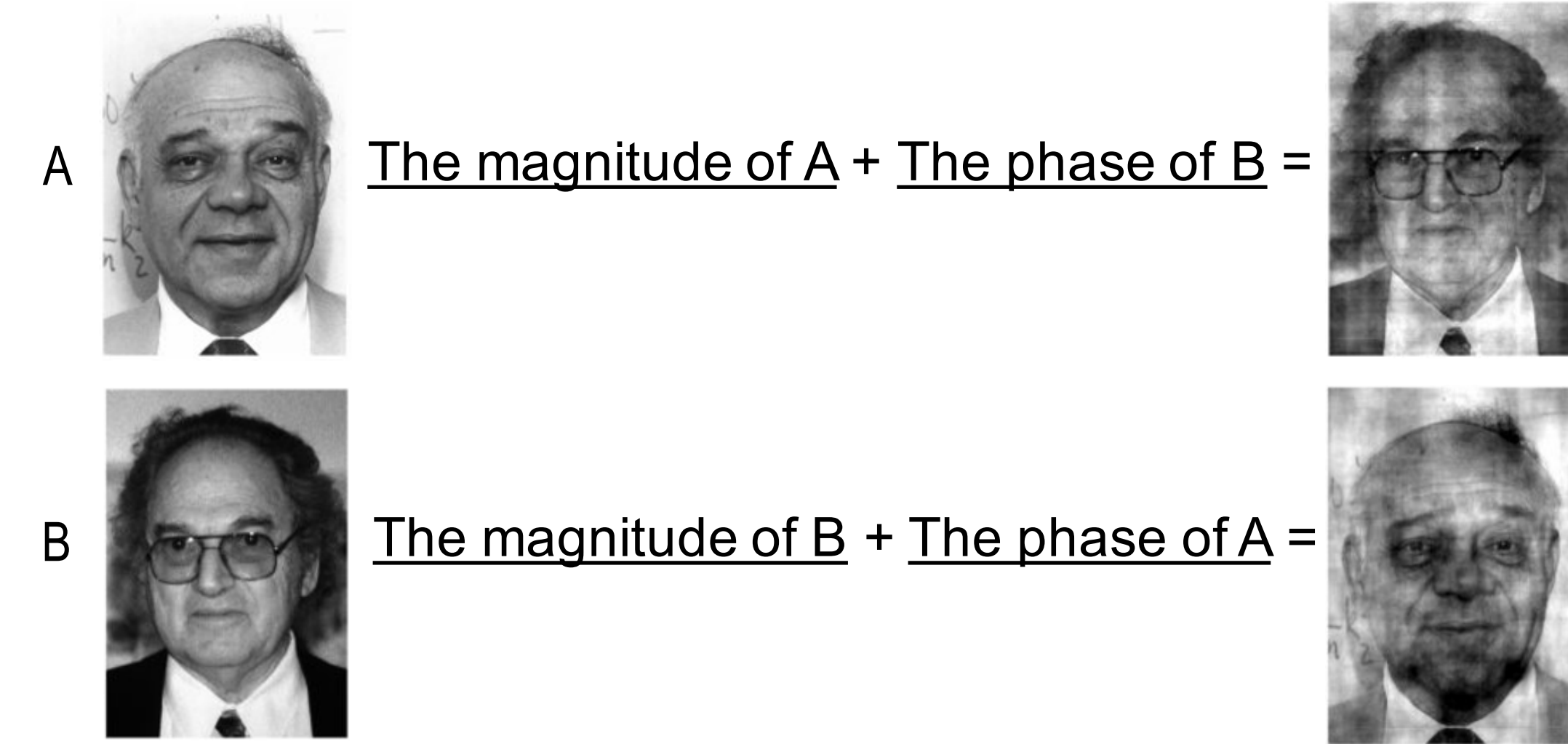


Introduction

Motivation I

- 2D Fourier Transform for images



- Phase contains the crucial discriminative information!

Motivation II

- Angles (Phase) usually give bounded output, avoiding covariate shift problem and stabilize the network training.
- For example, we usually use the cosine function of angles, which produces output from -1 to 1. The internal covariate shift can be largely prevented.

