

Key Concepts in Modern Deep Learning: ResNet, 3D Convolutions, and Attention Modules

ResNet (Residual Network)?

ResNet stands for **Residual Network**, introduced by **He et al. (2015)**.

It's a deep convolutional neural network architecture designed to **train very deep networks** without running into problems like vanishing gradients.

Key Idea

Instead of learning a direct mapping from input to output, ResNet learns a **residual mapping**, i.e., how much the output differs from the input.

Mathematically:

$$y = F(x) + x$$

Where:

- x = input
- $F(x)$ = residual function (what the network learns)
- y = output

So the network essentially learns **how to adjust the input** rather than rebuild it from scratch!

Why Use ResNet?

1. **Enables very deep networks**

- Before ResNet, deeper networks performed worse due to vanishing gradients.
- ResNet allows 50, 101, even 152+ layers to be trained efficiently.

2. **Solves the vanishing/exploding gradient problem**

- The skip connection helps gradients flow directly back through layers.

3. **Improves accuracy**

- Deep networks with residual connections outperform plain networks in classification tasks like ImageNet.

4. **Faster convergence**

- Easier to optimize because the identity shortcut helps retain information.

5. **Generalization**

- Deeper models learn richer representations without overfitting as easily.

How Does It Work?

► **The Residual Block**

A typical residual block looks like this:

Input → Conv → BatchNorm → ReLU → Conv → BatchNorm → Add (Input) → ReLU → Output

- The output of the second convolution is **added to the input** → identity shortcut.
- The network learns the difference (residual) instead of the entire transformation.

What is a Skip Connection?

A **skip connection** is a shortcut that allows the input of a layer to bypass one or more intermediate layers and be added directly to the output.

Why it's important:

- It allows the network to learn a residual function, i.e., only what's different from the input.
- Helps gradients flow more easily during backpropagation → avoids vanishing gradients.
- Makes it easier for deeper networks to converge.

► Example of a Skip Connection

Imagine the following block:

Input → Conv → BN → ReLU → Conv → BN → Add(Input) → ReLU → Output

Mathematically:

$$y = F(x) + x$$

Where:

- x is the input.
- $F(x)$ is the result of intermediate layers.
- y is the final output after adding xx .

If $F(x) = 0$, then the block simply passes the input unchanged → identity mapping.

✓ How is it Better?

Plain Deep Network	ResNet
Struggles with vanishing gradients in very deep models	Solves it via skip connections
Hard to train	Easier and faster to train
Accuracy drops as depth increases	Accuracy improves with more layers
Cannot learn identity function easily	Learns adjustments to identity mapping efficiently

✓ What Issue Does It Solve?

1. Vanishing/Exploding Gradients

- When networks get deeper, gradients shrink or explode during backpropagation → training stalls.
- Skip connections ensure gradients flow smoothly through layers.

2. Degradation Problem

- Deeper networks sometimes have worse performance → ResNet helps by allowing layers to learn only what's necessary.
-

✓ Variants of ResNet

- **ResNet-18, 34, 50, 101, 152** → Increasing depth
 - **ResNeXt, DenseNet** → Extensions improving efficiency or connectivity
 - **Wide ResNet** → Expands width instead of depth
 - Used in classification, object detection, segmentation, and even NLP!
-

✓ Summary

- **ResNet = Residual Network** that uses skip connections to learn residual functions.
 - It solves the **vanishing gradient problem** and allows training **very deep networks**.
 - It's faster, more accurate, and generalizes better than plain deep networks.
 - Skip connections ensure information and gradients flow uninterrupted.
-

Architecture of ResNet

ResNet is built by stacking **residual blocks**. There are two types:

1. **Basic Block** – used in ResNet-18, ResNet-34

2. **Bottleneck Block** – used in ResNet-50, ResNet-101, ResNet-152

► **Basic Residual Block (ResNet-18/34)**

Input → Conv(3×3) → BN → ReLU → Conv(3×3) → BN → Add(input) → ReLU → Output

If input and output dimensions differ → use a **1×1 convolution** on the input to match dimensions.

