

Online Normalization for Training Neural Networks

Ilya Sharapov

Sofía Samaniego de la Fuente













Motivation and Background

- Normalization accelerates learning
- Normalization transforms a neural network from a function to a statistical operator that depends on its input distribution.
- Batches are commonly used to approximate the input distribution[1]
- We propose an online process that eliminates batches to:
- Decrease memory usage
- Compute unbiased gradients
- Provide training/inference symmetry
- Online Norm is compatible with automatic differentiation,
- Online | Batch Norm Network ResNet-50, ImageNet, theory ResNet-50, ImageNet, measured a 3D U-Net, $150 \times 150 \times 150$ voxels, theory 3D U-Net, $250 \times 250 \times 250$ voxels, theory 2D U-Net, 1024×1024 pixels, theory 2D U-Net, 2048×2048 pixels, theory

Ryan Reece

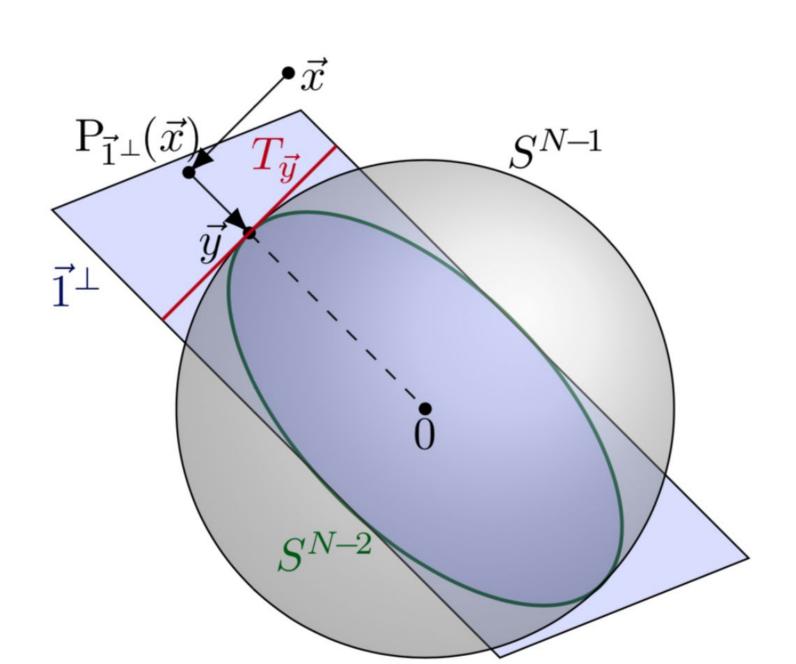
PyTorch implementation stores multiple copies of activations for improved performance.

so it can be integrated into any deep learning framework.

Principles of Normalization

 Normalization and its derivative are statistical operators \supset Notation: $(\cdot)' = \nabla_{(\cdot)} L$

$$y = f_{\mathbb{X}}[x] \equiv (x - \mu[x])/\sigma[x]$$
 and $x' = (\nabla_x f_{\mathbb{X}}[x])y', \quad x \sim \mathbb{X}$

