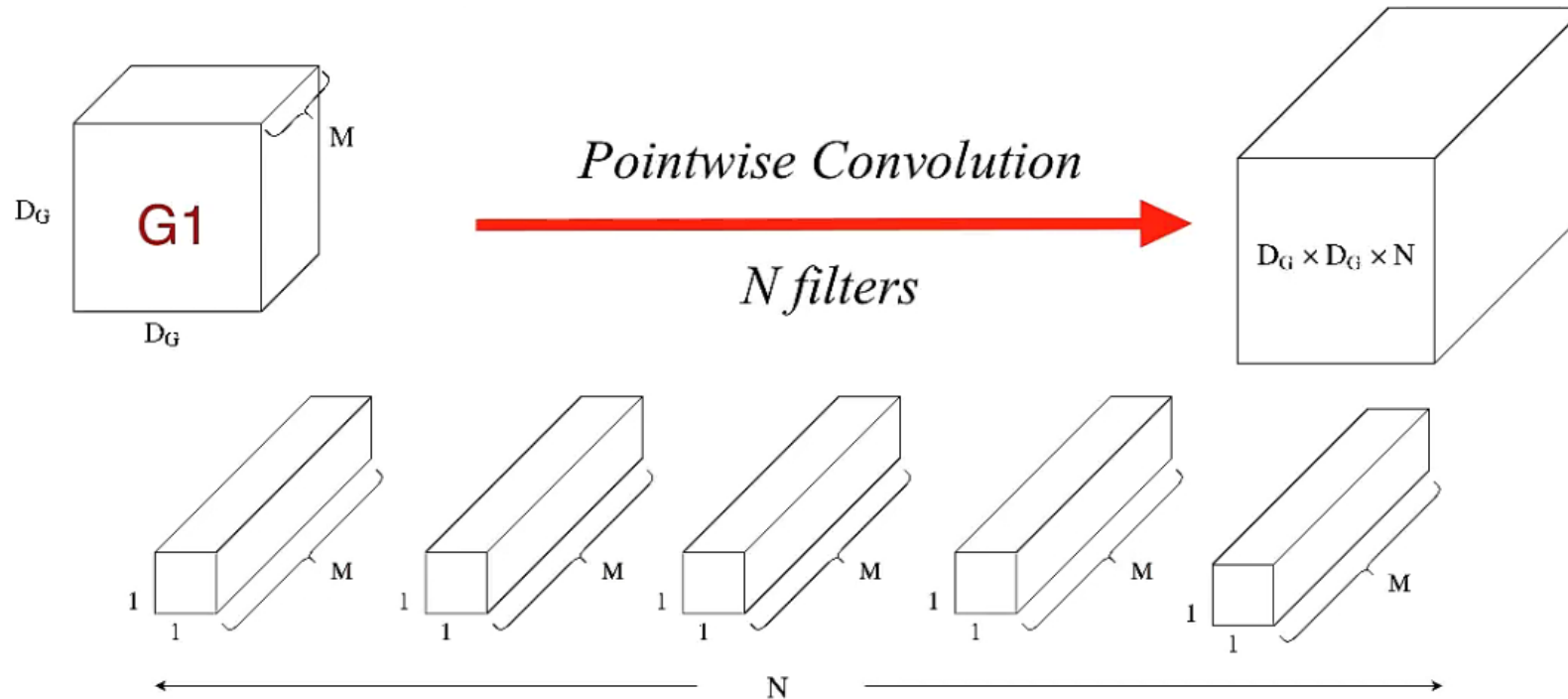
# How Efficient is Depthwise Separable Conv?

- **FLOPs Comparison: Standard vs. Depthwise Separable**
  - FLOPs Breakdown of Depthwise Separable Convolution
    - o Step-by-step FLOP Calculation
      - ✓ **Step 2. Pointwise Convolution (1×1 conv)**



➢ Each of the $D_F \times D_F$ positions applies **N** 1×1 filters across **M** channels

→ Total: $M \times N \times D_F^2$

# How Efficient is Depthwise Separable Conv?

- **FLOPs Comparison: Standard vs. Depthwise Separable**

  - **Total FLOPs**
    - $FLOPs = D_F^2 \times D_K^2 \times M \ (i.e., Step 1) + M \times N \times D_F^2 \ (i.e., Step \ 2)$

  - Compare with **Standard Convolution**
    - $FLOPs_{standard} = D_G^2 \times D_K^2 \times M \times N$

  - **Efficiency Gain**
    - $$\frac{Depthwise \ Separable \ Convolution}{Standard \ Convolution} = \frac{D_F^2 \times D_K^2 \times M + M \times N \times D_F^2}{D_G^2 \times D_K^2 \times M \times N} = \frac{M \times D_F^2 (D_K^2 + N)}{D_G^2 \times D_K^2 \times M \times N}$$

    $$= \frac{M \times D_G^2 (D_K^2 + N)}{D_G^2 \times D_K^2 \times M \times N} = \frac{D_K^2 + N}{D_K^2 \times N} = \frac{1}{N} + \frac{1}{D_K^2}$$

# How Efficient is Depthwise Separable Conv?

- **FLOPs Comparison: Standard vs. Depthwise Separable**

  - **Example – Efficiency Gain**

    - Efficiency Gain: $\dfrac{1}{N} + \dfrac{1}{D_K^2}$

    - For $D_K = 3$ (size of filter), $N = 256$ (number of channels)

      - Relative Cost $= \dfrac{1}{256} + \dfrac{1}{3^2} = \dfrac{1}{256} + \dfrac{1}{9} = \dfrac{265}{2304} = 0.11501 \ldots$

      - The calculation cost is reduced by about 9 times.