Motivation III

- By only using the angular information in the network learning, we could largely reduce the learning space of parameters, which could accelerate the training process and speed up the convergence.

SphereNet: A Neural Network Learned on Hyperspheres

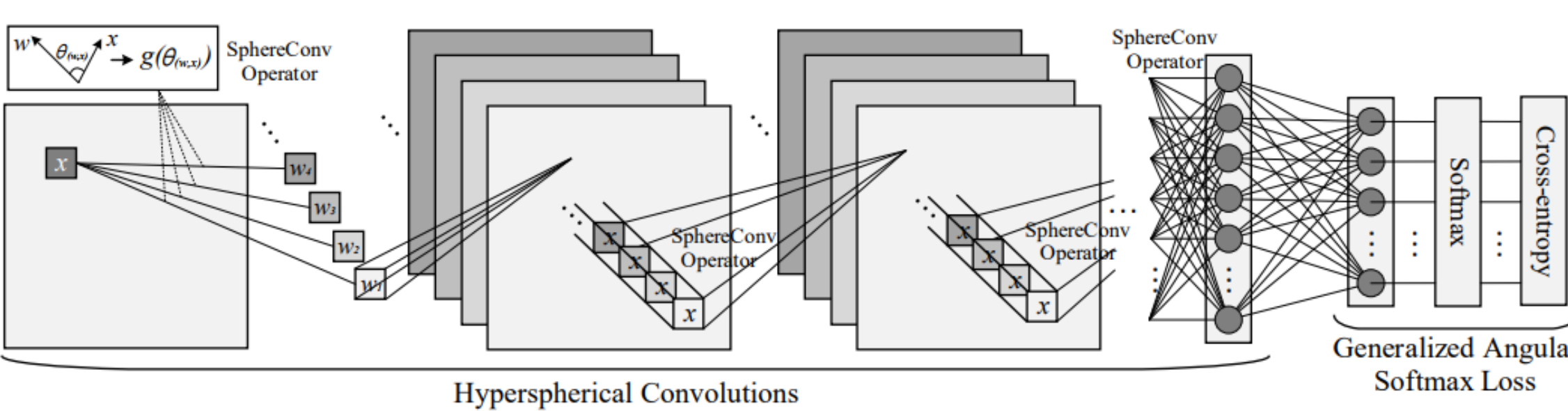
Hyperspherical Convolutional (SphereConv) Operator:

$$\mathcal{F}_s(w, x) = g(\theta_{(w, x)}) + b_{\mathcal{F}_s}$$

where $\theta_{(w, x)}$ is the angle between the kernel parameter w and the local patch x . A simple example is cosine SphereConv:

$$g(\theta_{(w, x)}) = \cos(\theta_{(w, x)})$$

We use this SphereConv operator to replace the original inner product based convolutional operator in the CNNs, and propose the *SphereNet*. (SphereNet comes from that angle can be viewed as the geodesic distance on a unit hypersphere)



Hyperspherical Convolutional Operator

- Three SphereConv operators:

linear SphereConv

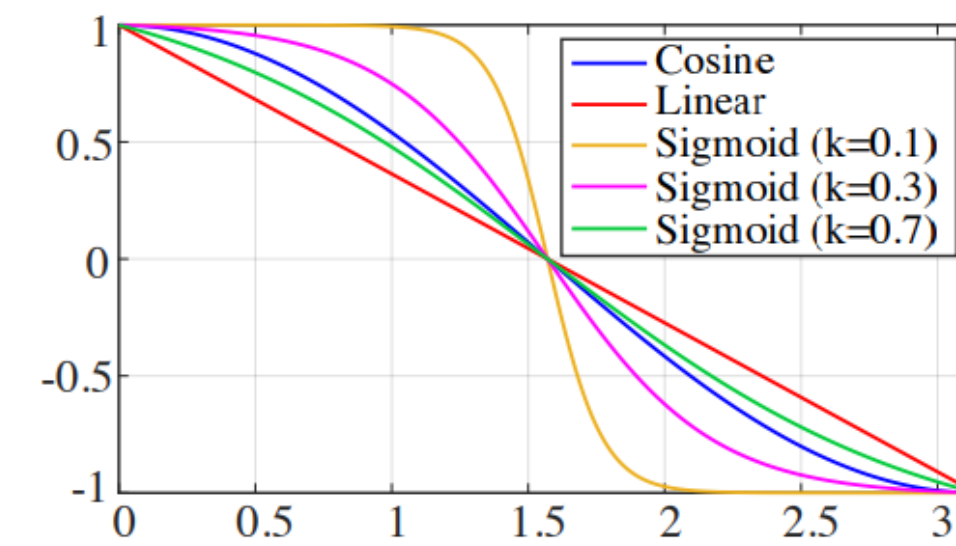
$$g(\theta_{(w, x)}) = a\theta_{(w, x)} + b$$

