

Normalization

CoE197Z/EE298Z (Deep Learning)

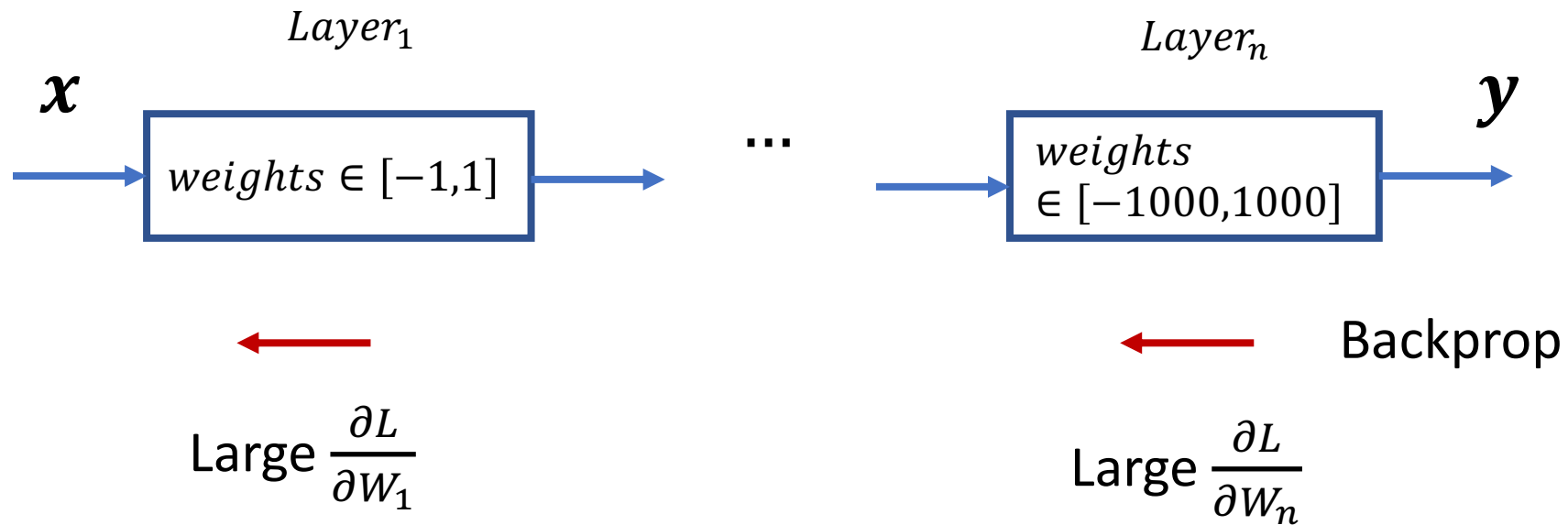
Rowel Atienza, Ph.D.

rowel@eee.upd.edu.ph

github.com/roatienza

Basic Idea

What if you have layers with different ranges of parameters?



