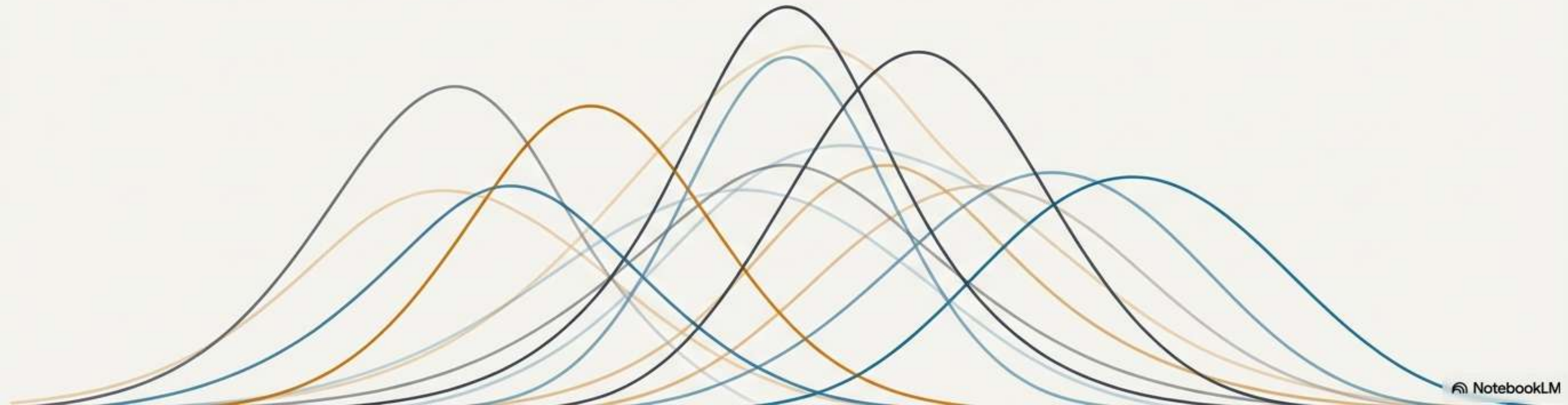


# The Shifting Target

Understanding and Mitigating Internal  
Covariant Shift in Deep Neural Networks





# We Begin With a Familiar Problem: Covariant Shift

**Covariant shift** occurs when the distribution of input data changes between the training and testing phases, while the conditional distribution of the output given the input remains unchanged.

## Key Idea

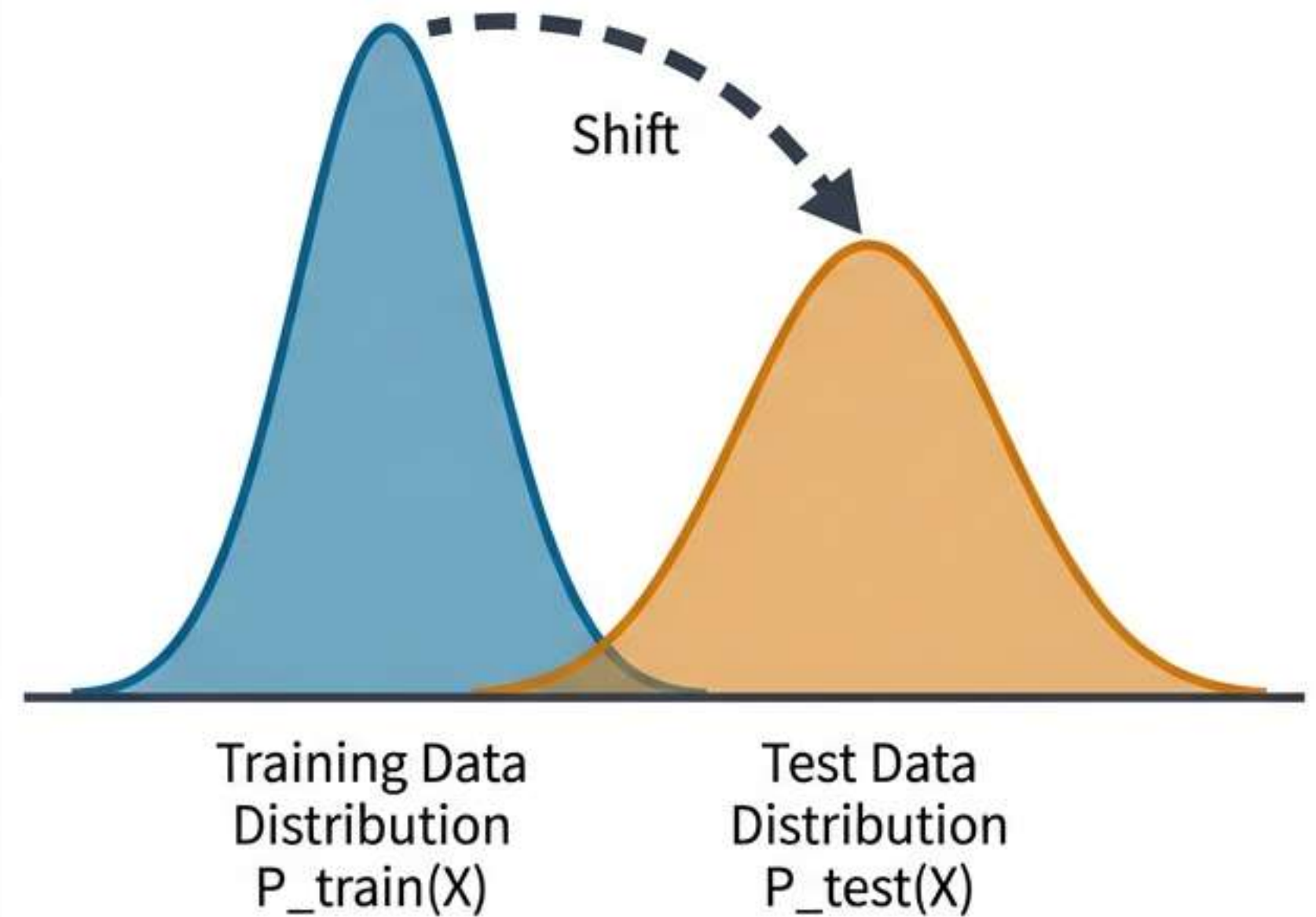
The model learns from one world (training data) but is expected to perform in another (test data).

