# Computer Vision

## Naeemullah Khan

naeemullah.khan@kaust.edu.sa

KAUST Academy
King Abdullah University of Science and Technology

February 17, 2025

# Learning Outcomes

▶ Understand the basic concepts of Computer Vision and its real-world applications.

▶ Describe how images are represented and processed in a computer.

▶ Explain the fundamental building blocks of Convolutional Neural Networks (CNNs).

▶ Differentiate between popular CNN architectures (AlexNet, VGG, InceptionNet, ResNet, EfficientNet, and MobileNet) and their key innovations.
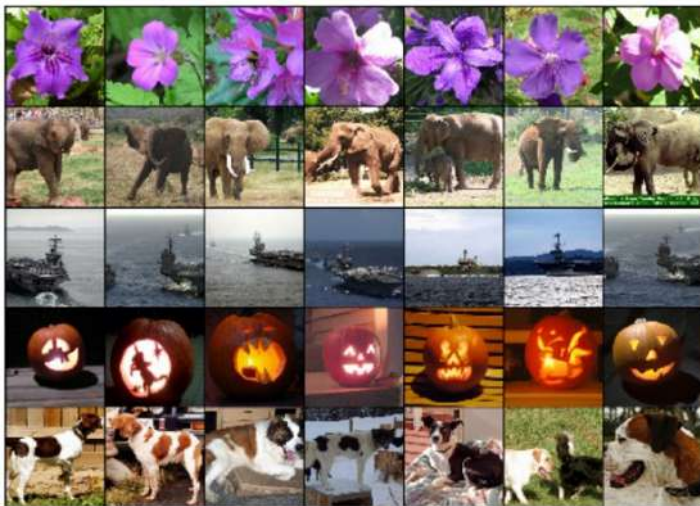
Building artificial systems that process, perceive, and reason about visual data
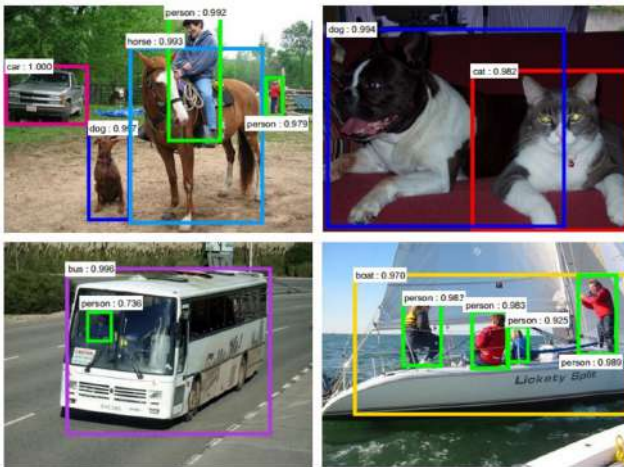
# Computer Vision is Everywhere

# Image Classification

## Image Retrieval

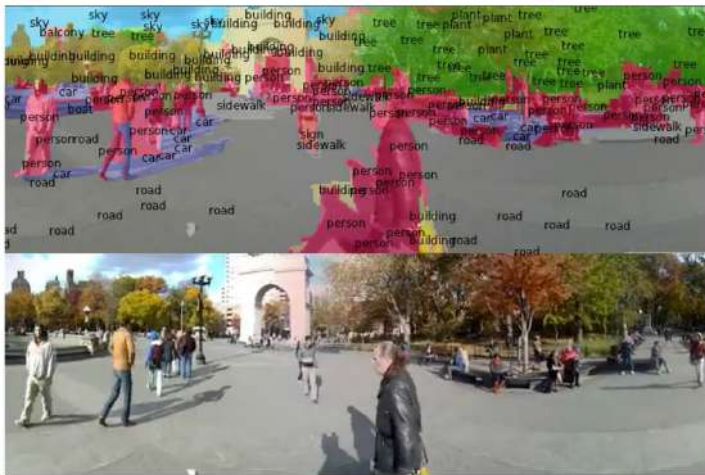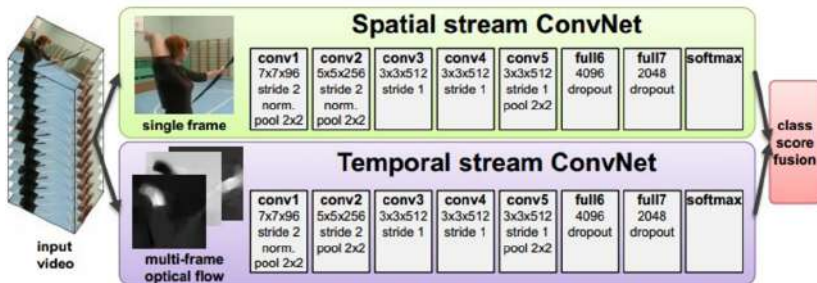## Object Detection