cosine SphereConv

$$g(\theta_{(w, x)}) = \cos(\theta_{(w, x)})$$

sigmoid SphereConv

$$g(\theta_{(w, x)}) = \frac{1 + \exp(-\frac{\pi}{2k})}{1 - \exp(-\frac{\pi}{2k})} \cdot \frac{1 - \exp(\frac{\theta_{(w, x)}}{k} - \frac{\pi}{2k})}{1 + \exp(\frac{\theta_{(w, x)}}{k} - \frac{\pi}{2k})}$$



- Besides the predefined SphereConv operators, we further consider the *learnable SphereConv*.
- SphereConv can also be used to the fully connected layers, recurrent layers, etc.
- We also design angular loss functions for *SphereConv*, i.e., generalized angular softmax (GA-Softmax) loss

Theoretical Insights

- Suppose the observation is $F = U^*V^{*\top}$ (ignore the bias), where $U^* \in \mathbb{R}^{n \times k}$ is the weight, $V^* \in \mathbb{R}^{m \times k}$ is the input that embeds weights from previous layers.

Scaling Issue of Neural Nets:

- Consider the objective:

$$\min_{U \in \mathbb{R}^{n \times k}, V \in \mathbb{R}^{m \times k}} \mathcal{G}(U, V) = \frac{1}{2} \|F - UV^\top\|_F^2$$

Lemma1: Consider a pair of global optimal points U, V satisfying $F = UV^\top$ and $\text{Tr}(V^\top V \otimes I_n) \leq \text{Tr}(U^\top U \otimes I_m)$. For any real $c > 1$, let $\tilde{U} = cU$ and $\tilde{V} = V/c$, then we have $\kappa(\nabla^2 \mathcal{G}(\tilde{U}, \tilde{V})) = \Omega(c^2 \kappa(\nabla^2 \mathcal{G}(U, V)))$, where $\kappa = \frac{\lambda_{\max}}{\lambda_{\min}}$ is the restricted condition number with λ_{\max} being the largest and λ_{\min} being the smallest nonzero eigenvalues.

Insensitiveness to Scaling for SphereConv:

- Consider our proposed cosine SphereConv operator, an equivalent problem is:

$$\min_{U \in \mathbb{R}^{n \times k}, V \in \mathbb{R}^{m \times k}} \mathcal{G}_S(U, V) = \frac{1}{2} \|F - D_U UV^\top D_V\|_F^2$$

where $D_U = \text{diag}(\frac{1}{\|U_{1,:}\|_2}, \dots, \frac{1}{\|U_{n,:}\|_2}) \in \mathbb{R}^{n \times n}$ and $D_V = \text{diag}(\frac{1}{\|V_{1,:}\|_2}, \dots, \frac{1}{\|V_{m,:}\|_2}) \in \mathbb{R}^{m \times m}$ are diagonal matrices.

Lemma2: For any real $c > 1$, let $\tilde{U} = cU$ and $\tilde{V} = V/c$, then we have $\lambda_i(\nabla^2 \mathcal{G}_S(\tilde{U}, \tilde{V})) = \lambda_i(\nabla^2 \mathcal{G}_S(U, V))$ for all $i \in [(n+m)k] = \{1, 2, \dots, (n+m)k\}$ and $\kappa(\nabla^2 \mathcal{G}(\tilde{U}, \tilde{V})) = \kappa(\nabla^2 \mathcal{G}(U, V))$, where κ is defined as in Lemma1.