► **Bottleneck Block (ResNet-50/101/152)**

Input → Conv(1×1) → BN → ReLU → Conv(3×3) → BN → ReLU → Conv(1×1) → BN → Add(input) → ReLU → Output

- Reduces dimensions, processes, and expands → efficient deeper networks.
-

► **ResNet Overall Architecture (ResNet-50 example)**

Input Image (224×224×3)

↓

Conv(7×7, stride=2) + BN + ReLU

↓

MaxPool(3×3, stride=2)

↓

[Residual Block x3] → Layer 1

[Residual Block x4] → Layer 2

[Residual Block x6] → Layer 3

[Residual Block x3] → Layer 4

↓

Global Average Pooling

↓

Fully Connected Layer (e.g. 1000 classes for ImageNet)

Code Snippet in PyTorch

► Implementation of BasicBlock

```
import torch
import torch.nn as nn
import torch.nn.functional as F

class BasicBlock(nn.Module):
    expansion = 1 # For compatibility if needed

    def __init__(self, in_channels, out_channels, stride=1):
        super().__init__()

        # First convolution
        self.conv1 = nn.Conv2d(in_channels, out_channels, kernel_size=3,
                                stride=stride, padding=1, bias=False)
        self.bn1 = nn.BatchNorm2d(out_channels)

        # Second convolution
        self.conv2 = nn.Conv2d(out_channels, out_channels, kernel_size=3,
                                stride=1, padding=1, bias=False)
        self.bn2 = nn.BatchNorm2d(out_channels)

        # Shortcut connection
        self.shortcut = nn.Sequential()
        if stride != 1 or in_channels != out_channels:
            # 1x1 convolution to match dimensions
            self.shortcut = nn.Sequential(
                nn.Conv2d(in_channels, out_channels, kernel_size=1,
                           stride=stride, bias=False),
                nn.BatchNorm2d(out_channels)
            )

    def forward(self, x):
```

```

out = F.relu(self.bn1(self.conv1(x)))
out = self.bn2(self.conv2(out))
out += self.shortcut(x) # Skip connection
out = F.relu(out)
return out

```

► Simple ResNet using BasicBlock

```

class SimpleResNet(nn.Module):
    def __init__(self, block, num_blocks, num_classes=10):
        super().__init__()
        self.in_channels = 64

        # Initial convolution layer
        self.conv1 = nn.Conv2d(3, 64, kernel_size=3, stride=1, padding=1, bias=F
alse)
        self.bn1 = nn.BatchNorm2d(64)

        # Stacked residual layers
        self.layer1 = self._make_layer(block, 64, num_blocks[0], stride=1)
        self.layer2 = self._make_layer(block, 128, num_blocks[1], stride=2)
        self.layer3 = self._make_layer(block, 256, num_blocks[2], stride=2)
        self.layer4 = self._make_layer(block, 512, num_blocks[3], stride=2)

        # Final classifier
        self.fc = nn.Linear(512 * block.expansion, num_classes)

    def _make_layer(self, block, out_channels, num_blocks, stride):
        strides = [stride] + [1]*(num_blocks-1)
        layers = []
        for s in strides:
            layers.append(block(self.in_channels, out_channels, stride=s))
            self.in_channels = out_channels * block.expansion
        return nn.Sequential(*layers)

```

```
def forward(self, x):
    out = F.relu(self.bn1(self.conv1(x)))
    out = self.layer1(out)
    out = self.layer2(out)
    out = self.layer3(out)
    out = self.layer4(out)
    out = F.adaptive_avg_pool2d(out, (1,1))
    out = torch.flatten(out, 1)
    out = self.fc(out)
    return out

# Example usage:
model = SimpleResNet(BasicBlock, [2, 2, 2, 2]) # ResNet-18 style
print(model)
```

How Gradients Flow in ResNet

► The Problem in Plain Networks

In deep networks:

- Gradients shrink exponentially → layers at the beginning hardly get updated → training stalls.