Training Distribution:  $P_{\text{train}}(X)$

Test Distribution:  $P_{\text{test}}(X)$

**The Problem:**  $P_{\text{train}}(X) \neq P_{\text{test}}(X)$

The Constant:  $P(Y|X)$  remains the same.

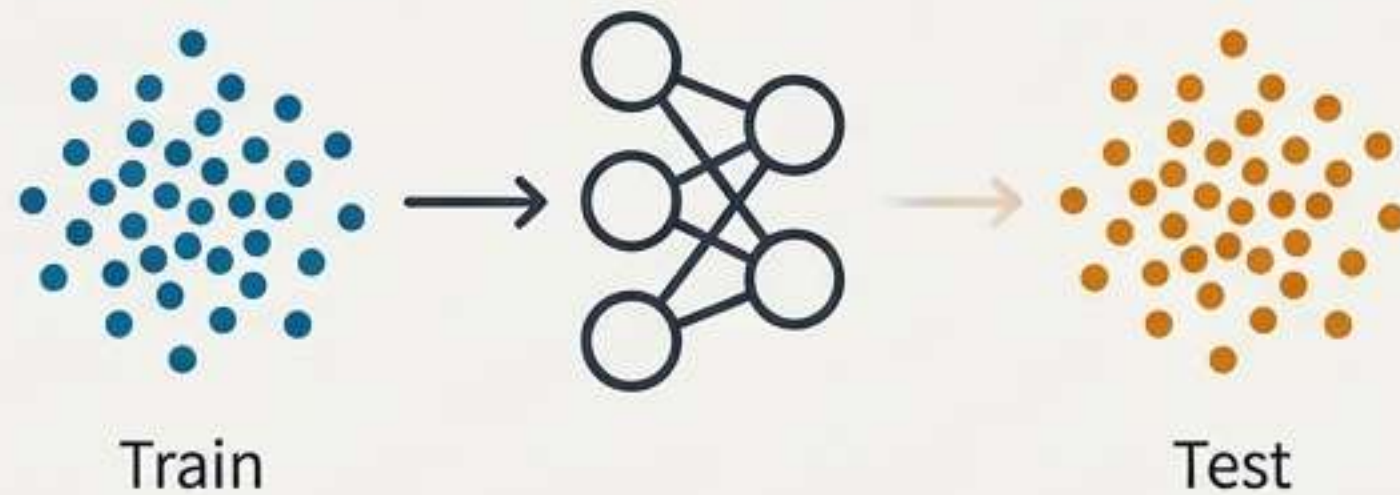




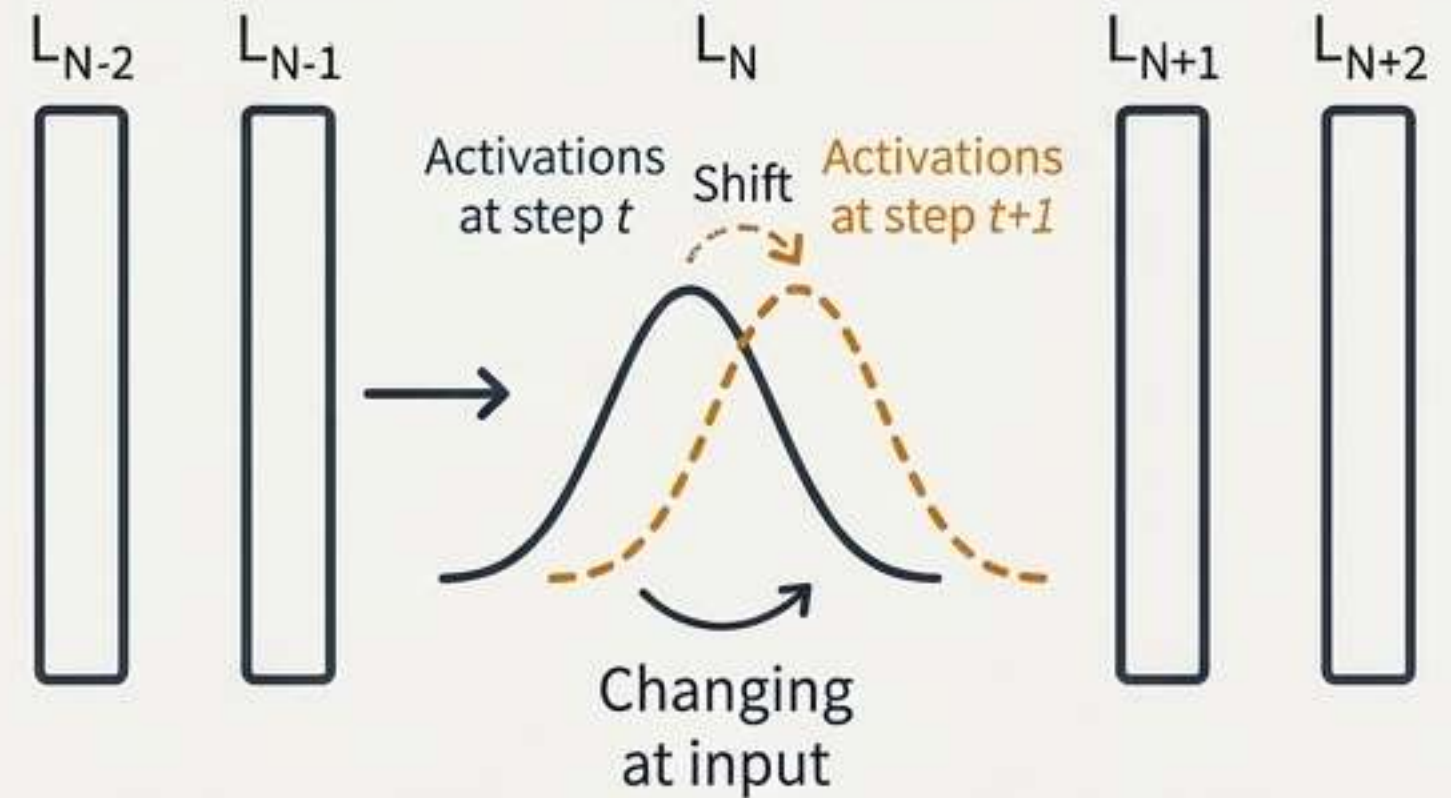
# Now, Imagine This Shift Happening *Inside* the Network, On Every Step

Internal Covariant Shift is the **change in the distribution of network activations** due to updates in network parameters during training. Each layer must continuously adapt to a new input distribution from the layer before it.