End result: Unstable training

Notations

Batch size: N

Feature Map Width: W

Feature Map Height: H

Feature Map Channels: C

Feature Map: \mathbf{x}_{nchw}

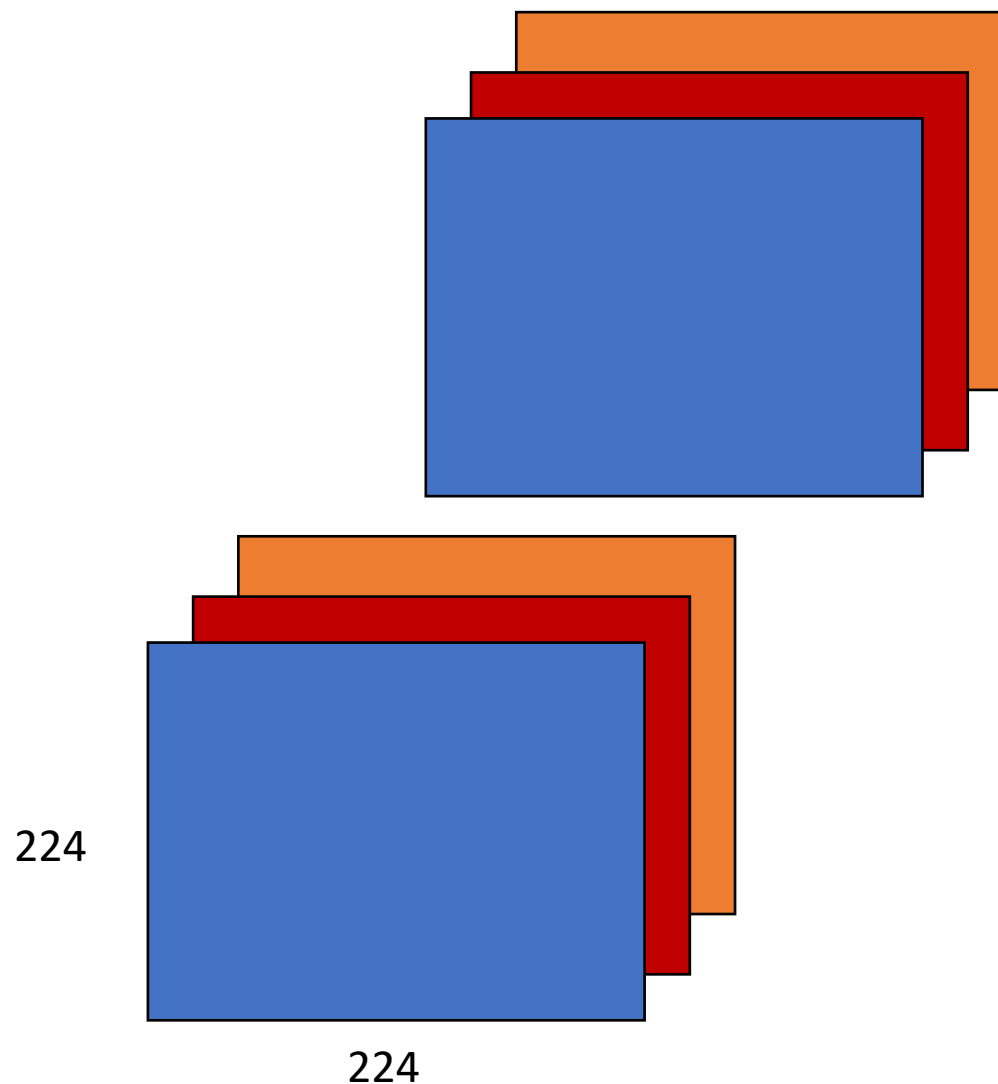
Example: \mathbf{x}_{nchw}

$$n = 1 \dots N = 2,$$

$$c = 1 \dots C = 3,$$

$$c = 1 \dots W = 224,$$

$$h = 1 \dots H = 224$$



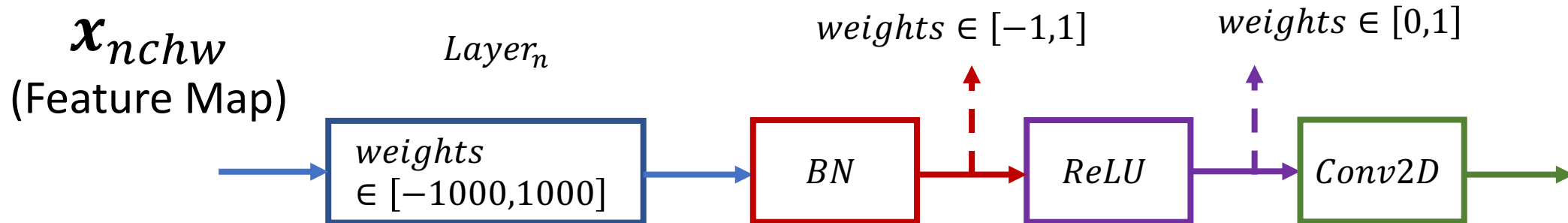
Batch Normalization (BN)

Ioffe, Sergey, and Christian Szegedy. "Batch normalization: Accelerating deep network training by reducing internal covariate shift." *arXiv preprint arXiv:1502.03167* (2015).

Batch Normalization (BN)

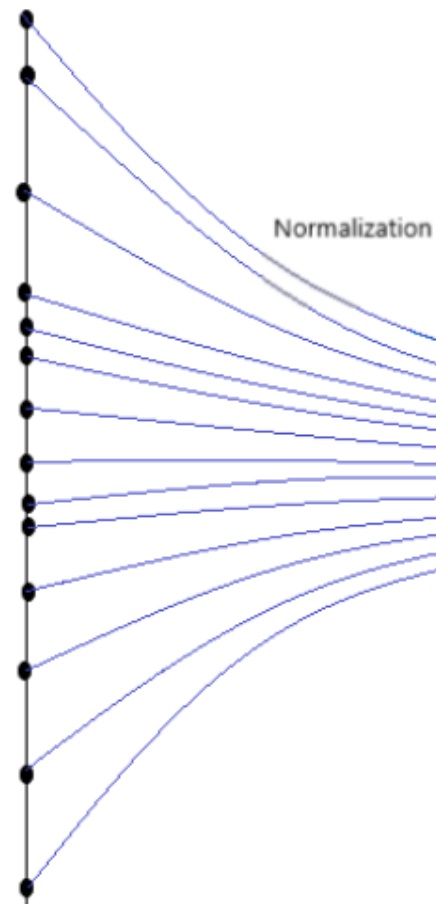
Batch Normalization normalizes the mean and standard deviation for each individual feature channel/map

