# Revisiting Batch Norm Initialization
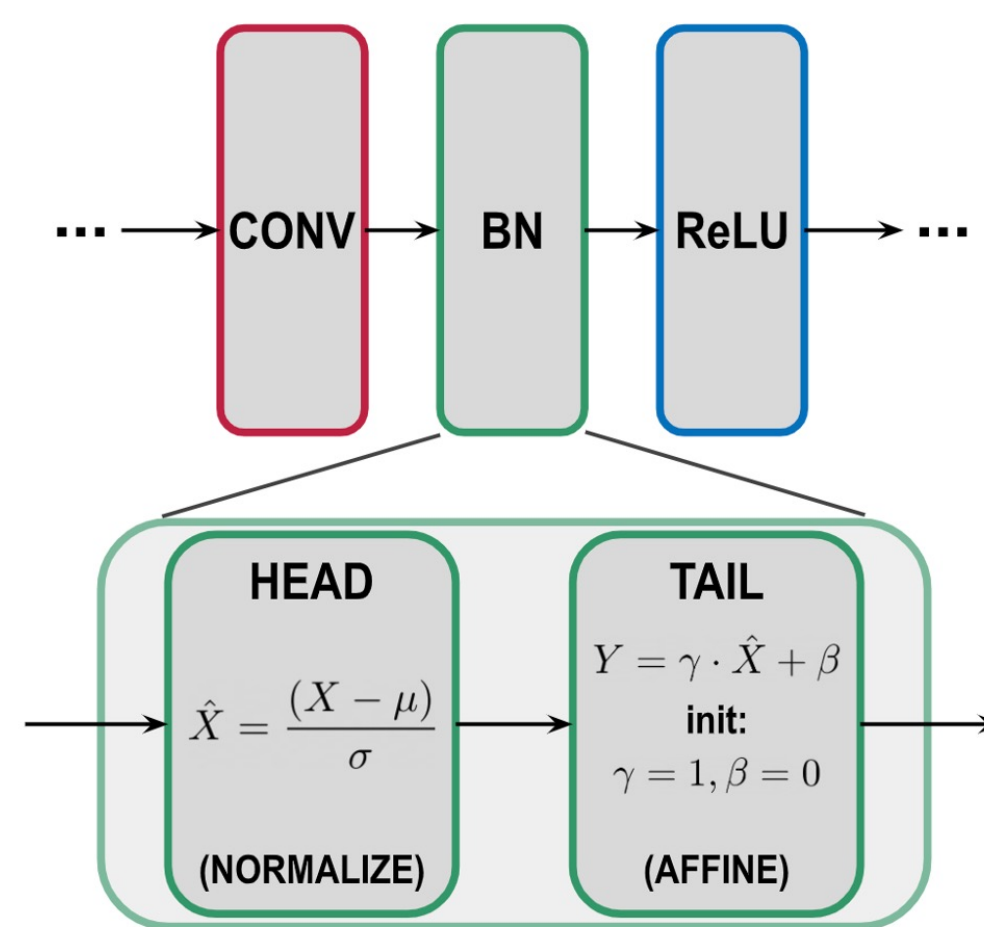
Jim Davis, **Logan Frank**

Department of Computer Science and Engineering, Ohio State University

## 1. Batch Normalization (BN) in Deep Neural Networks

**Two per-channel operations**
- **Head:** Normalizes data
- **Tail:** Learnable affine transformation

Constrains intermediate features, enabling smoother and faster optimization, and stochasticity of batch statistics can benefit generalization



## 2. BN: Forward Formulation

- Compute mean ( $\mu$ ) and variance ( $\sigma^2$ ) across batch dimension
- Use computed statistics to normalize the data ( $\mu = 0$, $\sigma^2 = 1$ )
- Apply an affine transformation to the normalized data using learnable parameters: scale ( $\gamma$ ) and shift ( $\beta$ )