External Covariant Shift



Internal Covariant Shift





# What Causes This Internal Instability?

## Changes in Data Distribution

**Core Driver:** As network parameters (weights and biases) are updated, the output distribution of each layer changes. This change propagates through the network.

### Examples:

- **\*\*Batch Training\*\***: Each mini-batch has a slightly different data distribution, causing continuous shifts in activations.
- **\*\*Data Augmentation\*\***: Techniques like rotation or scaling alter input distributions, which in turn changes activation distributions layer by layer.

## The Role of Training Dynamics

**Core Driver:** The optimization process itself creates a moving target for each layer.

### Examples:

- **\*\*Non-Linear Activations (ReLU, Sigmoid)\*\***
- **\*\*Non-Linear Activations (ReLU, Sigmoid)\*\***: Can amplify small changes in inputs, leading to dramatic shifts in output distributions.
- **\*\*Layer Depth\*\***: In deep networks, small shifts in early layers are amplified as they propagate, causing significant instability in later layers.



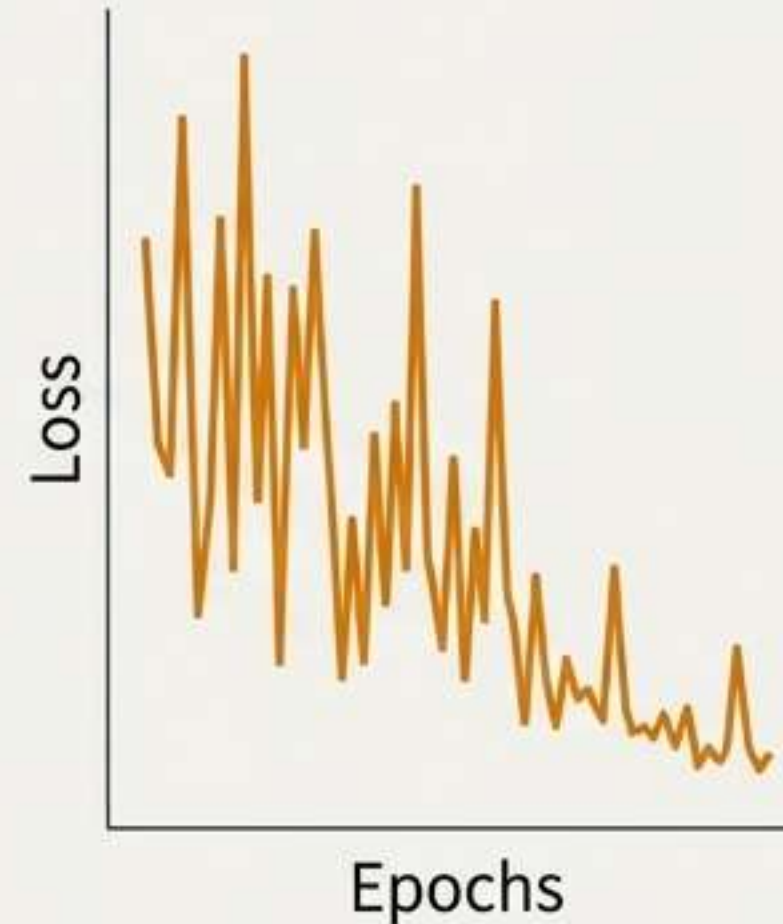
# The Consequences of a Constantly Shifting Target

## Impact on Training

- **Slower Convergence:** Layers must constantly re-adapt to their changing inputs, increasing the number of epochs required to train.
- **Training Instability:** The learning process can oscillate, struggling to find a stable path to minimizing the loss function.
- **Vanishing/Exploding Gradients:** ICS can exacerbate gradient issues. As activation statistics shift, gradients can become extremely small or large, hindering effective weight updates.
- **Poor Generalization:** Learned representations may be less robust because layers are always chasing a moving target, leading to poorer performance on new data.