Useful in deep networks (eg ResNet)

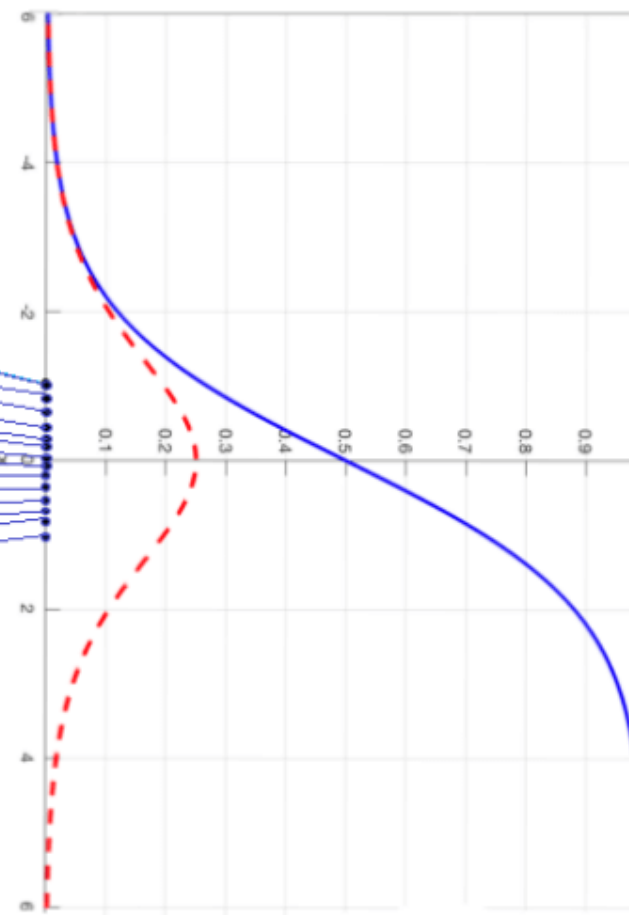


BN is a
mapping

Activation Inputs



Sigmoid Activation and Gradient



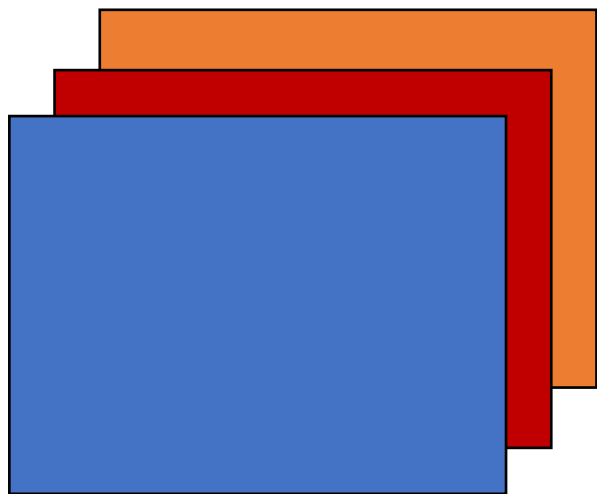
BN: The statistics are computed across the batch and the spatial dims

$$BN(\mathbf{x}_{nchw}) = \left(\frac{\mathbf{x}_{nchw} - \mu_c(\mathbf{x}_{nchw})}{\sigma_c(\mathbf{x}_{nchw})} \right)$$

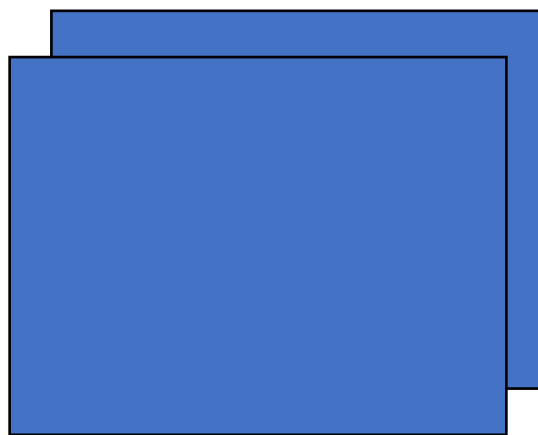
$$\mu_c(\mathbf{x}_{nchw}) = \frac{1}{NHW} \sum_{n=1}^N \sum_{h=1}^H \sum_{w=1}^W \mathbf{x}_{nchw}$$

$$\sigma_c(\mathbf{x}_{nchw}) = \sqrt{\frac{1}{NHW} \sum_{n=1}^N \sum_{h=1}^H \sum_{w=1}^W (\mathbf{x}_{nchw} - \mu_c(\mathbf{x}_{nchw}))^2}$$