Ren, He, Girshick, and Sun, 2015

Image Segmentation
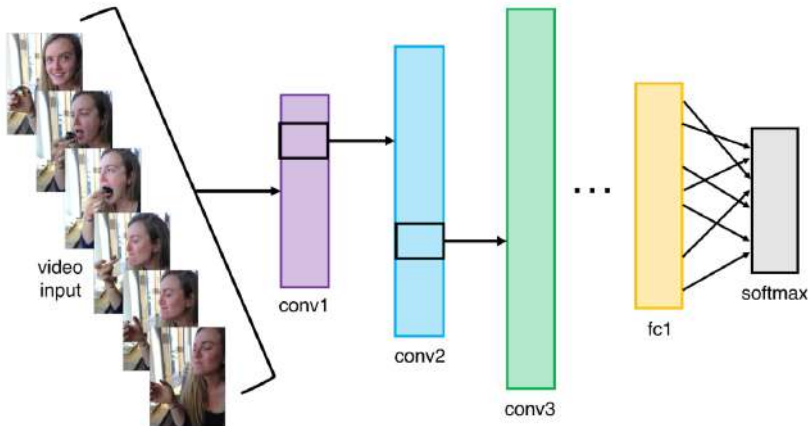


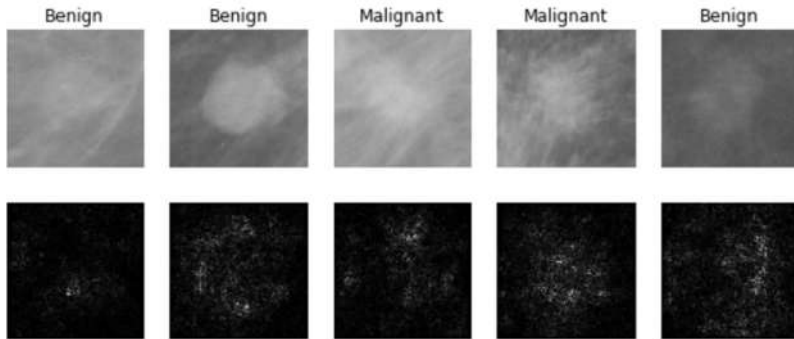Fabaret et al, 2012

## Video Classification



Simonyan et al, 2014

Activity Recognition

Pose Recognition (Toshev and Szegedy, 2014)

Medical Imaging

Image Captioning
Vinyals et al, 2015
Karpathy and Fei-Fei, 2015

A white teddy bear sitting in the grass

A man in a baseball uniform throwing a ball

A woman is holding a cat in her hand

A man riding a wave on top of a surfboard

A cat sitting on a suitcase on the floor

A woman standing on a beach holding a surfboard

Image Generation



"Teddy bears working on new
AI research underwater with
1990s technology"