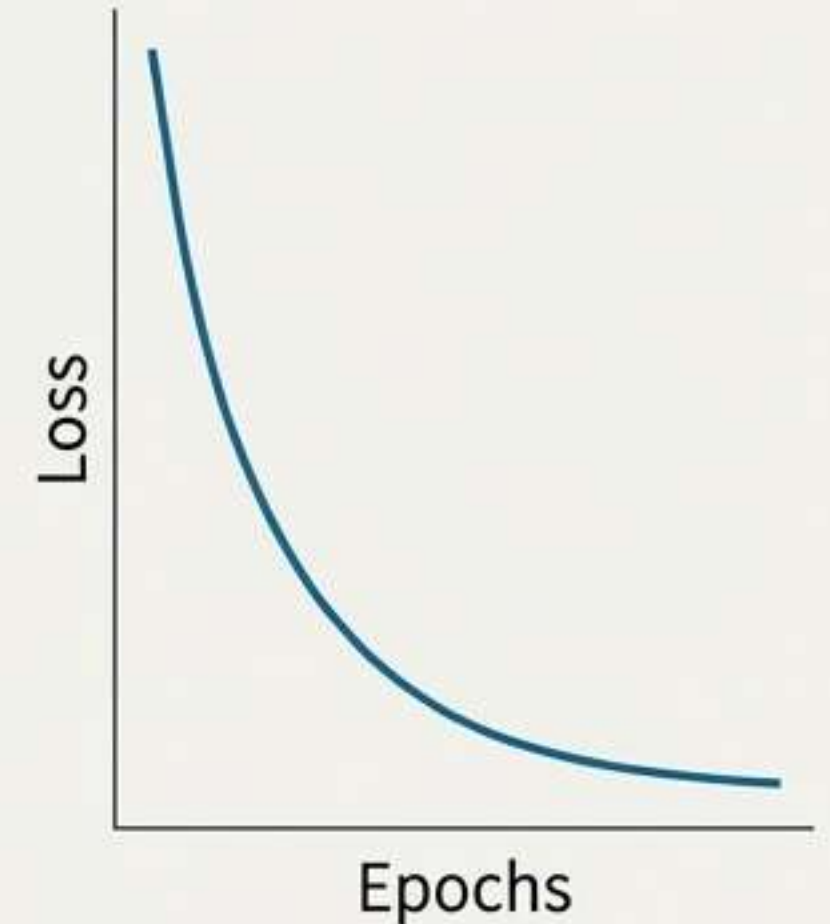
## Unstable Training

Loss vs. Epochs



## Stabilized Training

Loss vs. Epochs





# The Foundational Solution: Batch Normalization

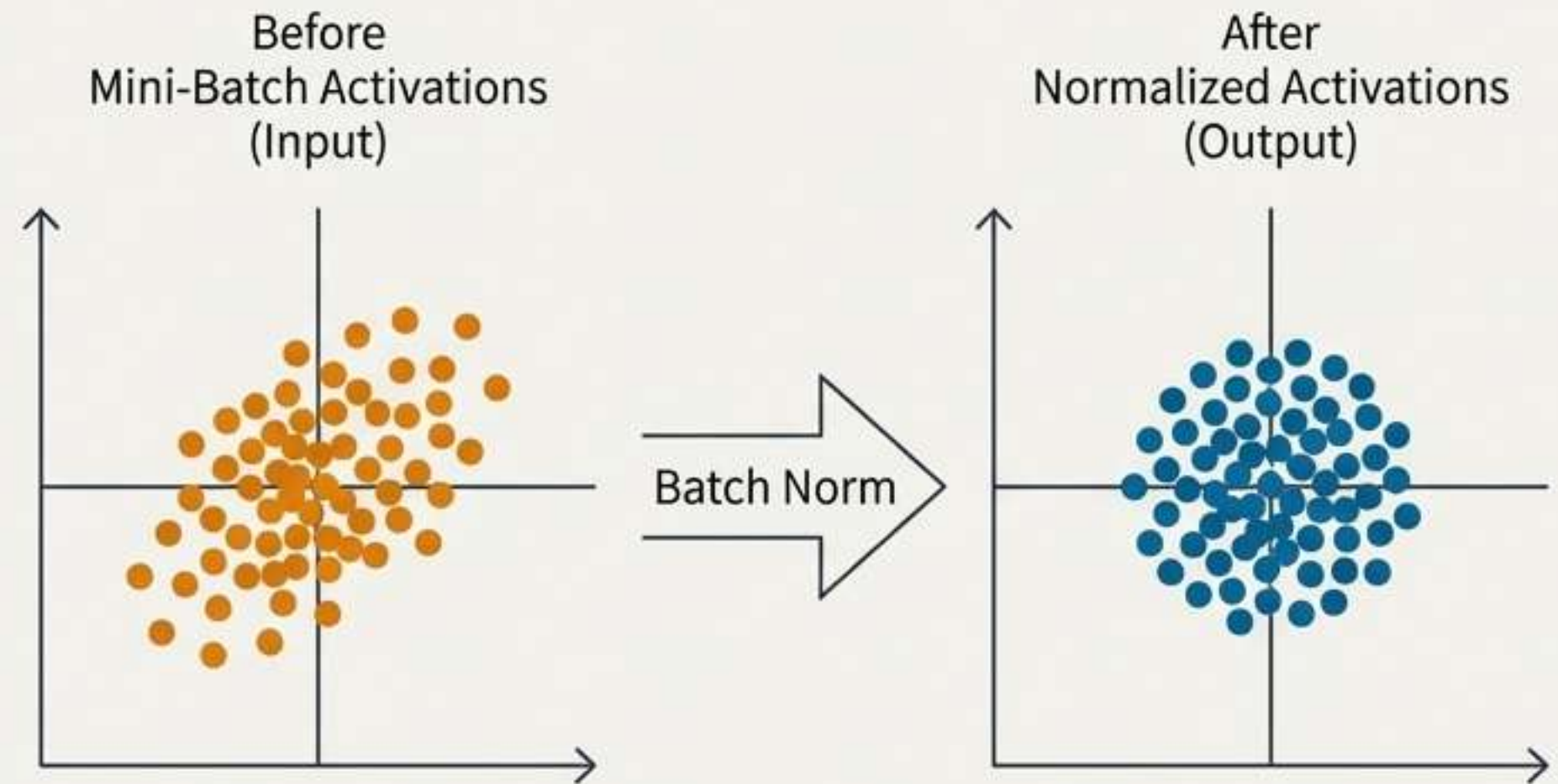
Normalizing activations across the mini-batch to stabilize learning.

## Core Concept

Batch Normalization re-centers and re-scales the inputs to each layer for every mini-batch. It aims to maintain a consistent distribution of activations throughout training.

## Key Benefit

By reducing the internal covariant shift, Batch Norm stabilizes the learning process, allowing for faster convergence and the use of higher learning rates.