$$\mu_B = \frac{1}{m}\sum_{i=1}^{m} x_i$$

$$\sigma_B^2 = \frac{1}{m}\sum_{i=1}^{m}(x_i - \mu_B)^2$$

$$\hat{X} = \frac{X - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} \longrightarrow Y = \gamma \cdot \hat{X} + \beta$$
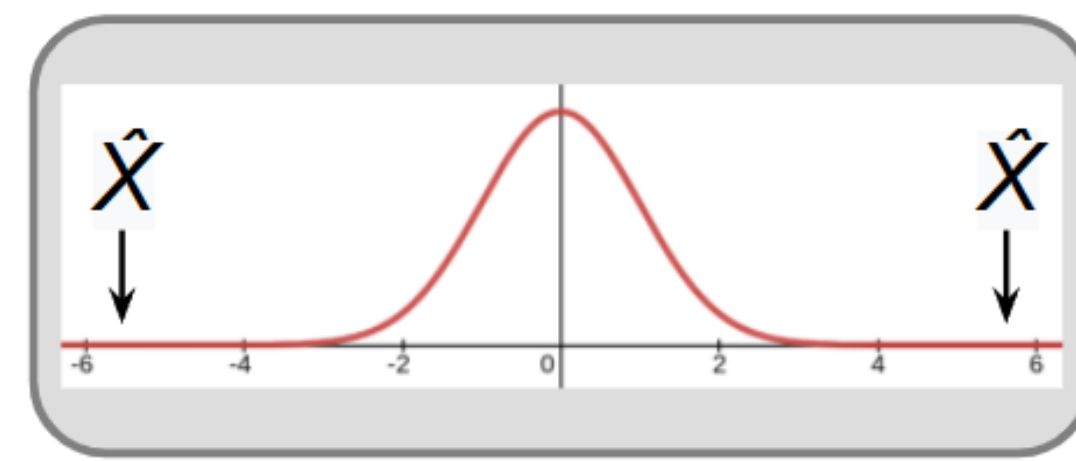
## 3. Observed Issues with BN

Learnable BN parameters are initialized to $\gamma = 1$ and $\beta = 0$ (identity function)

We observed that the final learned parameter values tend to remain close to their initialization

Furthermore, we observed that the BN normalization head can yield overly large values ( $\pm 6\sigma$ ) for the proceeding layer, which can be undesirable for training

$$
\gamma_{\text{init}} \qquad \gamma_{\text{learned}}
$$
$$
\begin{bmatrix}[1.0],\\ [1.0],\\ [1.0],\\ [1.0]\end{bmatrix} \xrightarrow{\text{TRAIN}} \begin{bmatrix}[0.99],\\ [1.03],\\ [1.17],\\ [0.95]\end{bmatrix}
$$



## 4. Proposed Adjustments to BN Scale Parameter $\gamma$

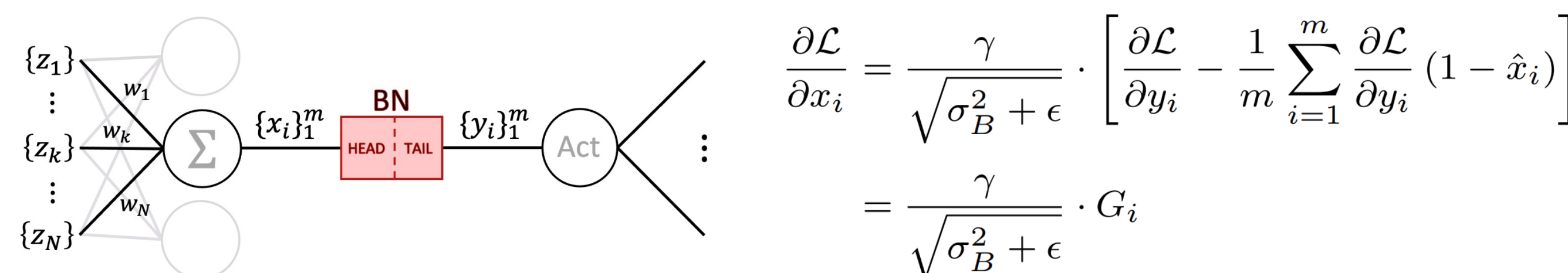**Initialize $\gamma$ to a value in $(0, 1]$**
- Directly addresses overly large values after normalization by immediately scaling down the data (with *no additional parameters*)
- Enables BN shift parameter $\beta$ to have a broader reach on scaled data before the proceeding activation function (in many cases ReLU)