DALL-E 2

Style Transfer

3D Vision
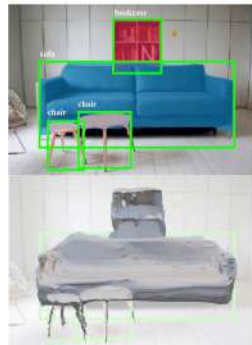


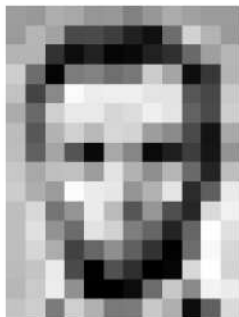Choy et al., 3D-R2N2: Recurrent Reconstruction Neural Network (2016)



Zhou et al., 3D Shape Generation and Completion through Point-Voxel Diffusion (2021)



Gkioxari et al., "Mesh R-CNN", ICCV 2019

# How to represent an image?

- Images are represented as Matrices with elements in $[0, 255]$
- Grayscale images have one channel.

# How to represent an image?

▶ RGB images have 3 channels.

▶ So far, we've worked with tabular data, where each sample consists of structured features. We processed this data using Fully Connected Neural Networks.

**Deep Neural Network**



Figure 12.2 Deep network architecture with multiple layers.

▶ But can we apply the same approach to images?

▶ Let's try.

▶ Let's consider each pixel as a feature.



▶ This should work! But...

# Drawbacks of Fully-Connected Neural Networks

▶ The number of trainable parameters becomes extremely large



256 weights

26 wheights

x1

x25

x256

**Input Image**
16 * 16

100 hiden unit

node

A

Z

$$25600 + 100 + 2600 + 26 = 28326$$

▶ Little or no invariance to shifting, scaling, and other forms of distortion

▶ In other words, the Fully connected NN cannot recognize that the two images (original and shifted) represent the same object (letter "A").



Figure 2: Original "A" character.



Figure 3: Shifted "A" character.

# Drawbacks of Fully-Connected Neural Networks (cont.)

▶ The topology of the input data is completely ignored (it treats pixels as independent features rather than structured patterns.)

▶ For a $32 \times 32$ image, we have

- Black and white patterns: $2^{32*32} = 2^{1024}$

- Grayscale patterns: $256^{32*32} = 256^{1024}$

▶ We need a model that can:

- **Find Patterns in Images**: Recognize small features like edges and shapes in different parts of the image.

- **Work Efficiently**: Avoid looking at every pixel separately by focusing on groups of pixels together.

- **Handle Changes**: Still recognize the same object even if it moves, rotates, or looks slightly different.

# What do we need?

▶ We need a model that can:

- **Find Patterns in Images**: Recognize small features like edges and shapes in different parts of the image.

- **Work Efficiently**: Avoid looking at every pixel separately by focusing on groups of pixels together.

- **Handle Changes**: Still recognize the same object even if it moves, rotates, or looks slightly different.

- We need **Convolutional Neural Networks (CNNs)!**

# Convolutional Neural Networks (CNNs)

▶ CNNs are neural networks designed for image processing and consist of two main parts:

- **Feature Extractor**: Automatically learns patterns (features) such as edges, textures, and shapes from the image.

- **Classifier**: The extracted features are flattened into a vector and passed to a standard neural network (like the ones we used before) for classification.



Figure 4: Convolutional Neural Network

# Convolutional Neural Networks (CNNs)

▶ The feature extractor consists of three essential components:

- **Convolution Layers**: Detect local patterns such as edges, textures, and shapes by applying filters to the image.

- **Activation Function (ReLU)**: Adds non-linearity.

- **Pooling Layers**: Reduce the spatial size of extracted features.

▶ Let's start with the convolution layer.

▶ Let's consider this image.

1. The image is divided into small overlapping tiles (regions).

2. We process each tile using the weights matrix of a small neural network. We call this matrix a **kernel** or **filter**.

Processing a single tile

3. Finally, the filter slides across the image (with the same weights) and processes all the tiles, creating an output **feature map**.