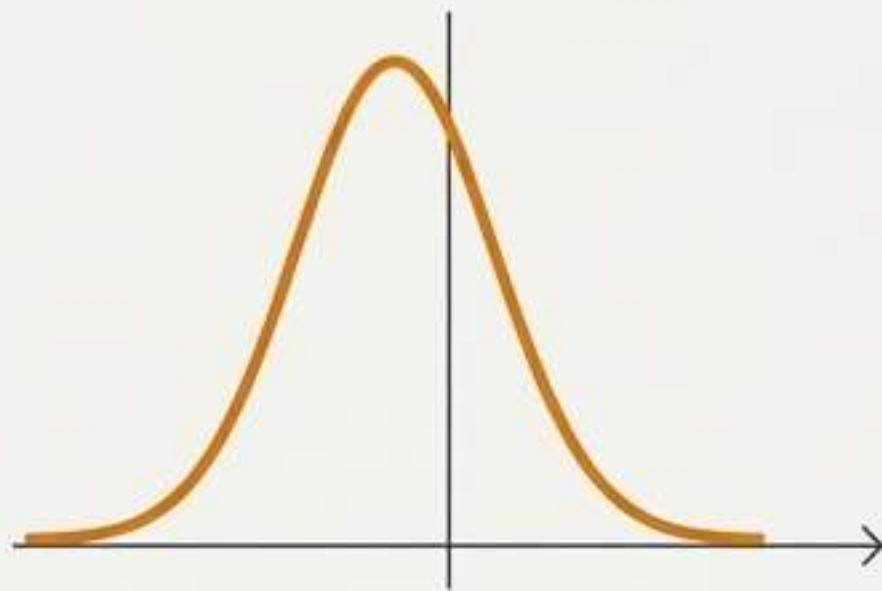




# How Batch Normalization Works, Step-by-Step

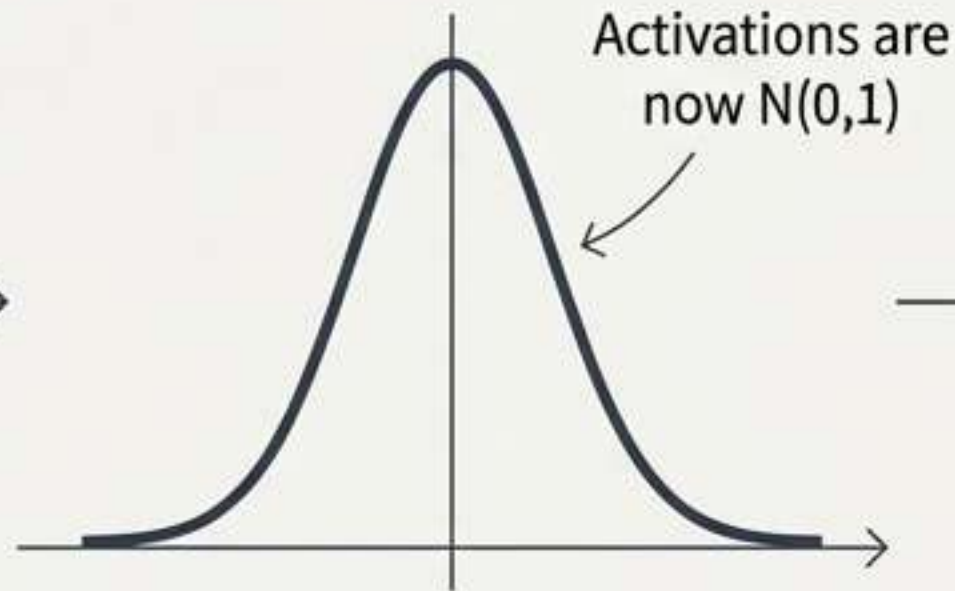
For a mini-batch of inputs to a layer:

## Step 1: Normalize

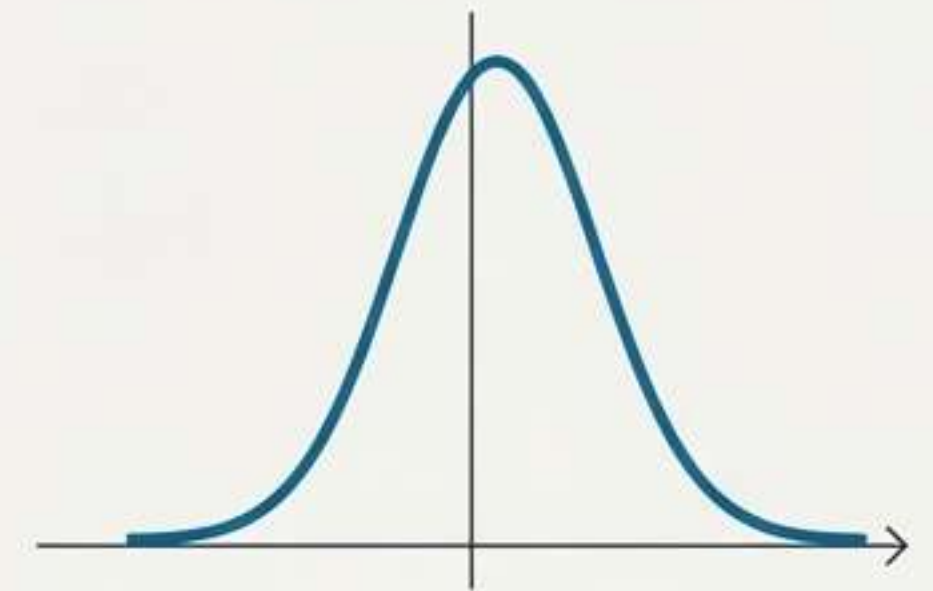


Calculate the mean ( $\mu$ ) and variance ( $\sigma^2$ ) of the activations across the mini-batch.

Normalize each activation:  $\frac{x - \mu}{\sqrt{\sigma^2 + \epsilon}}$ .



## Step 2: Scale and Shift



Introduce two learnable parameters per feature: a scale parameter ( $\gamma$ ) and a shift parameter ( $\beta$ ).

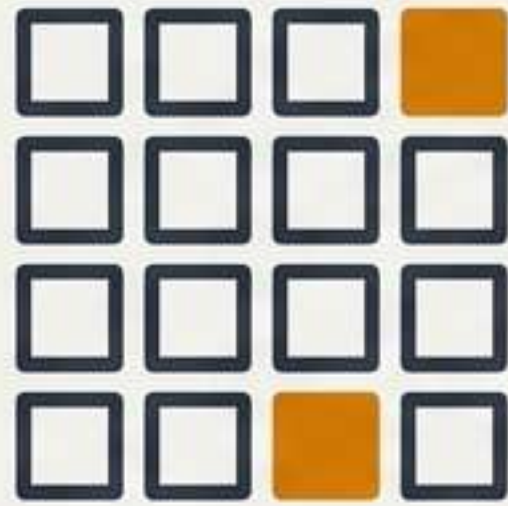
The final output is:  $\gamma * (\text{normalized\_x}) + \beta$

**Crucial Insight:** The Scale and Shift step is vital. It allows the network to learn the optimal distribution for each layer's activations, preventing the normalization from reducing the network's representational power. The network can learn to undo the normalization if needed.