► Skip Connection Helps by Allowing Direct Gradient Flow

$$y = F(x) + x$$

During backpropagation, the gradient of y w.r.t x is:

$$\frac{\partial y}{\partial x} = \frac{\partial F(x)}{\partial x} + 1$$

This "+1" term ensures:

- Even if $\frac{\partial F(x)}{\partial x}$ is very small (vanishing gradient), the gradient can still flow backward through the identity path.

- The network can easily learn to skip layers if deeper processing isn't necessary → prevents degradation.
-

► Intuition

1. The input can bypass intermediate layers → no risk of losing information.
 2. The network can choose to learn the residual or preserve the input → stable learning.
 3. Deep networks can be trained more easily → better performance.
-

✓ Summary

- ✓ **Skip connections** allow information and gradients to bypass layers → avoid vanishing gradients.
 - ✓ ResNet is built from **residual blocks**, stacking them to create deep networks.
 - ✓ ResNet solves the **degradation problem** by making deeper models trainable.
 - ✓ In PyTorch, skip connections are implemented by **adding the input to the output** using simple layers like `nn.Conv2d` and `nn.BatchNorm2d`.
 - ✓ During backpropagation, the "+1" gradient ensures smooth flow → deeper networks converge.
-

2D Convolution to 3D Convolution

✓ What is 3D Convolution?

A **3D Convolution** applies a 3-dimensional filter (kernel) over an input that has three spatial dimensions — typically **depth, height, and width**.

► Use cases:

- Video data (frames across time → temporal dimension)
 - Medical imaging (e.g., CT scans → volumetric data)
 - Scientific data (3D simulations, point clouds)
-

Input shape for 3D convolution:

$$\text{Input} = (D_{in}, H_{in}, W_{in}, C_{in})$$

Where:

- $D_{in} \rightarrow$ Depth (e.g. frames in a video, slices in CT scan)
- $H_{in} \rightarrow$ Height
- $W_{in} \rightarrow$ Width
- $C_{in} \rightarrow$ Channels (e.g. RGB \rightarrow 3)

Concepts in 2D vs 3D Convolution

Concept	2D Convolution	3D Convolution
Input shape	$H \times W \times C$	$D \times H \times W \times C$
Kernel shape	$k_H \times k_W \times C_{in}$	$k_D \times k_H \times k_W \times C_{in}$
Stride	Moves across height and width	Moves across depth, height, and width
Padding	Adds rows/columns to edges	Adds slices, rows, and columns to edges
Pooling	2×2 or 3×3 pooling over spatial dims	$2 \times 2 \times 2$ or $3 \times 3 \times 3$ pooling over depth+spatial dims
Output	2D feature maps	3D feature volumes
Applications	Images	Videos, CT/MRI scans, scientific simulations

Detailed Explanation of Concepts

► 1. Kernel (Filter)

- **2D kernel:**
Looks at a patch across height & width only.
- **3D kernel:**

Looks at a patch across depth, height, and width.

Example:

- 2D kernel $\rightarrow 3 \times 3$ (covers nearby pixels)
- 3D kernel $\rightarrow 3 \times 3 \times 3$ (covers nearby frames + pixels)

It learns spatiotemporal patterns or volumetric features.

► 2. Padding

- **2D padding:** Add rows/columns at the edges \rightarrow controls output size.
- **3D padding:** Add slices along depth, and rows/columns along height and width \rightarrow preserves temporal or volumetric structure.

Example:

For a video of shape $8 \times 64 \times 64$:

- Padding = 1 \rightarrow output keeps depth as 8 instead of shrinking.
-

► 3. Stride

- **2D stride:** Moves the kernel by `s_H` and `s_W` in spatial directions.
- **3D stride:** Moves the kernel by `s_D`, `s_H`, and `s_W` in depth and spatial directions.

Example:

Stride = (2, 2, 2) \rightarrow skips frames and pixels \rightarrow reduces resolution in all directions.