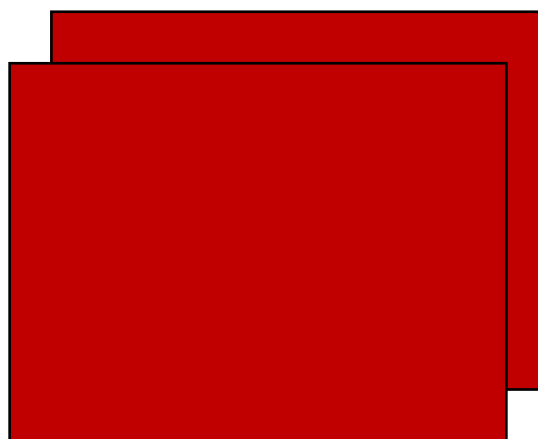
$\mathbf{x} = \mathbf{x}_{nchw}$
(Feature Map)



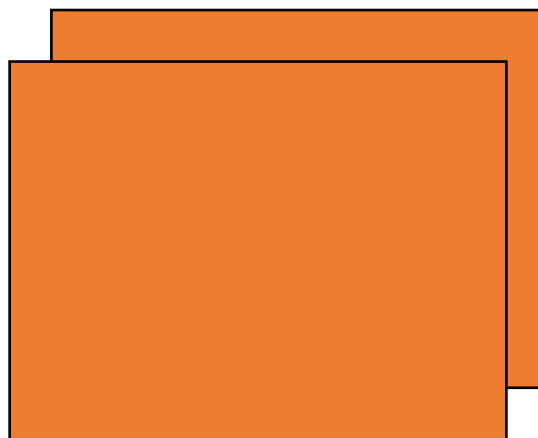
BN: sort
by
channel



$\mu_1(\mathbf{x}), \sigma_1(\mathbf{x}) \in \mathbb{R}$



$\mu_2(\mathbf{x}), \sigma_2(\mathbf{x}) \in \mathbb{R}$



$\mu_3(\mathbf{x}), \sigma_3(\mathbf{x}) \in \mathbb{R}$

1 per channel dim

Benefits of BN

BN accelerates the training of deep neural networks

Implicit regularization

Enhances gradient flow

Now ok to use of saturating non-linear functions (eg sigmoid or tanh)

Danger of BN

Small batch sizes leads to inaccurate batch statistics

Problematic with variable batch size



How Does Batch Normalization Help Optimization?

Shibani Santurkar*, Dimitris Tsipras*, Andrew Ilyas*, Aleksander Madry



Layer Normalization (LN)

Ba, Jimmy Lei, Jamie Ryan Kiros, and Geoffrey E. Hinton. "Layer normalization."
arXiv preprint arXiv:1607.06450 (2016).

Layer Normalization

The statistics are computed across all channels and spatial dims

The statistics are independent of the batch

Useful in transformers

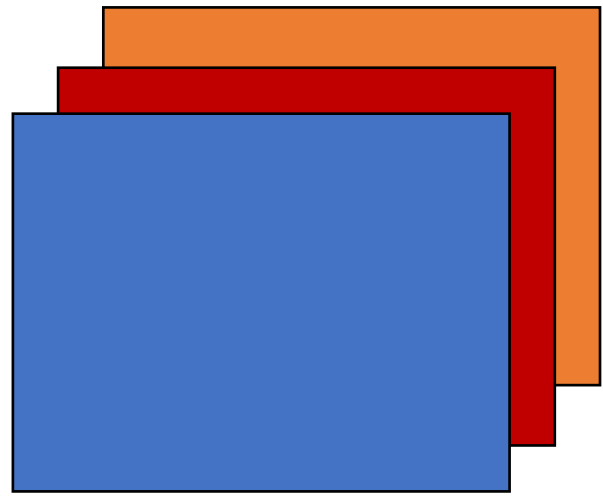
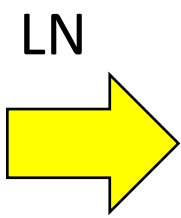
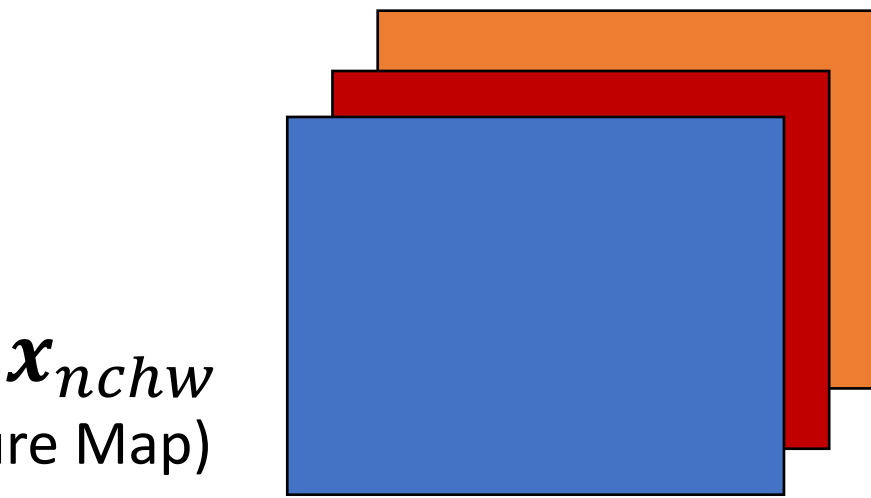
LN: The statistics are computed across all channels and the spatial dims