- Regular Neural Nets: scales as $\Omega(c^2)$
- SphereConv: **insensitive** to scaling

Learnable SphereConv Operator

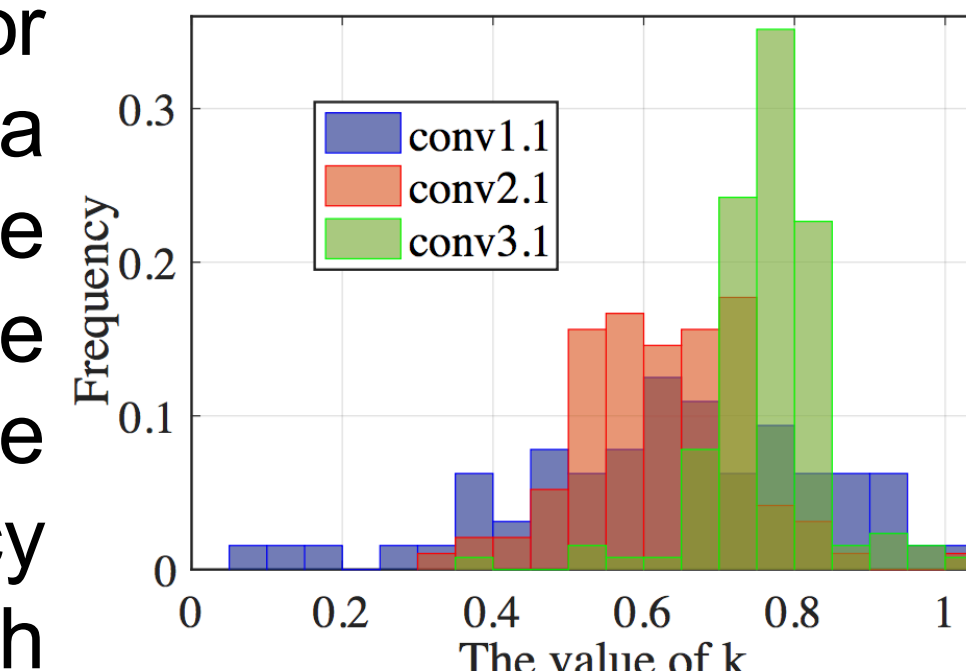
- With sigmoid SphereConv, we naturally come up with a learnable SphereConv. Specifically, we propose to learn the parameter k in the sigmoid SphereConv.

$$g(\theta_{(w, x)}) = \frac{1 + \exp(-\frac{\pi}{2k})}{1 - \exp(-\frac{\pi}{2k})} \cdot \frac{1 - \exp(\frac{\theta_{(w, x)}}{k} - \frac{\pi}{2k})}{1 + \exp(\frac{\theta_{(w, x)}}{k} - \frac{\pi}{2k})}$$

k is learned by back-prop!