▶ What we did in the last step is called the **Convolution Operation**.

▶ We convolved the kernel with the image, which means sliding the kernel over the entire image and computing the dot product between the kernel and small regions of the image at each step.



Kernel

Input Image

Output Image

▶ This can be represented mathematically by:

$$z = W * x_{i,j} = \sum_{a=0}^{m-1} \sum_{b=0}^{n-1} W_{ab} X_{(i+a)(j+b)}$$

- $z$: Output of the convolution at $(i, j)$.

- $W$: Filter (kernel) matrix of size $m \times n$.

- $x_{i,j}$: Input value at position $(i, j)$.

- $W_{ab}$: Weight of the filter at $(a, b)$.

- $X_{(i+a)(j+b)}$: Input value in the small region at $(i + a, j + b)$.

The **kernel** slides across the image and produces an output value at each position

The **kernel** slides across the image and produces an output value at each position

The **kernel** slides across the image and produces an output value at each position

- **Receptive Field**: The region of the input image that the kernel operates on at each step.

- **Feature Map**: The features extracted from the input image.

We convolve multiple kernels and obtain multiple feature maps or **channels**

▶ **Filters Example**: Different filters detect patterns like edges, textures, or smoothness, producing corresponding feature maps.

$$\begin{bmatrix} 1 & 0 & -1 \\ 0 & 0 & 0 \\ -1 & 0 & 1 \end{bmatrix} \qquad \begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$$



$$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix} \qquad \begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix} \qquad \frac{1}{16}\begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$$

▶ Applying Convolution as such reduces the size of the borders.

▶ Sometimes this is not desirable.

▶ Solution: We can use **Padding**, which pads the border with zeros.

▶ it has two types:

1. **Same Convolution**: Padding is added so the output size equals the input size.

# Controlling the Convolution Process (cont.)

2. **Full Convolution**: Padding is added so the kernel covers the entire input, including edges.

output size = input size + kernel size - 1

# Controlling the Convolution Process (cont.)

▶ **Strided Convolution**: Kernel slides along the image with a step $> 1$
▶ Larger strides reduce output size.

- **Strided Convolution**: Kernel slides along the image with a step $> 1$
- Larger strides reduce output size.

▶ **Pooling**: Compute mean or max over small windows to reduce resolution and extract more general features.

## Single depth slice

| 1 | 1 | 2 | 4 |
|---|---|---|---|
| 5 | 6 | 7 | 8 |
| 3 | 2 | 1 | 0 |
| 1 | 2 | 3 | 4 |

x

y

max pool with 2x2 filters
and stride 2

$\longrightarrow$

| 6 | 8 |
|---|---|
| 3 | 4 |

- No learnable parameters
- Introduces spatial invariance

▶ **Dilated Convolution**: Kernel is spread out, step $> 1$ between kernel elements.

▶ It expands the receptive field without increasing the number of parameters, which makes it efficient.

# CNN Output size

$$n_{out} = \left\lfloor \frac{n_{in} + 2p - k}{s} \right\rfloor + 1$$

$n_{in}$ : number of input features

$n_{out}$ : number of output features

$k$ : convolution kernel size

$p$ : convolution padding size

$s$ : convolution stride size

# Activation

▶ Just like Fully-Connected Neural Networks, we can apply an activation over convolutional layer outputs

▶ It helps break linearity

▶ For example, Rectified Linear Unit (ReLU): $\sigma(x) = max(0, x)$

**Feature Map**

| 9 | 3 | 5 | -8 |
|---|---|---|----|
| -6 | 2 | -3 | 1 |
| 1 | 3 | 4 | 1 |
| 3 | -4 | 5 | 1 |

ReLU Layer

| 9 | 3 | 5 | 0 |
|---|---|---|---|
| 0 | 2 | 0 | 1 |
| 1 | 3 | 4 | 1 |
| 3 | 0 | 5 | 1 |

# Components of a CNN



Convolution Layers

Pooling Layers

Fully-Connected Layers

Activation Function

Normalization

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

▶ Over time, researchers built advanced CNN architectures to improve performance and efficiency. These architectures introduced key innovations:

- **AlexNet [Krizhevsky et al. 2012]**: The first CNN to achieve breakthrough performance on image classification.