$$LN(\mathbf{x}_{nchw}) = \left(\frac{\mathbf{x}_{nchw} - \mu_n(\mathbf{x}_{nchw})}{\sigma_n(\mathbf{x}_{nchw})} \right)$$

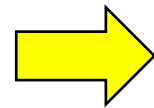
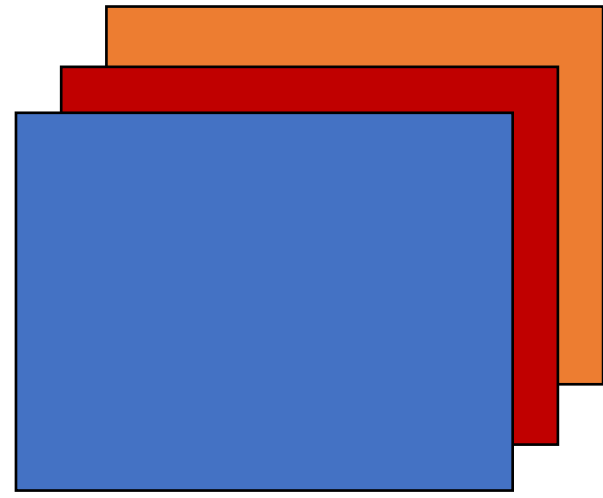
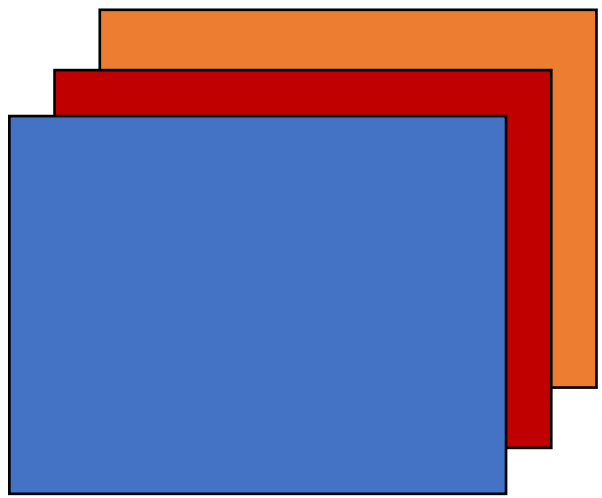
$$\mu_n(\mathbf{x}_{nchw}) = \frac{1}{CHW} \sum_{c=1}^C \sum_{h=1}^H \sum_{w=1}^W \mathbf{x}_{nchw}$$

$$\sigma_n(\mathbf{x}_{nchw}) = \sqrt{\frac{1}{CHW} \sum_{c=1}^C \sum_{h=1}^H \sum_{w=1}^W (\mathbf{x}_{nchw} - \mu_n(\mathbf{x}_{nchw}))^2}$$

$\mathbf{x} = \mathbf{x}_{nchw}$
(Feature Map)



$\mu_1(\mathbf{x}), \sigma_1(\mathbf{x}) \in \mathbb{R}$

A yellow arrow pointing from the first output feature map stack to the second output feature map stack.