**Reduce the learning rate $\alpha$ on $\gamma$**
- Divide learning rate on $\gamma$ by constant $c$ ( $\alpha_\gamma = \alpha/c$ )
- Allows for fine-grained search near initialization value
- Leave $\beta$ with original learning rate, enabling it to have a broader and now more stable search of the normalized and scaled data

## 5. BN: Gradients and Insights

Using a fully-connected layer as illustration (below, left), we use the gradients given in the original BN paper to derive the gradient of the loss with respect to the input $\partial\mathcal{L}/\partial x_i$ (below, right)



$$\frac{\partial\mathcal{L}}{\partial x_i} = \frac{\gamma}{\sqrt{\sigma_B^2 + \epsilon}} \cdot \left[\frac{\partial\mathcal{L}}{\partial y_i} - \frac{1}{m}\sum_{i=1}^{m}\frac{\partial\mathcal{L}}{\partial y_i}(1 - \hat{x}_i)\right]$$

$$= \frac{\gamma}{\sqrt{\sigma_B^2 + \epsilon}} \cdot G_i$$

**Insights**
- No effects introduced by $\gamma < 1$ for the first backward pass
- *Negligible* effects for remainder of training
- More insights shown in the paper

$$\sigma_B^2 = \sigma_{act}^2 \cdot \gamma_{prev}^2 \sum_{k=1}^{N} w_k^2 \longrightarrow \frac{\gamma_{curr}}{\gamma_{prev} \cdot \sqrt{\sigma_{act}^2 \sum w_k^2 + \epsilon}} = \frac{1}{\sqrt{\sigma_{act}^2 \sum w_k^2 + \epsilon}}$$

## 6. Training Details

**BN scale initialization:** $\gamma \in \{0.01, 0.05, 0.1, 0.25, 0.5, 0.75, 1.0\}$
- Examine subset of values after initial CIFAR-10 experiments

**BN scale learning rate reduction factor:** $c = 100$

## 7. Statistical Significance

Different RNG seeds can cause variations in final score (accuracy)

For **each** experiment, we conduct 15 runs with different seeds, aggregate the results of each run (to report a mean and standard deviation), and compare to a baseline (or related approach) using a one-sided paired t-test (using a p-value of 0.05)

## 8. Results

- Significant improvements across multiple initial values of $\gamma$ and learning rates for CIFAR-10 (T1.a), as well as CIFAR-100, CUB-200, and Stanford Cars (T2.b)

- Even greater gains with deeper network architectures (T2.a)

- Outperforms other existing related approaches which require additional parameters and computations (T2.b)

| $\gamma$ | Learning Rate ($\alpha$) | | |
|---|---|---|---|
| | 0.1 | 0.01 | 0.001 |
| 0.01 | $85.50_{\pm0.39}$ | $87.11_{\pm0.23}$ | $\underline{80.37}_{\pm0.58}$ |
| 0.05 | $90.19_{\pm0.32}$ | $88.84_{\pm0.32}$ | $76.98_{\pm0.71}$ |
| 0.10 | $\mathbf{90.80}_{\pm0.20}$ | $87.31_{\pm0.37}$ | $74.48_{\pm0.55}$ |
| 0.25 | $90.32_{\pm0.24}$ | $85.33_{\pm0.43}$ | $73.83_{\pm0.64}$ |
| 0.50 | $90.17_{\pm0.19}$ | $84.60_{\pm0.35}$ | $72.80_{\pm0.68}$ |
| 0.75 | $90.19_{\pm0.18}$ | $84.43_{\pm0.30}$ | $72.01_{\pm0.58}$ |
| 1.00 | $89.81_{\pm0.46}$ | $84.48_{\pm0.33}$ | $71.15_{\pm0.56}$ |
| BASE | $89.44_{\pm0.45}$ | $84.64_{\pm0.25}$ | $71.32_{\pm0.60}$ |