- **VGGNet [Simonyan and Zisserman, 2014]**: Used very deep networks (up to 19 layers).

- **InceptionNet (GoogLeNet) [Szegedy et al., 2014]**: Used multiple filter sizes per layer (Inception modules).

- **ResNet [He et al., 2015]**: Introduced skip connections for training very deep networks.

- **EfficientNet [Tan and Le, 2019]**: Found a scaling method that simultaneously scales a CNN's depth, width, and resolution optimally using a single scaling coefficient.

# AlexNet

▶ First big improvement in image classification.

▶ Made use of CNN, pooling, dropout, ReLU and training on GPUs.

▶ 5 convolutional layers, followed by max-pooling layers; with three fully connected layers at the end

# VGGNet

- **Improvement over AlexNet:** Uses a deeper network with small filters instead of a shallow network with larger filters.

- A stack of 3x3 conv layers (vs. 11x11 conv) has same receptive field, more non-linearities, fewer parameters and deeper network representation.

- Two variants: VGG16 or VGG19 conv layers plus 3 FC layers.

# InceptionNet

- ▶ Going Deep: 22 layers
- ▶ Only 5 million parameters! ($12\times$ fewer than AlexNet, $27\times$ fewer than VGGNet).
- ▶ Introduced efficient "Inception module"
- ▶ Introduced "bottleneck" layers that use 1x1 convolutions to reduce the number of channels and mix information across them.

# InceptionNet (cont.)

▶ **Inception module:** Uses multiple filter sizes (1×1, 3×3, 5×5), in parallel, to capture different features, then combines their outputs.

▶ **Problem:** Making networks deeper does not always improve accuracy.

▶ **Why?** In very deep networks, gradients become extremely small as they move backward through layers, making learning slow or stopping it altogether (**vanishing gradient problem**).

▶ **Solution:** Residual Network (ResNet) introduces **skip connections (residuals)**, allowing information to flow more easily.

# ResNet

- ▶ Very deep networks using residual connections
- ▶ 152-layer model for ImageNet
- ▶ Stacked Residual Blocks
- ▶ **Residual:** A shortcut connection that helps the network pass information through layers more easily.

► What happens when we continue stacking deeper layers on a "plain" convolutional neural network?

▶ What happens when we continue stacking deeper layers on a "plain" convolutional neural network?

▶ What happens when we continue stacking deeper layers on a "plain" convolutional neural network?



▶ 56-layer model performs worse on both test and training error

▶ What happens when we continue stacking deeper layers on a "plain" convolutional neural network?



▶ 56-layer model performs worse on both test and training error

▶ The deeper model performs worse, but it's not caused by overfitting!

▶ **Fact:** Deep models have more representation power (more parameters) than shallower models.

- ▶ **Fact:** Deep models have more representation power (more parameters) than shallower models.

- ▶ **Hypothesis:** The problem is an optimization problem, deeper models are harder to optimize

# ResNet

▶ **Fact:** Deep models have more representation power (more parameters) than shallower models.

▶ **Hypothesis:** The problem is an optimization problem, deeper models are harder to optimize

▶ **Solution:** Use network layers to fit a residual mapping instead of directly trying to fit a desired underlying mapping

# ResNet

▶ **Fact:** Deep models have more representation power (more parameters) than shallower models.