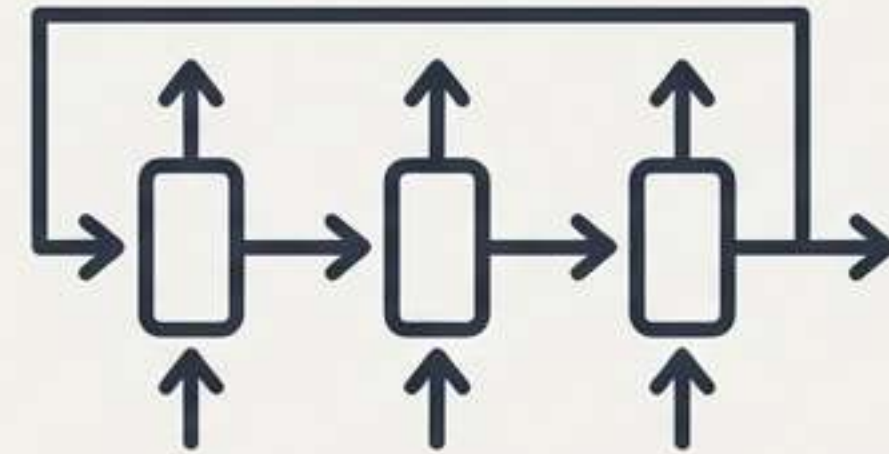
# Beyond the Batch: When Batch Norm Isn't the Answer

Batch Normalization's effectiveness is tied to batch statistics. This creates challenges in specific scenarios:



## Small Mini-Batch Sizes

The calculated batch mean and variance become noisy and unreliable, harming model performance. This is common in tasks like object detection or when using large models that require small batches due to memory constraints.



## Recurrent Neural Networks (RNNs)

Applying Batch Norm to RNNs is complex because the statistics need to be computed differently for each time step.

*How do we stabilize activations when batch-level statistics are unreliable or inapplicable?*



# Alternative I: Layer Normalization

Normalizing activations across the features within a single training example.

## How it Works

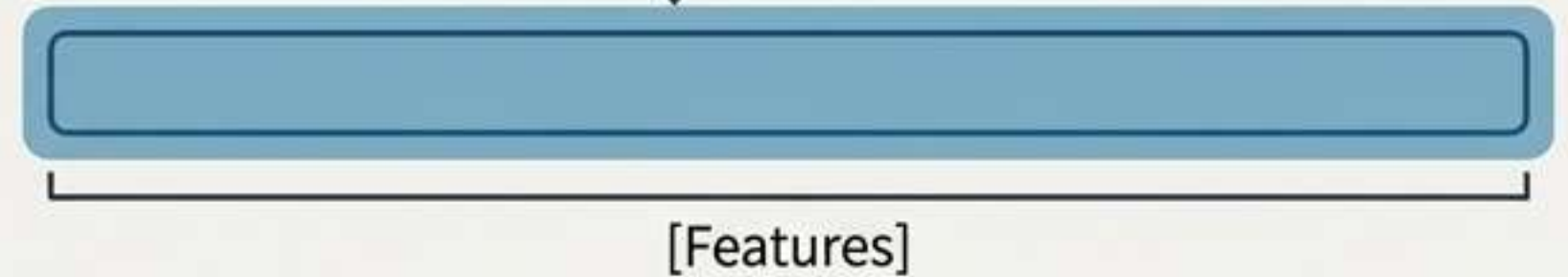
Layer Norm computes the mean and variance used for normalization from all of the summed inputs to the neurons in a layer for a single data sample. It is independent of independent of the batch.

## Key Advantage

Its performance is not dependent on the mini-batch size, making it highly effective for RNNs and scenarios with small batches.

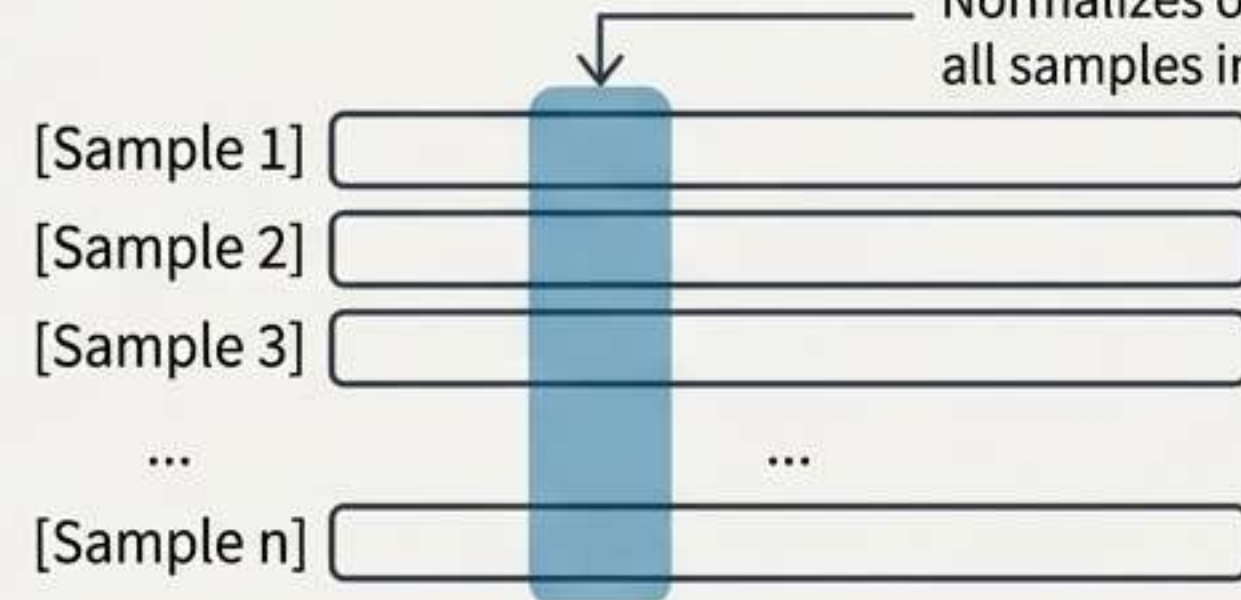
Layer Norm

Layer Norm: Normalizes across all features for one sample.