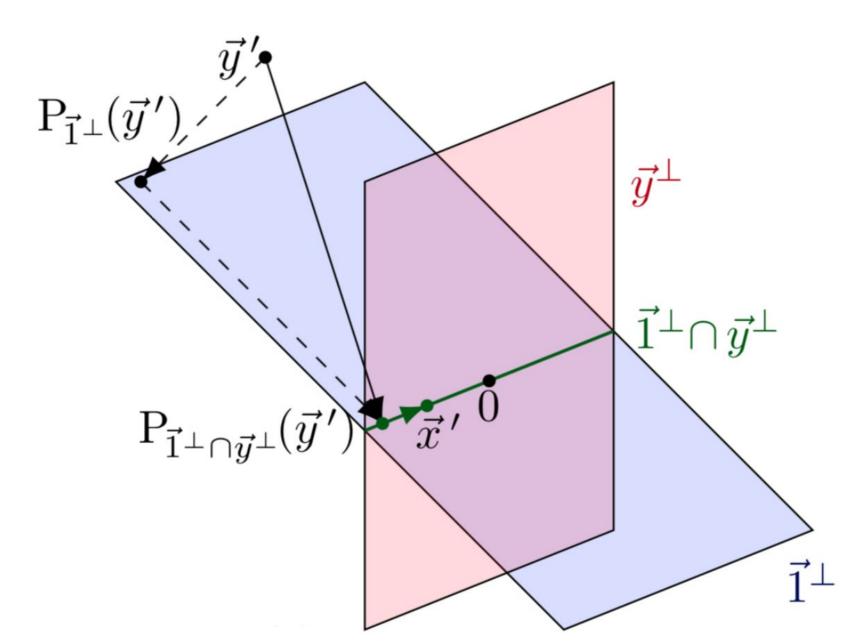


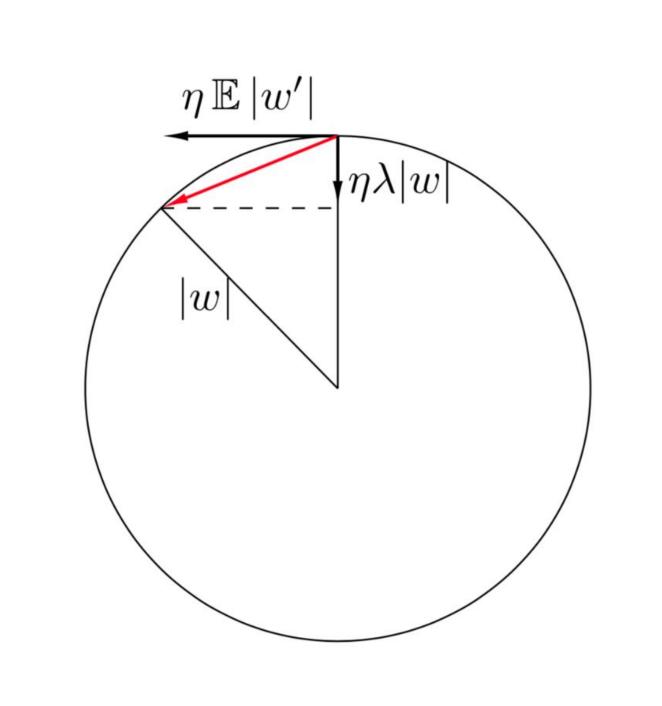
FWD: Projection and rescaling

s.t. normalized output lies in

the zero-centered unit sphere:

 $\mu[y] = 0 \quad \text{and} \quad \mu[y^2] = 1$





BWD: Two projections [2] st. gradient satisfies the orthogonality conditions that follow from the backward eqs: $\mu[x'] = 0 \quad \text{and} \quad \mu[x'y] = 0$