$\mu_2(\mathbf{x}), \sigma_2(\mathbf{x}) \in \mathbb{R}$

1 per batch dim

Instance Normalization (IN)

Instance Normalization (IN)

Instance Normalization (IN) is computed only across the features' spatial dimensions

It is independent for each channel and sample

Useful in style transfer

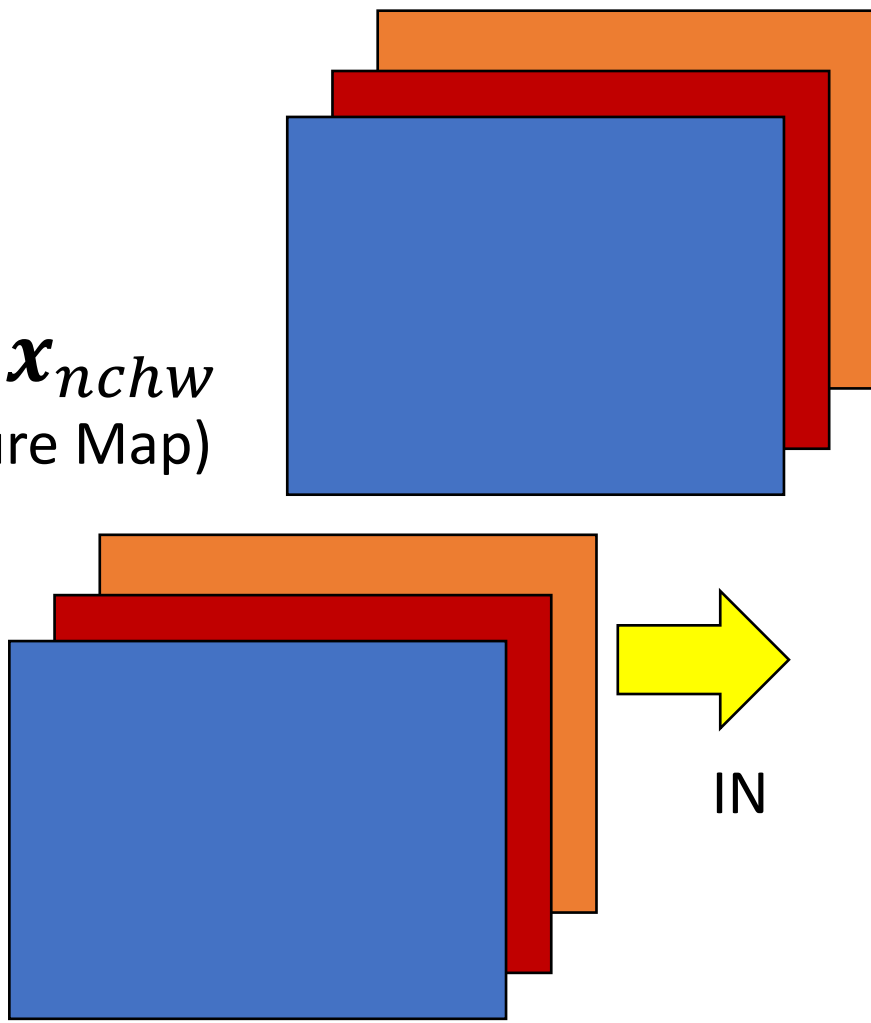
IN: The statistics are computed across spatial dims of each feature map

$$IN(\mathbf{x}_{nchw}) = \left(\frac{\mathbf{x}_{nchw} - \mu_{nc}(\mathbf{x}_{nchw})}{\sigma_{nc}(\mathbf{x}_{nchw})} \right)$$

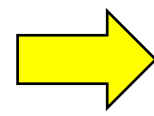
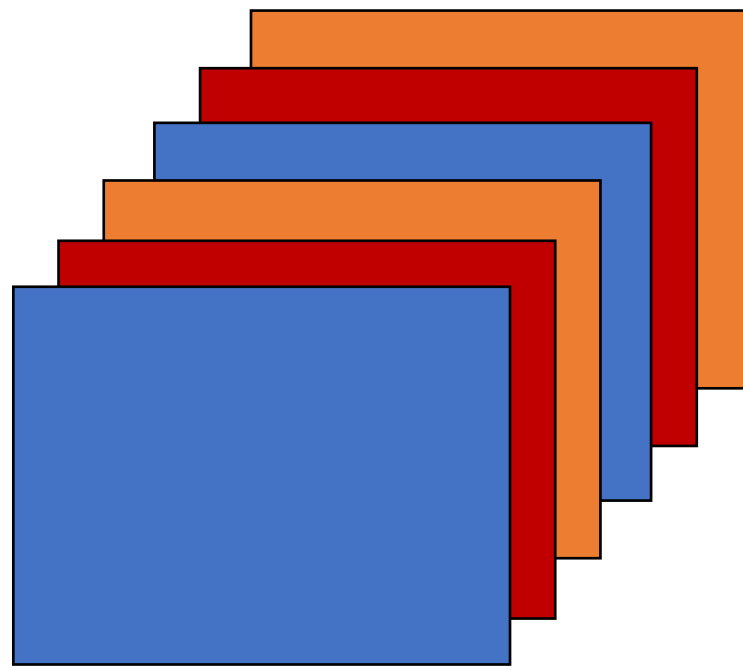
$$\mu_{nc}(\mathbf{x}_{nchw}) = \frac{1}{HW} \sum_{h=1}^H \sum_{w=1}^W \mathbf{x}_{nchw}$$

$$\sigma_{nc}(\mathbf{x}_{nchw}) = \sqrt{\frac{1}{HW} \sum_{h=1}^H \sum_{w=1}^W (\mathbf{x}_{nchw} - \mu_{nc}(\mathbf{x}_{nchw}))^2}$$