Batch Norm (for comparison)

Batch Norm (for comparison): Normalizes one feature across all samples in the batch.





# Alternative II: Group Normalization

A hybrid approach, normalizing within groups of channels.

## How it Works

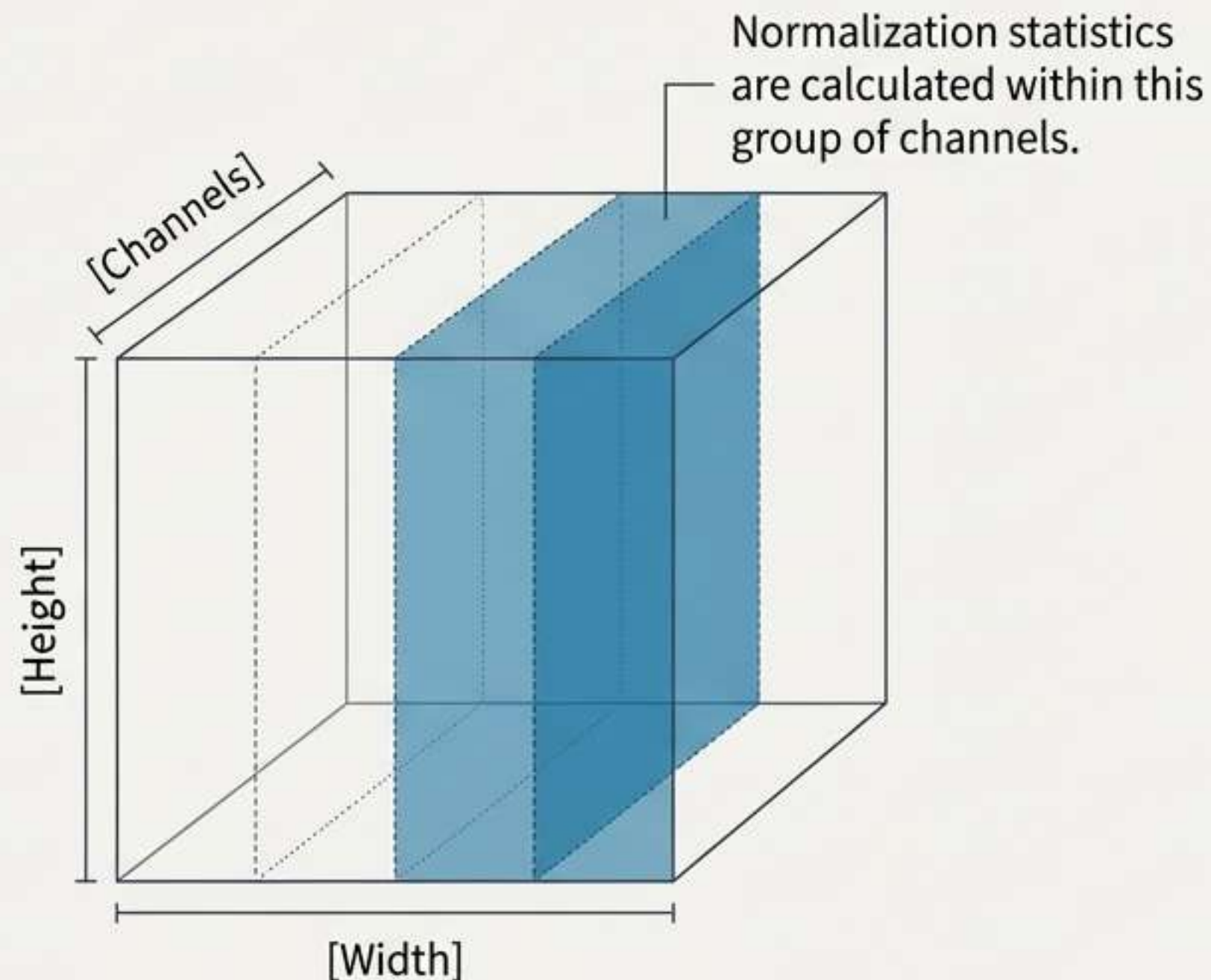
Group Normalization divides the channels into a number of groups and computes the mean and variance for normalization within each group for a single data sample.

## The Middle Ground

It sits between Layer Normalization (one big group) and Instance Normalization (one group per channel). Its performance is stable across a wide range of batch sizes.

## Key Advantage

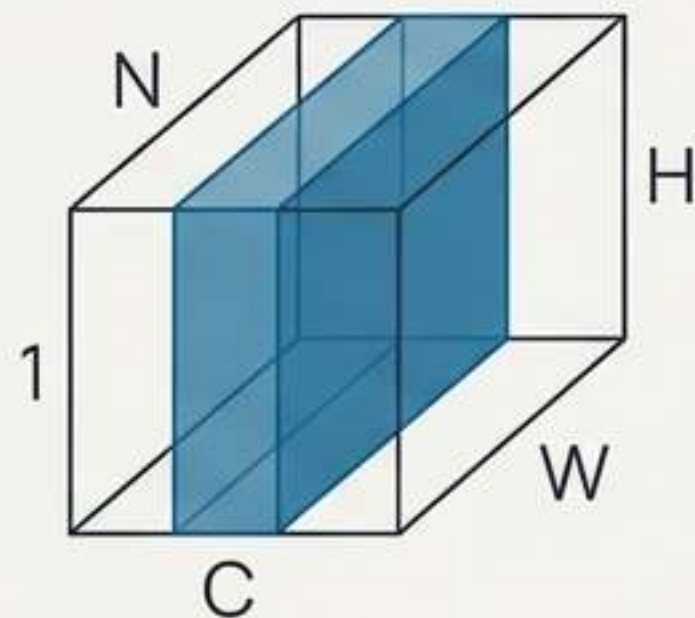
Offers a robust alternative to Batch Norm, especially in computer vision tasks where channel dependencies are important and batch sizes can vary.