Normalized networks are invariant to gradient scale

Weight decay is required to prevent magnitude growth

Online Normalization

Vitaliy Chiley

Forward Pass: Tracking process

- Exponential moving average
- Maintains mean zero, unit variance in expectation

Layer-wide correction

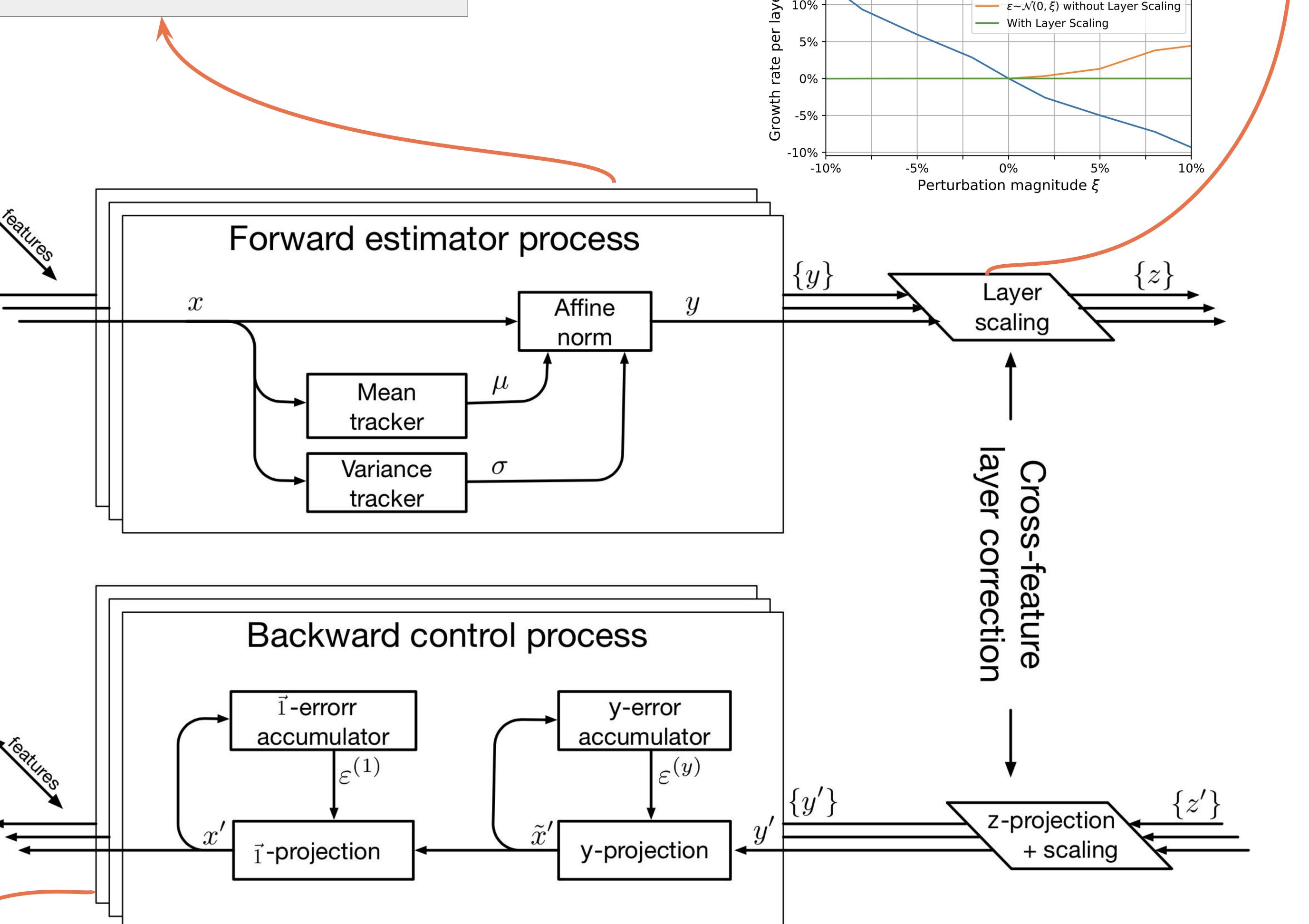
Vishal Subbiah

Urs Köster

Michael James

Atli Kosson

- Cross-Feature second moment correction
- Avoids exponential growth of activations