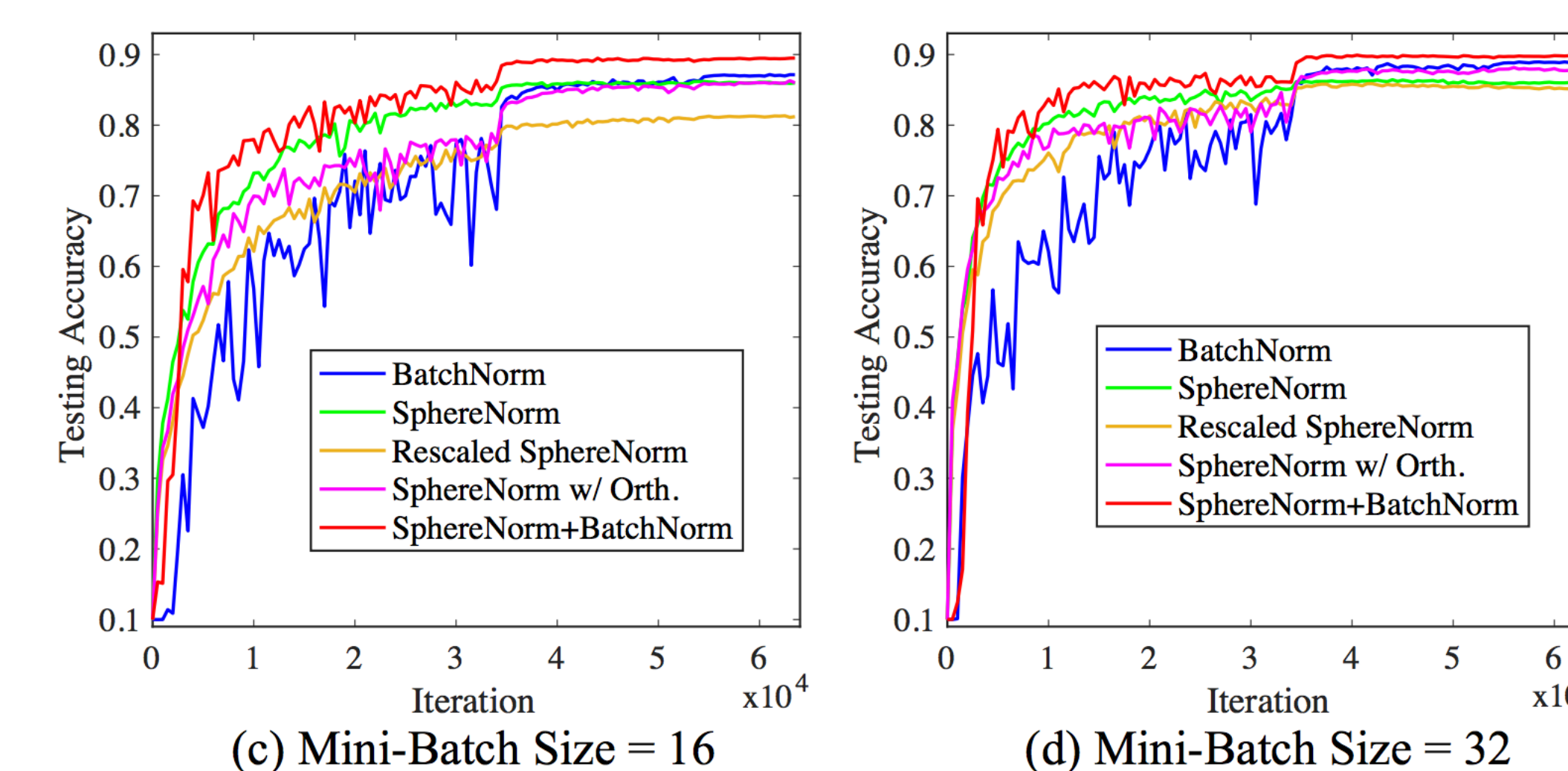
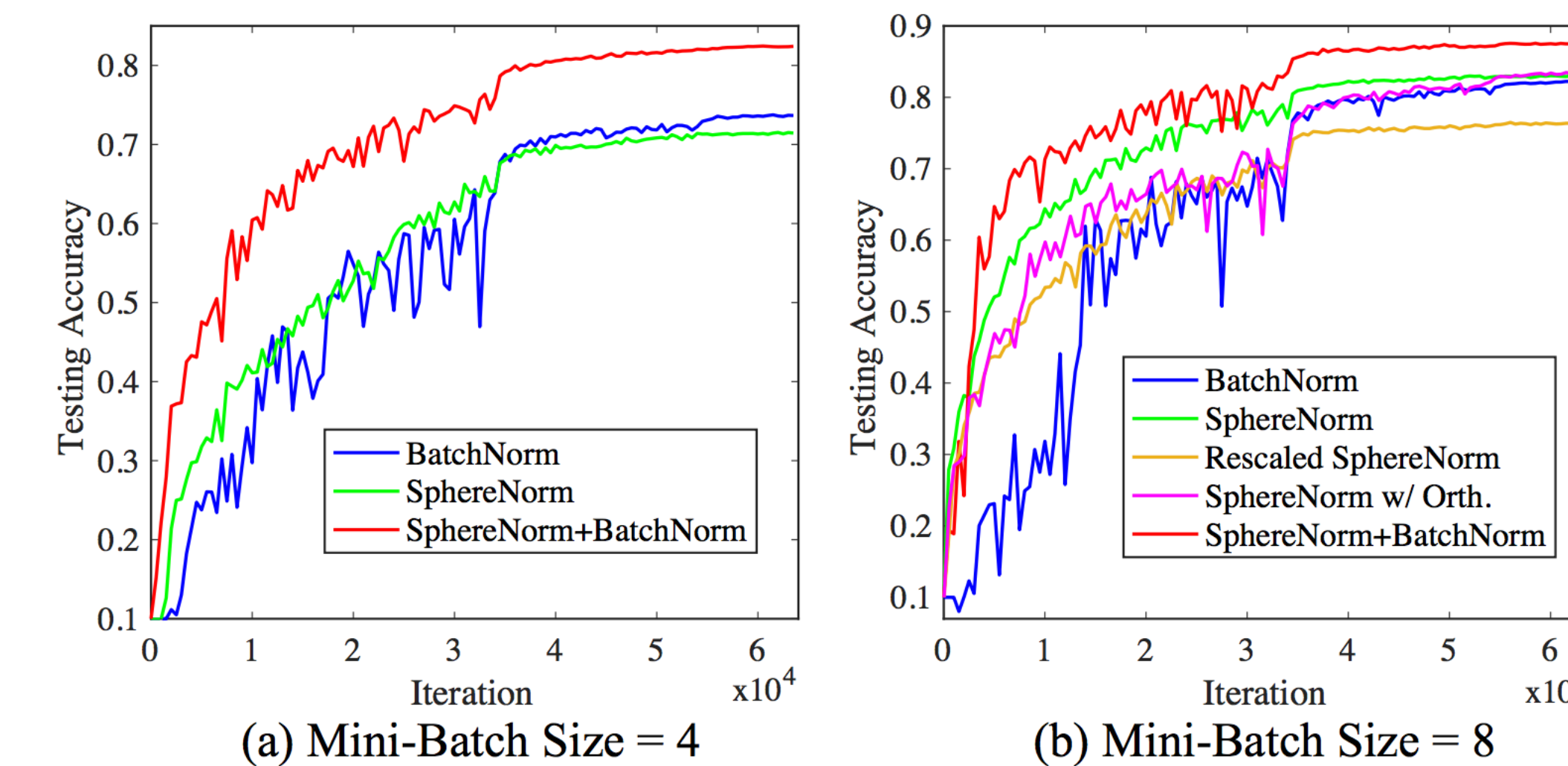
- K is updated using $k^{t+1} = k^t + \eta \frac{\partial L}{\partial k}$ where t denotes the iteration and $\frac{\partial L}{\partial k}$ is computed by the chain rule.