$\mathbf{x} = \mathbf{x}_{nchw}$
(Feature Map)



IN



$$\mu_1(\mathbf{x}), \sigma_1(\mathbf{x}) \in \mathbb{R}$$

$$\mu_2(\mathbf{x}), \sigma_2(\mathbf{x}) \in \mathbb{R}$$

...

$$\mu_6(\mathbf{x}), \sigma_6(\mathbf{x}) \in \mathbb{R}$$

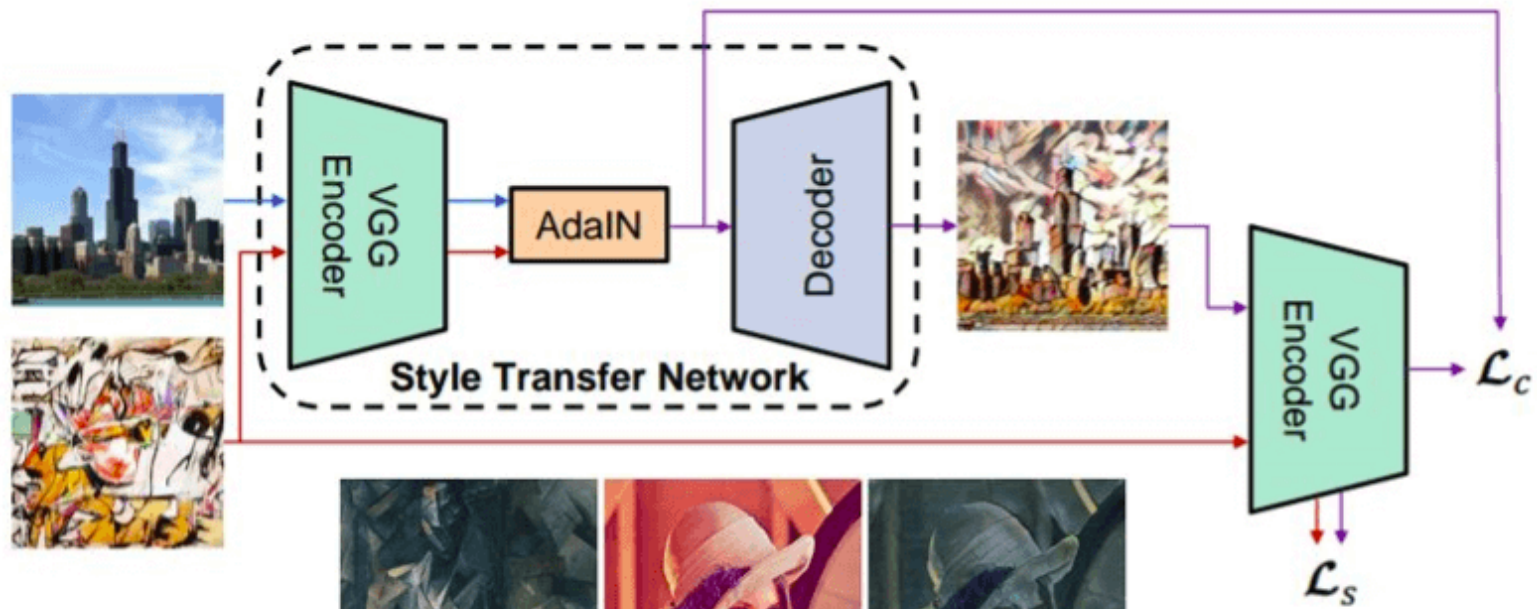
1 per instance

Adaptive Instance Normalization (AIN)

Adaptive Instance Normalization (AIN)

Adaptive Instance Normalization (AdaIN) receives an input image \mathbf{x} (content) and a style input \mathbf{y} , and simply aligns the channel-wise mean and variance of \mathbf{x} to match those of \mathbf{y} .

$$AIN(\mathbf{x}_{nchw}) = \sigma_{nc}(\mathbf{y}_{nchw}) \left(\frac{\mathbf{x}_{nchw} - \mu_{nc}(\mathbf{x}_{nchw})}{\sigma_{nc}(\mathbf{x}_{nchw})} \right) + \mu_{nc}(\mathbf{y}_{nchw})$$



Style

Content

Ours

Group Normalization

Group Normalization

Group normalization (GN) divides the channels into groups and computes the first-order statistics within each group.