▶ **Hypothesis:** The problem is an optimization problem, deeper models are harder to optimize

▶ **Solution:** Use network layers to fit a residual mapping instead of directly trying to fit a desired underlying mapping



Identity mapping:
$H(x) = x$ if $F(x) = 0$

# ImageNet

- ▶ The most extensive data for Image Classification
- ▶ 3 RGB channels from 0 to 255
- ▶ 14,197,122 images
- ▶ 1000 classes

▶ **Problem:** To improve accuracy, we can:

- Increase the number of layers (**depth**)

- Increase the number of neurons in each layer (**width**)

- Use higher-resolution images (**resolution**)

But finding the right balance of these three was largely based on trial and error.

▶ **Solution:** EfficientNet introduced a **compound scaling** method—a mathematical formula to systematically find the optimal balance among **depth**, **width**, and **resolution**.

Figure 5: EfficientNet Scaling.

▶ **Compound Scaling** Uses a single **scaling coefficient** ($\phi$) to control:

- **Network Depth** ($\alpha^\phi$) $\rightarrow$ More layers

- **Network Width** ($\beta^\phi$) $\rightarrow$ More channels per layer

- **Input Resolution** ($\gamma^\phi$) $\rightarrow$ Larger input images

▶ The goal: find $\alpha, \beta, \gamma$ that balance accuracy & efficiency, then scale up optimally by increasing the global coefficient $\phi$.

▶ EfficientNet optimizes depth, width, and resolution using this constraint:
$$\alpha \cdot \beta^2 \cdot \gamma^2 \approx 2$$

▶ **Why this equation?**

- Increasing **depth** ($\alpha$) increases FLOPs **linearly**.

- Increasing **width** ($\beta$) increases FLOPs **quadratically** ($\beta^2$).

- Increasing **resolution** ($\gamma$) increases FLOPs **quadratically** ($\gamma^2$).

▶ To double total FLOPs, the three factors must be balanced together.

▶ The authors of EfficientNet searched for the best scaling factors on a small baseline model.

▶ They found:

$$\alpha = 1.2, \quad \beta = 1.1, \quad \gamma = 1.15$$

▶ The EfficientNet family (B0 to B7) is generated using:

$$\text{Depth} = 1.2^\phi, \quad \text{Width} = 1.1^\phi, \quad \text{Resolution} = 1.15^\phi$$

| Model | $\phi$ | Depth ($\alpha^\phi$) | Width ($\beta^\phi$) |
|---|---|---|---|
| B0 | 0 | $1.2^0 = 1.0$ | $1.1^0 = 1.0$ |
| B1 | 1 | $1.2^1 = 1.2$ | $1.1^1 = 1.1$ |
| B2 | 2 | $1.2^2 = 1.44$ | $1.1^2 = 1.21$ |
| B3 | 3 | $1.2^3 = 1.73$ | $1.1^3 = 1.33$ |
| B4 | 4 | $1.2^4 = 2.07$ | $1.1^4 = 1.46$ |
| B5 | 5 | $1.2^5 = 2.49$ | $1.1^5 = 1.61$ |
| B6 | 6 | $1.2^6 = 2.99$ | $1.1^6 = 1.77$ |
| B7 | 7 | $1.2^7 = 3.58$ | $1.1^7 = 1.94$ |

Table 1: Scaling EfficientNet from B0 to B7

▶ We multiply these scaling factors by the baseline EfficientNet-B0 values and round to the nearest integer to get the new depth, width, and resolution for each model.

▶ EfficientNet models achieve state-of-the-art accuracy with significantly fewer parameters and FLOPs.

▶ EfficientNet models achieve state-of-the-art accuracy with significantly fewer parameters and FLOPs.

▶ EfficientNet-B7 reaches 84.4% Top-1 and 97.3% Top-5 accuracy on ImageNet.

# EfficientNet

- EfficientNet models achieve state-of-the-art accuracy with significantly fewer parameters and FLOPs.

- EfficientNet-B7 reaches 84.4% Top-1 and 97.3% Top-5 accuracy on ImageNet.