► 4. Pooling

- **2D pooling:** Reduces size by summarizing regions, e.g., max pooling over 2×2 .
- **3D pooling:** Reduces depth, height, and width simultaneously, e.g., max pooling over $2 \times 2 \times 2$.

Used to compress spatiotemporal or volumetric data.

► 5. Kernel Sliding

- **2D sliding:** Kernel moves across the image grid.
- **3D sliding:** Kernel moves across volume or video frames → at each position, it summarizes depth + spatial info.

Example:

A 3×3×3 kernel slides across slices (depth), and each frame (height × width).

► 6. Output Size Calculation

For **2D**, output dimensions:

$$H_{out} = \frac{H_{in} + 2P_H - k_H}{s_H} + 1$$

$$W_{out} = \frac{W_{in} + 2P_W - k_W}{s_W} + 1$$

For **3D**, output dimensions:

$$D_{out} = \frac{D_{in} + 2P_D - k_D}{s_D} + 1$$

$$H_{out} = \frac{H_{in} + 2P_H - k_H}{s_H} + 1$$

$$W_{out} = \frac{W_{in} + 2P_W - k_W}{s_W} + 1$$

Where P = padding, k = kernel size, s = stride.

✓ PyTorch Example of 3D Convolution

```
import torch
import torch.nn as nn

# Example input: batch_size=1, channels=3, depth=8, height=64, width=64
input_tensor = torch.randn(1, 3, 8, 64, 64)
```

```
# 3D Convolution
conv3d = nn.Conv3d(in_channels=3, out_channels=16, kernel_size=(3, 3, 3), s
tride=1, padding=1)

# Apply convolution
output = conv3d(input_tensor)

print("Input shape:", input_tensor.shape)
print("Output shape:", output.shape)
```

Output explanation:

- Input $\rightarrow [1, 3, 8, 64, 64]$
- Output $\rightarrow [1, 16, 8, 64, 64]$ (since padding=1, stride=1 \rightarrow depth & spatial dimensions preserved)

✓ Summary Table – 2D vs 3D Convolution

Concept	2D Convolution	3D Convolution
Input	$H \times W \times C$	$D \times H \times W \times C$
Kernel	$k_H \times k_W$	$k_D \times k_H \times k_W$
Stride	2D steps	3D steps (depth + spatial)
Padding	Add rows/columns	Add slices/rows/columns
Pooling	Spatial pooling only	Depth + spatial pooling
Applications	Images	Videos, CT scans, volumetric data
Output size	2D grid	3D volume

✓ Final Notes

- ✓ 3D convolutions generalize 2D convolutions to handle data with an extra dimension (temporal or volumetric).
- ✓ All operations — padding, stride, pooling, kernel sliding — extend from 2D to 3D by considering depth alongside spatial dimensions.

✓ 3D convolutions are widely used in video processing, medical imaging, and other scientific domains.

✓ PyTorch's `nn.Conv3d` makes it easy to implement these operations with minimal code.

Can we use 2D convolution on video data?

Why We Might Use 2D Convolution on Video Data

1. Treat each frame independently

- A video is a sequence of images (frames).
- We can process each frame with a 2D convolution to extract spatial features like edges, textures, shapes.
- The temporal relationship (between frames) is ignored in this setup.

Example:

You can apply a CNN to each frame separately for tasks like frame-wise classification.

2. Combine frame features after extraction

After applying 2D convolutions on individual frames, you can aggregate information over time using:

- ✓ **Recurrent Neural Networks (RNNs)** – like LSTMs, GRUs
- ✓ **Temporal pooling** – averaging or max pooling over time
- ✓ **Attention mechanisms** – learning to focus on key frames

This hybrid approach is used in many practical pipelines.

3. Use 2D convolutions with frame stacking

- Stack several consecutive frames along the channel dimension → treat it like a multi-channel image.
- Apply 2D convolutions over the combined frame stack.

Example:

Input shape $\rightarrow H \times W \times (3 \times T)$, where T is the number of frames stacked.

This approach captures short-term temporal cues.