Independent of batch sizes

Accuracy is more stable than BN in a wide range of batch sizes

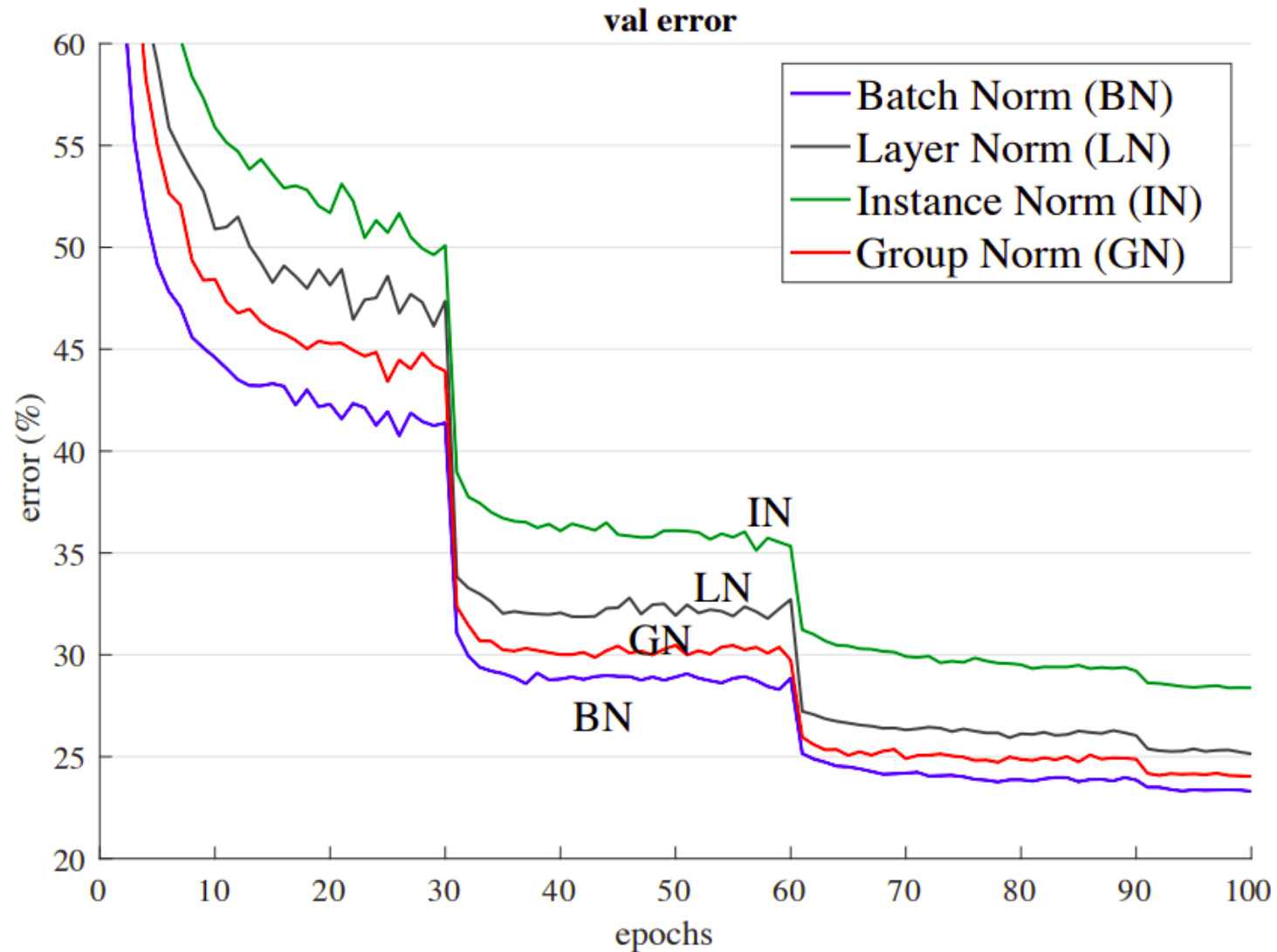
Groups = number of channels : Equivalent to instance normalization

Groups = 1: Equivalent to layer normalization

Group Normalization

ECCV 2018, Munich

Yuxin Wu, Kaiming He
Facebook AI Research (FAIR)



Comparison using a batch size of 32 images per GPU in ImageNet. Validation error VS the numbers of training epochs is shown. The model is ResNet-50. Source: [Group Normalization](#)

In Summary: Different problems require different normalizations

<https://youtu.be/UfffxcCQMPQ>