- Preliminary results: we learn a parameter k independently for each kernel and draw a frequency histogram for the value of k . Note that, we initialize all k with the same constant 0.5. The final accuracy can be further boosted with learnable SphereConv.



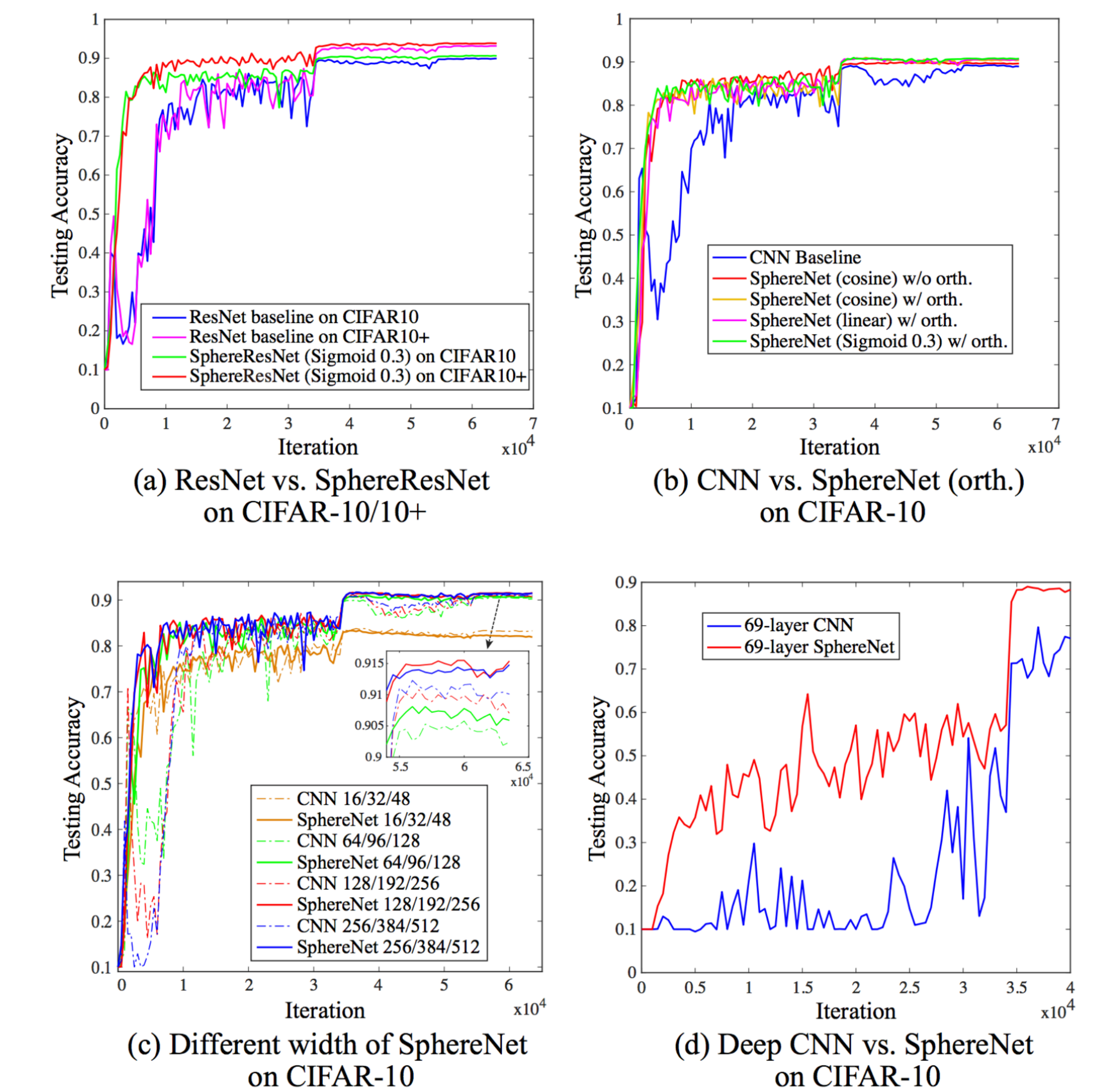
SphereNorm: a New Normalization Method

- Similar to batch normalization (BatchNorm), we note that the hyperspherical learning can also be viewed as a way of normalization, because SphereConv constrain the output value in $[-1, 1]$ ($[0, 1]$ after ReLU).
- Different from BatchNorm, SphereNorm normalizes the network based on spatial information and the weights, so it has nothing to do with the mini-batch statistic.
- SphereNorm and BatchNorm are complimentary to each other and could be used simultaneously.



Experiments on CIFAR-10 and CIFAR100

- On both CIFAR-10 and CIFAR-100, we observe the faster convergence on multiple network architectures like plain CNNs and ResNets.



- Our SphereResNet uses only 34 layers to perform comparably to the 1001-layer ResNet.

Method	CIFAR-10+	CIFAR-100
ELU [2]	94.16	72.34
FitResNet (LSUV) [14]	93.45	65.72
ResNet-1001 [7]	95.38	77.29
Baseline ResNet-32 (softmax)	93.26	72.85
SphereResNet-32 (S-SW)	94.47	76.02
SphereResNet-32 (L-LW)	94.33	75.62
SphereResNet-32 (C-CW)	94.64	74.92
SphereResNet-32 (S-G)	95.01	76.39

Experiments on Imagenet-2012

- Our SphereResNet also shows much faster convergence on large-scale dataset like Imagenet-2012.
- The error is the single central crop error.