✓ Limitations of Using 2D Convolutions on Video

Limitation	Explanation
No temporal modeling	Relationships between frames (motion, flow) aren't explicitly learned
Short-term context only	Even if you stack frames, you can only capture small temporal windows
Computational redundancy	Same 2D filters applied to each frame separately without temporal context

✓ Why Use 3D Convolutions Instead?

- ✓ They process both spatial and temporal patterns at once
 - ✓ Learn motion features, changes over time, and depth information
 - ✓ More suitable for video classification, action recognition, and video generation
- But they are computationally heavier and require more data.

✓ Real-world Use Cases of 2D Convolutions on Video

- ✓ **Video surveillance** – detect events frame by frame
- ✓ **Object detection in videos** – per-frame detection pipelines
- ✓ **Action recognition (simplified)** – stacking frames as input
- ✓ **Video compression** – frame-level processing before encoding
- ✓ **Medical imaging sequences** – slice-by-slice segmentation

✓ Example Approaches Using 2D Convolution on Video

1. Frame-wise CNN + LSTM

Extract spatial features per frame → sequence modeling over time.

2. Temporal pooling on frame features

Average features from multiple frames → reduce to a single prediction.

3. 3-channel stacking (early fusion)

Combine frames → learn joint spatial patterns → used in lightweight models.

Summary

- ✓ Yes, you can apply **2D convolutions on video data** by treating frames as independent or by stacking them along the channel axis.
 - ✓ It's useful for simpler tasks or when computational resources are limited.
 - ✓ However, it doesn't capture long-term temporal dependencies unless paired with other methods like RNNs or attention mechanisms.
 - ✓ For richer video understanding (action recognition, motion capture), **3D convolutions** or advanced architectures like transformers are more appropriate.
-

Visual Attention?

Definition

Visual attention is a technique that allows a model to selectively focus on the **most informative regions** in an image or video, rather than treating every pixel equally. This helps the model process only the important parts, improving efficiency and performance.

Key Concepts

1. Focus Mechanism

The model focuses on important regions in an image or video. For example, if you're trying to detect a cat, the model learns to pay attention to the cat's body rather than the background.

2. Feature Weighting

Attention assigns different importance (weights) to various spatial locations or feature channels. Important areas/features get higher weights, while irrelevant areas get lower weights.

3. Types of Visual Attention

- **Spatial Attention** → Focuses on important locations in the image (e.g., corners, objects).
- **Channel Attention** → Focuses on which feature channels (like edges, textures) are more useful.
- **Self-Attention** → Relates all pixels or patches to each other to understand global relationships, as seen in **Vision Transformers (ViTs)**.

4. Applications

Visual attention is used in:

- **Object detection** → Detect specific objects in images.
- **Image captioning** → Describe images using natural language.
- **Scene understanding** → Interpret scenes holistically.
- **Transformers** → Capture long-range relationships beyond what convolution layers can cover.

Why is Visual Attention Important?

- ✓ It allows the model to focus computational resources on important regions.
- ✓ It improves interpretability by showing what the model is paying attention to.
- ✓ It helps capture both local and global context, especially in large or complex scenes.
- ✓ It powers advanced models like Vision Transformers and other deep learning architectures.

Example

For instance, when processing a crowded street scene:

- **Spatial attention** might highlight the person crossing the road.
 - **Channel attention** might focus on features like edges and contours.
 - **Self-attention** could help understand relationships, like how objects relate across the entire image.
-

Basic Attention Module?

The **Basic Attention Module** is a neural network component designed to help models dynamically focus on the most important parts of input data — whether that's an image, video, or feature map.

Instead of treating all parts of the input equally, the attention module learns **where to look** or **what to emphasize**. It computes attention weights that highlight the regions or features most useful for solving a task, such as recognizing objects or detecting patterns.