Backward Pass: Control Process

 Basic controller to enforce backward pass orthogonality constraints:

$$\langle x', 1 \rangle = 0$$
 and $\langle x', y \rangle = 0$

 Leads to uniformly bounded error in gradient calculation

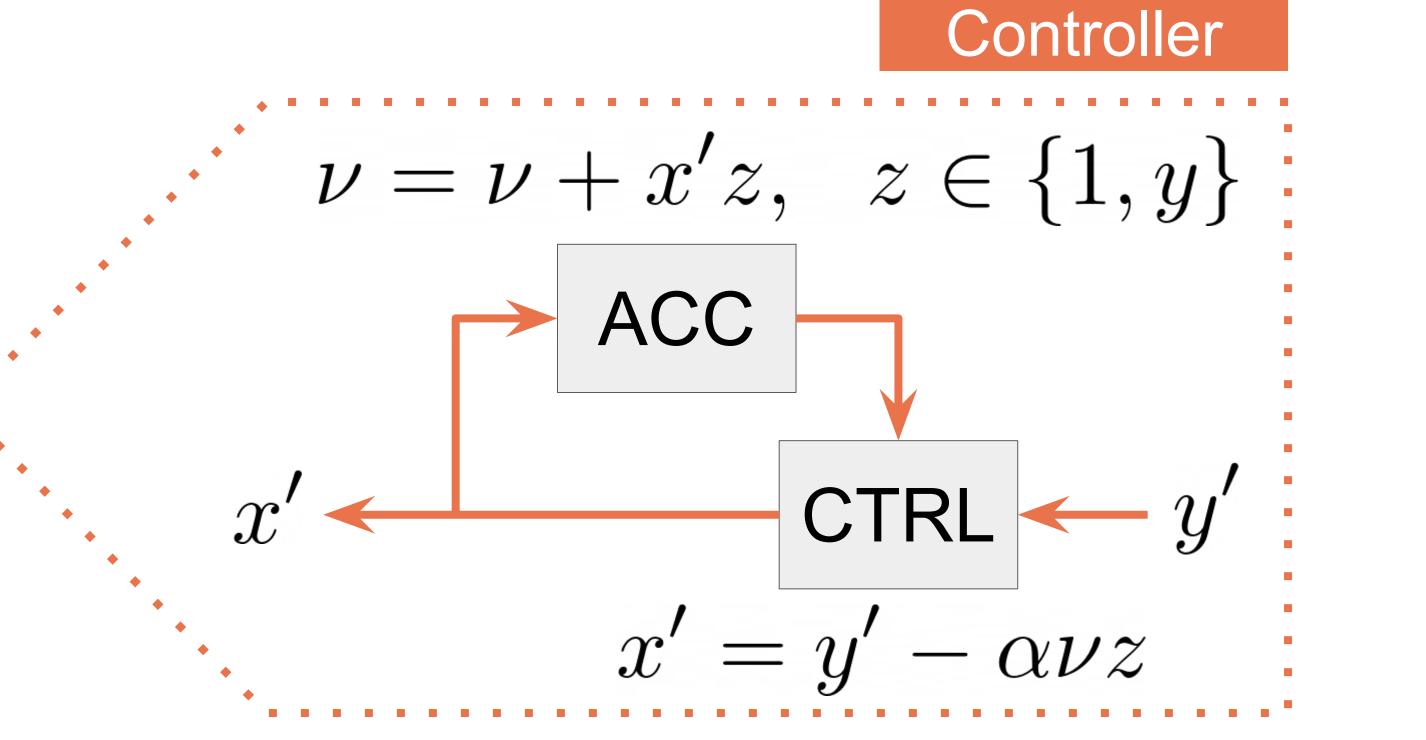
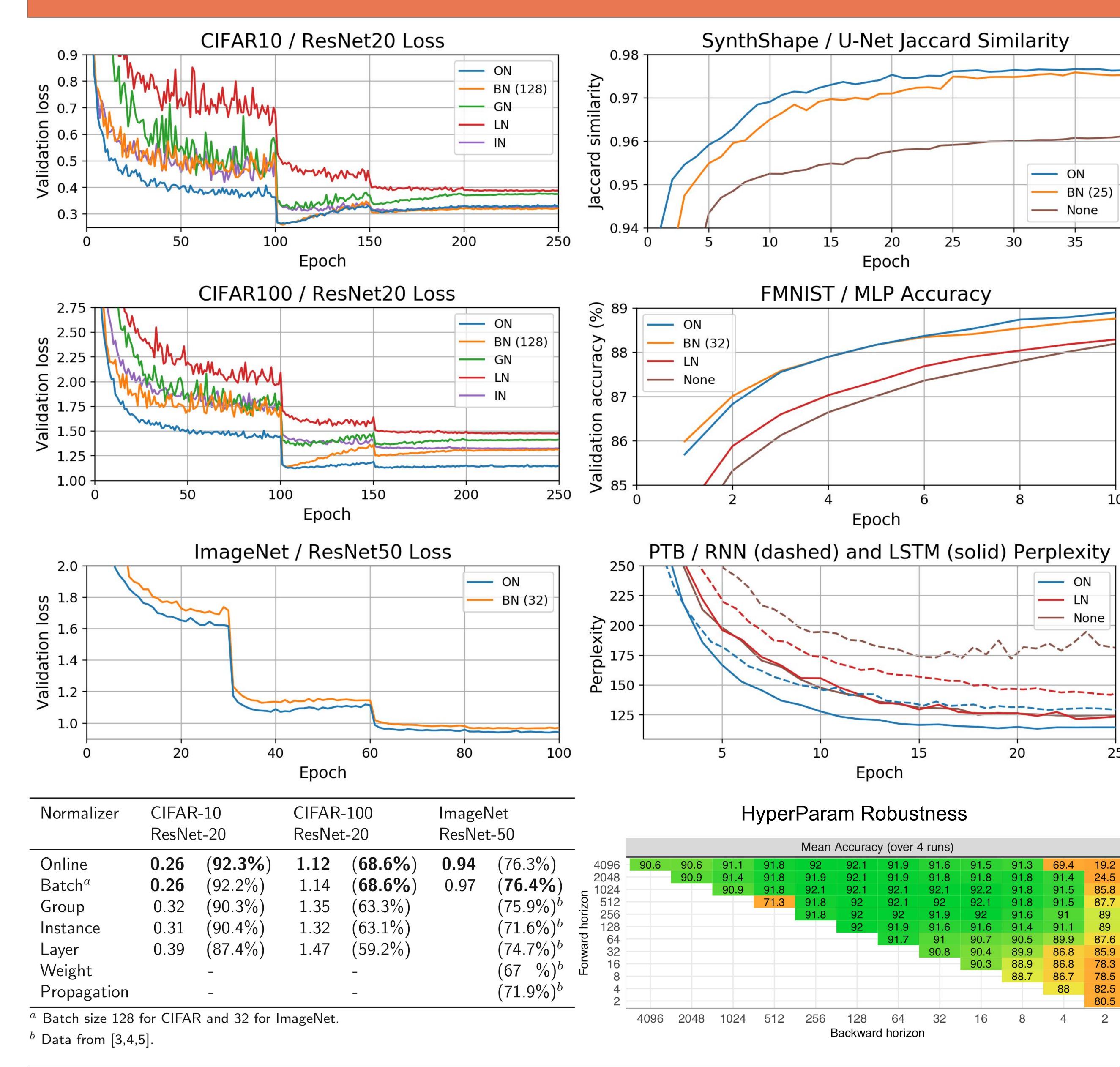


Image / Language / Generative Models



References

- [1] Sergey loffe and Christian Szegedy. "Batch normalization: Accelerating deep network training by reducing internal covariate shift."
- [2] Sergey loffe. "Batch renormalization: Towards reducing minibatch dependence in batch-normalized models." [3] Yuxin Wu and Kaiming He. "Group normalization."
- [4] Igor Gitman and Boris Ginsburg. "Comparison of batch normalization and weight normalization algorithms for the
- large-scale image classification." [5] Wenling Shang, Justin Chiu, and Kihyuk Sohn. "Exploring normalization in deep residual networks with concatenated rectified linear units."