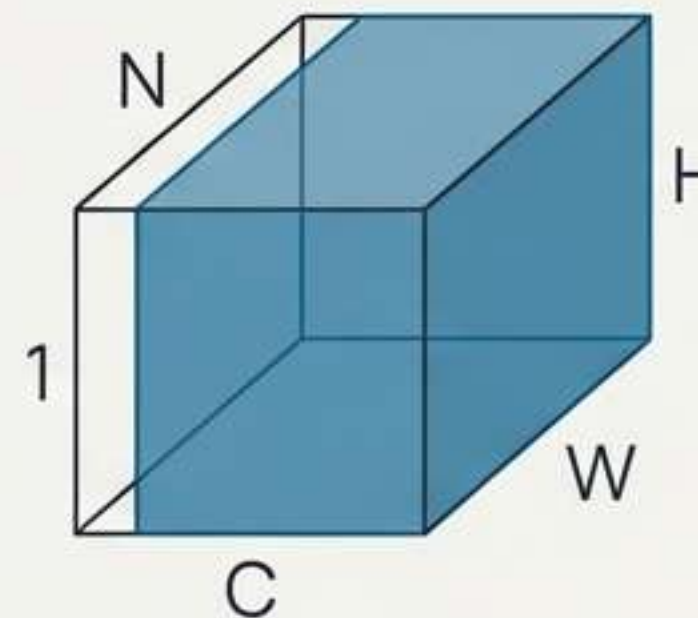


# A Visual Guide to Normalization Techniques

**Batch Norm**

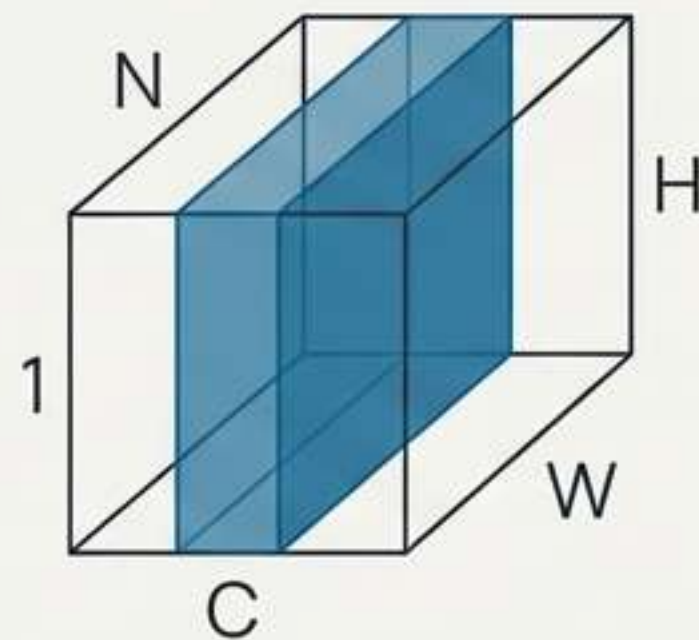


**Layer Norm**

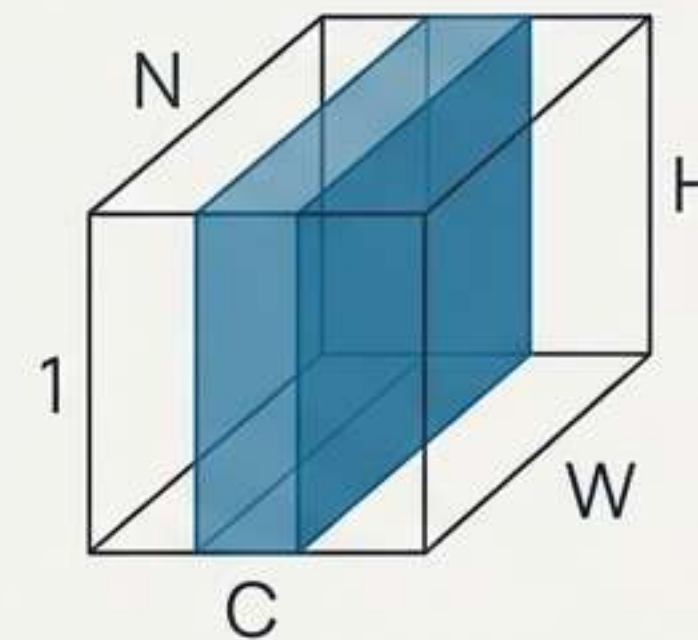


**Instance Norm**

(Common in style transfer)



**Group Norm**





# Complementary Strategies for a Stable Training Process

## Adaptive Learning Rates

### Techniques

Optimizers like Adam, AdaGrad, and RMSprop.



### Role in Mitigating ICS

They dynamically adjust the learning rate for each parameter based on past gradients. This makes the training process more resilient to the shifting distributions caused by ICS, though it doesn't solve the root problem.

## Data Augmentation

### Techniques

Applying transformations like rotations, flips, and color changes.



### Role in Mitigating ICS

By exposing the network to a wider range of input distributions during training, the model becomes more robust to the internal shifts that occur. It learns representations that are less sensitive to specific input statistics.



# Anchoring the Target, Stabilizing the Learning

## The Core Problem Recap

Internal Covariant Shift creates a 'moving target' for each layer in a deep network, slowing down and destabilizing the entire training process.



## The Core Solution Recap

Normalization techniques (Batch, Layer, Group) act as anchors. They control the distribution of activations, providing each layer with a more stable input.



## The Final Takeaway

Addressing internal covariant shift is not just an optimization—it is a critical step for building robust, efficient, and high-performing deep learning models. It enables faster training, better convergence, and more reliable results.