How It Works

► 1. Input Feature Map

- The process starts with an **input feature map**, which is the output of earlier convolutional layers.
 - It contains spatial and channel-wise information extracted from the image or video.
-

► 2. Convolutional Bottleneck

- The input passes through a **convolutional bottleneck**.
 - This reduces dimensionality or processes features in a more efficient way while retaining essential patterns.
 - It helps extract relevant information before computing attention.
-

► 3. Multi-Layer Perceptron (MLP)

- The bottleneck output is fed into a **multi-layer perceptron (MLP)**.

- The MLP computes interactions between features to learn which parts of the input are important.
 - It helps model non-linear relationships and refine the attention mechanism.
-

► 4. Sigmoid Activation

- The output of the MLP goes through a **sigmoid function**.
 - The sigmoid ensures that the attention weights are between 0 and 1, representing how much focus should be given to each part of the input.
-

► 5. 2D/3D Activation Map

- The result is an **activation map** (either spatial for 2D or spatiotemporal for 3D data).
 - It highlights regions or channels of the feature map that are important.
-

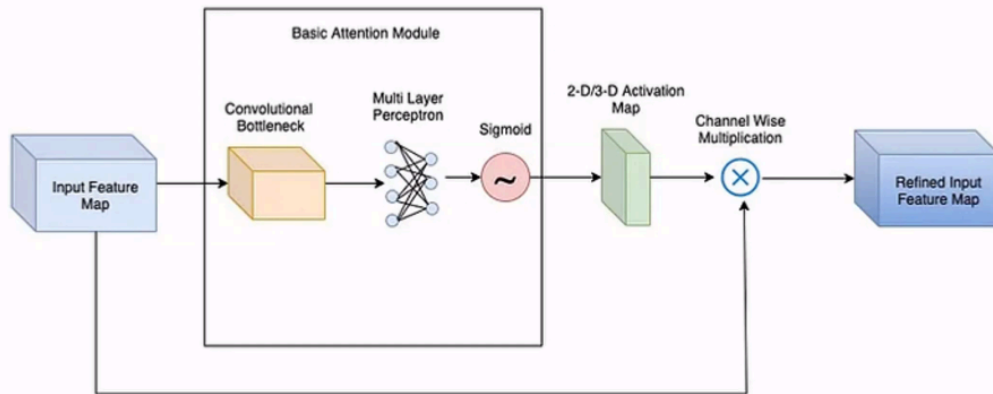
► 6. Channel-Wise Multiplication

- The activation map is multiplied with the original input feature map **channel-wise**.
 - This means the important regions are amplified while irrelevant areas are suppressed.
-

► 7. Refined Input Feature Map

- The output is a **refined feature map**, where the network can better focus on significant parts of the image or video.
- It improves the representational capacity and interpretability of the model.

Basic attention module



✓ When to Use the Basic Attention Module

You should consider using this module when:

Important regions/features vary dynamically

- For example, in images where the object of interest changes across samples or videos where motion patterns are inconsistent.

You want interpretability

- It helps visualize which parts of the input the model is focusing on.

Tasks involve noisy or cluttered data

- Attention helps the model ignore irrelevant information and focus on the signal.

You want to improve performance without drastically increasing model size

- It's an efficient way to guide learning without adding heavy computation.

Applications

- **Object detection** → focusing on regions where objects are likely to appear.
- **Medical imaging** → highlighting anomalies like tumors or lesions.

- **Image captioning** → picking key objects or actions for generating descriptions.
 - **Scene understanding** → interpreting complex images with multiple objects.
 - **Video recognition** → identifying moving or temporally important regions.
 - **Transformers and ViTs** → capturing long-range dependencies and global patterns.
-

✓ Summary

The **Basic Attention Module** is a lightweight yet powerful mechanism that helps models learn where to focus by computing attention weights over spatial or channel-wise features. It's useful when input data is complex, noisy, or has variable importance across regions. By guiding the model's focus, it enhances both interpretability and performance, making it widely applicable across computer vision tasks.

BAM and CBAM

BAM – Block Attention Module

► What is BAM?