- More efficient than previous CNN models—8.4x smaller and 6.1x faster than competitors.
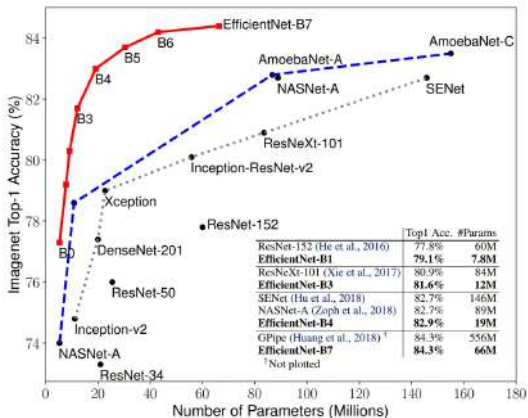
Figure 6: EfficientNet Performance on Imagenet.

► **Why MobileNets?**

- Small-sized models are crucial for mobile and embedded devices.

- MobileNets reduce computational cost and memory usage while maintaining good accuracy.

► **Key Idea:**

- Use **depthwise-separable convolutions** to significantly reduce computation compared to standard convolutions.

▶ **Computational cost of standard convolution:**

Cost = # filter params × # filter positions × # filters

▶ Filters operate on all input channels, increasing computation significantly.



$3 \times 3 \times 3$

$6 \times 6 \times 3$

Filter count = 5

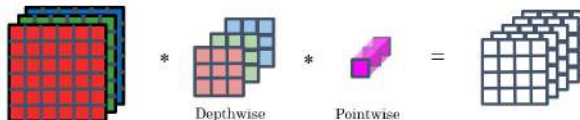$4 \times 4 \times 5$

# Depthwise-Separable Convolutions

▶ Split standard convolution into two steps:

  • **Depthwise Convolution:** Applies a single filter per input channel.

  • **Pointwise Convolution:** Combines outputs from depthwise convolution.

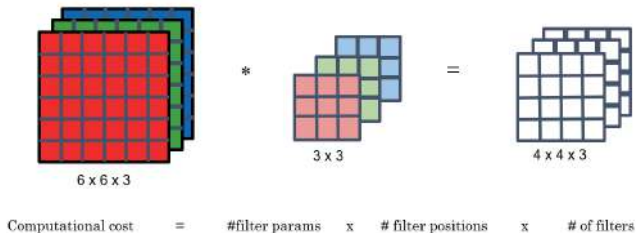▶ **Key Benefit:** Reduces computational cost significantly compared to standard convolution.

▶ Operates on each input channel separately.



$6 \times 6 \times 3$  $3 \times 3$  $4 \times 4 \times 3$

Computational cost = #filter params x # filter positions x # of filters

# Pointwise Convolution

▶ Combines outputs from depthwise convolution using $1 \times 1$ convolutions (mixes channels).



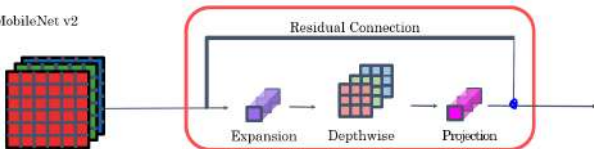|   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|
| Computational cost | = | #filter params | x | # filter positions | x | # of filters |

# MobileNet Architectures

- **MobileNet v2:**
  - Adds **residual connections**.
  - Introduces:
    - **Expansion step:** Expands input dimensions before depthwise convolution.
    - **Projection step:** Reduces dimensions after processing.

# References

These slides have been adapted from

- Fei-Fei Li, Yunzhu Li & Ruohan Gao, Stanford CS231n: Deep Learning for Computer Vision

- Assaf Shocher, Shai Bagon, Meirav Galun & Tali Dekel, WAIC DL4CV Deep Learning for Computer Vision: Fundamentals and Applications

- Justin Johnson, UMich EECS 498.008/598.008: Deep Learning for Computer Vision

- Sander Dieleman, Deepmind: Deep Learning Lecture Series 2020