Table 1.a

| Dataset | Method | Learning Rate ($\alpha$) | |
|---|---|---|---|
| | | 0.1 | 0.01 |
| CIFAR10 | RBN | $90.17_{\pm0.22}$ | $84.72_{\pm0.29}$ |
| | RBN⁻ | $90.11_{\pm0.24}$ | $84.50_{\pm0.36}$ |
| | IEBN | $90.18_{\pm0.26}$ | $85.34_{\pm0.39}$ |
| | IEBN⁻ | $90.15_{\pm0.24}$ | $85.29_{\pm0.35}$ |
| | Ours | $\mathbf{90.80}_{\pm0.20}$ | $\mathbf{88.84}_{\pm0.32}$ |
| | BASE | $89.44_{\pm0.45}$ | $84.64_{\pm0.25}$ |
| CIFAR100 | RBN | $66.95_{\pm0.57}$ | $58.95_{\pm0.42}$ |
| | RBN⁻ | $66.82_{\pm0.55}$ | $58.90_{\pm0.61}$ |
| | IEBN | $66.94_{\pm0.39}$ | $60.61_{\pm0.40}$ |
| | IEBN⁻ | $66.95_{\pm0.32}$ | $60.89_{\pm0.41}$ |
| | Ours | $\mathbf{68.80}_{\pm0.49}$ | $\mathbf{64.01}_{\pm0.54}$ |
| | BASE | $66.01_{\pm0.95}$ | $58.48_{\pm0.53}$ |
| CUB-200 | RBN | $48.68_{\pm1.56}$ | $44.68_{\pm0.59}$ |
| | RBN⁻ | $47.14_{\pm2.72}$ | $43.02_{\pm1.22}$ |
| | IEBN | $54.12_{\pm0.60}$ | $44.92_{\pm0.74}$ |
| | IEBN⁻ | $53.81_{\pm0.76}$ | $44.09_{\pm0.65}$ |
| | Ours | $\mathbf{58.52}_{\pm0.69}$ | $\mathbf{45.31}_{\pm0.59}$ |
| | BASE | $46.26_{\pm1.59}$ | $41.61_{\pm1.03}$ |
| ST-Cars | RBN | $68.17_{\pm1.84}$ | $51.87_{\pm1.34}$ |
| | RBN⁻ | $67.84_{\pm2.96}$ | $52.30_{\pm1.73}$ |
| | IEBN | $73.60_{\pm0.92}$ | $51.06_{\pm0.87}$ |
| | IEBN⁻ | $74.04_{\pm1.55}$ | $51.08_{\pm0.78}$ |
| | Ours | $\mathbf{78.29}_{\pm0.44}$ | $51.18_{\pm2.16}$ |
| | BASE | $64.73_{\pm2.87}$ | $51.86_{\pm1.80}$ |

Table 2.b

| Network | $\gamma$ | Learning Rate ($\alpha$) | |
|---|---|---|---|
| | | 0.1 | 0.01 |
| ResNet-50 | 0.05 | $91.23_{\pm0.20}$ | $\underline{89.60}_{\pm0.19}$ |
| | 0.10 | $\mathbf{91.28}_{\pm0.26}$ | $87.67_{\pm0.20}$ |
| | 0.50 | $89.49_{\pm0.27}$ | $84.74_{\pm0.37}$ |
| | BASE | $86.94_{\pm1.23}$ | $85.04_{\pm0.32}$ |
| ResNet-101 | 0.05 | $\mathbf{91.58}_{\pm0.22}$ | $\mathbf{90.02}_{\pm0.22}$ |
| | 0.10 | $91.26_{\pm0.18}$ | $88.35_{\pm0.28}$ |
| | 0.50 | $89.89_{\pm0.74}$ | $85.23_{\pm0.50}$ |
| | BASE | $88.28_{\pm1.39}$ | $84.74_{\pm0.56}$ |
| ResNet-152 | 0.05 | $91.20_{\pm0.16}$ | $\mathbf{90.00}_{\pm0.17}$ |
| | 0.10 | $90.89_{\pm0.41}$ | $88.31_{\pm0.33}$ |
| | 0.50 | $90.17_{\pm0.23}$ | $85.23_{\pm0.62}$ |
| | BASE | $88.73_{\pm0.62}$ | $84.15_{\pm0.79}$ |

Table 2.a

## 9. QR Codes:

Paper

GitHub