- It's a **hybrid attention mechanism** combining:
 1. **Spatial attention** → Focuses on important regions in the image.
 2. **Channel attention** → Highlights the most useful feature channels.
 - These are combined into a single module, allowing the network to learn both **where** and **what** to focus on.
-

► How BAM Works

1. **Input Feature Map (F)** → passed through both channel and spatial attention paths.
2. **Channel Attention (Mc(F))**

- Global average pooling → extracts summary statistics.
- Fully connected layers → learn which channels are important.
- Batch normalization → stabilizes learning.

3. Spatial Attention ($M_s(F)$)

- Convolutions with dilation → extract patterns across spatial locations.
- Focuses on the key regions of the image.

4. Combination

$$M(F) = \sigma(M_c(F) + M_s(F))$$

- The outputs of both attention paths are summed and passed through a sigmoid activation.

5. Feature Refinement

$$F' = F + F \otimes M(F)$$

- The attention map is multiplied channel-wise with the original feature map to refine it.

► When to Use BAM

- ✓ To improve CNNs like ResNet and DenseNet
- ✓ When you need both spatial and channel-wise attention
- ✓ For tasks such as classification, object detection, and segmentation
- ✓ When you want a lightweight, plug-and-play module that doesn't increase computational cost too much

CBAM – Convolutional Block Attention Module

► What is CBAM?

- CBAM extends BAM by applying attention **sequentially**, first over channels and then over spatial locations.

- It's designed to capture more fine-grained dependencies by separating the attention mechanisms.

► How CBAM Works

1. **Input Feature Map** → passed through two separate attention modules.
2. **Channel Attention Module ($M_c(F)$)**
 - Applies **average pooling** and **max pooling** along the spatial dimension.
 - Both pooled results are passed through shared MLPs and combined.
 - A sigmoid activation generates the channel attention map.
3. **Spatial Attention Module ($M_s(F)$)**
 - Uses average and max pooling along the channel dimension.
 - Combines the results and applies a convolutional layer.
 - The output attention map is used to focus on important spatial areas.
4. **Sequential Application**

$$F' = M_c(F) \otimes FF'' = M_s(F') \otimes F'$$

- The channel attention is applied first, followed by spatial attention.

► When to Use CBAM

- ✓ When fine-grained attention is needed in both dimensions
- ✓ For improving feature representation at multiple levels
- ✓ In tasks where both spatial and channel information is critical (e.g., medical imaging, video analysis)

✓ Differences Between BAM and CBAM

Feature	BAM	CBAM
Type	Hybrid (parallel)	Sequential (channel → spatial)

Feature	BAM	CBAM
Complexity	Simpler	Slightly more complex
Attention flow	Combined attention	Separate attentions applied step by step
Use cases	Lightweight attention enhancement	Fine-grained attention for critical tasks
Implementation	One sigmoid after sum	Two sigmoids after channel and spatial processing

✓ Summary

- ✓ **BAM** combines spatial and channel attention in a single step, helping CNNs focus on both *where* and *what* is important.
- ✓ **CBAM** applies attention sequentially, refining features in two steps to better capture dependencies across both dimensions.
- ✓ Both modules enhance performance in tasks like classification, detection, and segmentation without heavy computation.
- ✓ They are widely used in deep learning pipelines, especially in architectures like ResNet or DenseNet.

<https://medium.com/visionwizard/understanding-attention-modules-cbam-and-bam-a-quick-read-ca8678d1c671>

Class Activation Maps (CAM)

Class Activation Maps (CAM) are a visualization technique used in Convolutional Neural Networks (CNNs) to **highlight the regions in an input image that are most important for a model's prediction of a particular class**.

They help us answer: *"Which parts of the image did the CNN focus on to decide this is a cat, dog, car, etc.?"*

🔑 Core Idea

- CNNs learn feature maps at different layers.

- The last convolutional layer retains **spatial information** (where patterns are located in the image).
- The fully connected layers (after flattening) usually discard this spatial info.
- **CAM works by connecting the final convolutional feature maps directly to the class score** before softmax, using a weighted sum.
- This gives a heatmap showing **which regions contributed most** to that class prediction.