# MobileNet Architecture
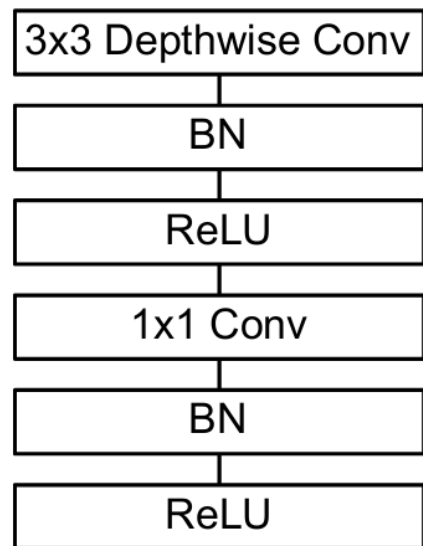
- **Building Blocks of MobileNet**
  - **All layers (except the first) use depthwise separable conv**

  - **Each block**
    - Depthwise conv → BN → ReLU →
      Pointwise conv ($1 \times 1$) → BN → ReLU

  - **Total: 28 layers**



| Type / Stride | Filter Shape | Input Size |
|---|---|---|
| Conv / s2 | $3 \times 3 \times 3 \times 32$ | $224 \times 224 \times 3$ |
| Conv dw / s1 | $3 \times 3 \times 32$ dw | $112 \times 112 \times 32$ |
| Conv / s1 | $1 \times 1 \times 32 \times 64$ | $112 \times 112 \times 32$ |
| Conv dw / s2 | $3 \times 3 \times 64$ dw | $112 \times 112 \times 64$ |
| Conv / s1 | $1 \times 1 \times 64 \times 128$ | $56 \times 56 \times 64$ |
| Conv dw / s1 | $3 \times 3 \times 128$ dw | $56 \times 56 \times 128$ |
| Conv / s1 | $1 \times 1 \times 128 \times 128$ | $56 \times 56 \times 128$ |
| Conv dw / s2 | $3 \times 3 \times 128$ dw | $56 \times 56 \times 128$ |
| Conv / s1 | $1 \times 1 \times 128 \times 256$ | $28 \times 28 \times 128$ |
| Conv dw / s1 | $3 \times 3 \times 256$ dw | $28 \times 28 \times 256$ |
| Conv / s1 | $1 \times 1 \times 256 \times 256$ | $28 \times 28 \times 256$ |
| Conv dw / s2 | $3 \times 3 \times 256$ dw | $28 \times 28 \times 256$ |
| Conv / s1 | $1 \times 1 \times 256 \times 512$ | $14 \times 14 \times 256$ |
| $5\times$ Conv dw / s1 | $3 \times 3 \times 512$ dw | $14 \times 14 \times 512$ |
| Conv / s1 | $1 \times 1 \times 512 \times 512$ | $14 \times 14 \times 512$ |
| Conv dw / s2 | $3 \times 3 \times 512$ dw | $14 \times 14 \times 512$ |
| Conv / s1 | $1 \times 1 \times 512 \times 1024$ | $7 \times 7 \times 512$ |
| Conv dw / s2 | $3 \times 3 \times 1024$ dw | $7 \times 7 \times 1024$ |
| Conv / s1 | $1 \times 1 \times 1024 \times 1024$ | $7 \times 7 \times 1024$ |
| Avg Pool / s1 | Pool $7 \times 7$ | $7 \times 7 \times 1024$ |
| FC / s1 | $1024 \times 1000$ | $1 \times 1 \times 1024$ |
| Softmax / s1 | Classifier | $1 \times 1 \times 1000$ |

# MobileNet Architecture

▪ **Width Multiplier & MobileNet's Efficiency**

- **MobileNet Trade-offs: Width Multiplier**

    o Scaling with Width Multiplier $\alpha$

    ✓Controls number of channels at each layer

    o Input/output channels → multiplied by $\alpha$ $(0 < \alpha \le 1)$

    ✓Reduces FLOPs and parameters by approximately $\alpha^2$

    o Common choices

    ✓$\alpha \in \{1.0, 0.75, 0.5, 0.25\}$

| Width Multiplier | ImageNet Accuracy | Million Mult-Adds | Million Parameters |
|---|---|---|---|
| 1.0 MobileNet-224 | 70.6% | 569 | 4.2 |
| 0.75 MobileNet-224 | 68.4% | 325 | 2.6 |
| 0.5 MobileNet-224 | 63.7% | 149 | 1.3 |
| 0.25 MobileNet-224 | 50.6% | 41 | 0.5 |

    o **Key Observations**

    ✓As **width multiplier** decreases

    ➢**(1)** Accuracy drops / **(2)** FLOPs and parameters drop **significantly**

    ✓Enables scaling model size for devices with limited resources