How CAM is computed (step by step):

Suppose we have a CNN trained for image classification.

1. Feature Maps Extraction

Let the last convolutional layer output K feature maps:

$$f_k(x, y), \quad k = 1, 2, \dots, K$$

where (x, y) are spatial locations.

2. Global Average Pooling (GAP)

For each feature map, compute the average activation:

$$F_k = \frac{1}{Z} \sum_{x, y} f_k(x, y)$$

where Z is the number of pixels in the feature map.

3. Fully Connected Layer Weights

Each class c has weights w_k^c connecting feature F_k to class score S_c .

So the class score is:

$$S_c = \sum_k w_k^c F_k$$

4. Back-project to Spatial Domain

Instead of averaging, we combine the **raw feature maps** using the learned weights:

$$M_c(x, y) = \sum_k w_k^c f_k(x, y)$$

$M_c(x, y)$ is the **Class Activation Map** for class c .

👉 This map shows *which pixels strongly contribute* to predicting class c.

Intuition

- If a CNN predicts "dog," the CAM heatmap will glow over the **dog's face/body** regions in the image.
- If it predicts "car," the heatmap highlights **car-shaped regions**.

It's like peeking inside the CNN's "attention," but unlike self-attention in transformers, this is based on **convolutional feature importance**.

Advantages

- Provides **interpretability** for CNN decisions.
 - Helps debug when models focus on wrong regions (e.g., background instead of object).
 - Useful for weakly supervised tasks (localization without bounding box labels).
-

Limitations

- Requires **Global Average Pooling (GAP)** in architecture (original CAM works only with this setup).
 - For architectures without GAP (e.g., ResNet, VGG), we use **Grad-CAM**, a generalization.
-

Use Cases

- Model interpretability
 - Debugging misclassifications
 - Weakly supervised object localization
-

Gradient-weighted Class Activation Mapping (Grad-CAM)

Gradient-weighted Class Activation Mapping (Grad-CAM) is a visualization technique that produces class-discriminative heatmaps by using the **gradients** of a target class flowing into the last convolutional layer.

It answers: *"For this prediction, which regions in the image had the greatest influence?"*

Core Idea

- CAM needed **Global Average Pooling (GAP)** to map feature maps \rightarrow class score.
 - Grad-CAM avoids this restriction by instead computing **importance weights using gradients** of the target class score w.r.t. feature maps.
-

How Grad-CAM is computed (step by step):

Suppose a CNN outputs a class score S_c for class c .

1. Forward pass

Get feature maps from the last convolutional layer:

$$A^k(x, y), \quad k = 1, \dots, K$$

2. Backward pass

Compute gradients of the score S_c w.r.t. each feature map:

$$\frac{\partial S_c}{\partial A^k(x, y)}$$

3. Compute importance weights

Take the **average gradient** over all spatial positions:

$$\alpha_k^c = \frac{1}{Z} \sum_x \sum_y \frac{\partial S_c}{\partial A^k(x, y)}$$

where Z is the number of pixels in the feature map.

👉 This weight α_k^c tells us **how important feature map k is for predicting class c**.

4. Weighted sum of feature maps

Combine feature maps using these weights:

$$L_{\text{Grad-CAM}}^c(x, y) = \text{ReLU} \left(\sum_k \alpha_k^c A^k(x, y) \right)$$

- The **ReLU** ensures only positive contributions are kept (we only highlight what supports the class, not what suppresses it).

Intuition

- If the CNN predicts "dog," the gradients show *which feature maps were most sensitive to the dog class*.
- We weight those maps and project them back spatially → heatmap.
- The red regions in the heatmap = "This part of the image made the CNN think it's a dog."

Advantages

- Works with **any CNN architecture** (no need for GAP).
- More flexible than CAM.
- Produces interpretable heatmaps aligned with human intuition.

Limitations

- Resolution limited by the last conv layer (coarse heatmaps).
 - Sometimes highlights broad regions, not sharp object boundaries.
 - Only shows **where**, not **why**, a model is focusing.
-