# Switchable Normalization for Learning-to-Normalize Deep Representation

Ping Luo, Ruimao Zhang, Jiamin Ren, Zhanglin Peng, Jingyu Li

**Abstract**—We address a learning-to-normalize problem by proposing Switchable Normalization (SN), which learns to select different normalizers for different normalization layers of a deep neural network. SN employs three distinct scopes to compute statistics (means and variances) including a channel, a layer, and a minibatch. SN switches between them by learning their importance weights in an end-to-end manner. It has several good properties. First, it adapts to various network architectures and tasks (see Fig.1). Second, it is robust to a wide range of batch sizes, maintaining high performance even when small minibatch is presented (*e.g.* 2 images/GPU). Third, SN does not have sensitive hyper-parameter, unlike group normalization that searches the number of groups as a hyper-parameter. Without bells and whistles, SN outperforms its counterparts on various challenging benchmarks, such as ImageNet, COCO, CityScapes, ADE20K, MegaFace and Kinetics. Analyses of SN are also presented to answer the following three questions: (a) Is it useful to allow each normalization layer to select its own normalizer? (b) What impacts the choices of normalizers? (c) Do different tasks and datasets prefer different normalizers? We hope SN will help ease the usage and understand the normalization techniques in deep learning. The code of SN has been released at https://github.com/switchablenorms.

**Index Terms**—Deep Learning, Normalization, Image/Video Classification, Object Detection, Semantic Segmentation and Face Verification

✦

## 1 INTRODUCTION

Normalization techniques are effective components in deep learning, advancing many research fields such as natural language processing, computer vision, and machine learning. In recent years, many normalization methods such as Batch Normalization (BN) [1], Instance Normalization (IN) [2], and Layer Normalization (LN) [3] have been developed. Despite their great successes, existing practices often employed the same normalizer in all normalization layers of an entire network, rendering suboptimal performance. Also, different normalizers are used to solve different tasks, making model design cumbersome.

To address the above issues, we propose *Switchable Normalization (SN)*, which combines three types of statistics estimated channel-wise, layer-wise, and minibatch-wise by using IN, LN, and BN respectively. SN switches among them by learning their importance weights. By design, *SN is adaptable to various deep networks and tasks*. For example, the ratios of IN, LN, and BN in SN are compared in multiple tasks as shown in Fig.1 (a). We see that using one normalization method uniformly is not optimal for these tasks. For instance, image classification and object detection prefer the combination of three normalizers. In particular, SN chooses BN more than IN and LN in image classification and the backbone network of object detection, while LN has larger weights in the box and mask heads. For artistic image style transfer [4], SN selects IN. For neural architecture search, SN is applied to LSTM where LN is preferable than group normalization (GN) [5], which is a variant of IN by dividing channels into groups.
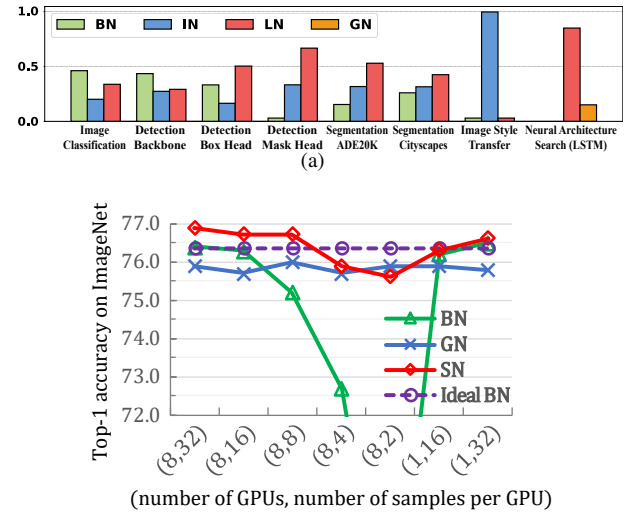


Fig. 1: (a) shows that SN adapts to various networks and tasks by learning importance ratios to select normalizers. In (a), a ratio is between 0 and 1 and all ratios of each task sum to 1. (b) shows the top-1 accuracies of ResNet50 trained with SN on ImageNet and compared with BN and GN in different batch settings. The gradients in training are averaged over all GPUs and the statistics of normalizers are estimated in each GPU. For instance, all methods are compared to an ideal case, 'ideal BN', whose accuracies are 76.4% for all settings. This ideal case cannot be obtained in practice. In fact, when the minibatch size decreases, BN's accuracies drop significantly, while SN and GN both maintain reasonably good performance. SN surpasses or is comparable to both BN and GN in all settings.

The selectivity of normalizers makes *SN robust to minibatch size*. As shown in Fig.1 (b), when training ResNet50 [6] on ImageNet [7] with different batch sizes, SN is close to the

Ping Luo is with Department of Computer Science, The University of Hong Kong, Hong Kong, China. (Email: pluo.lhi@gmail.com).
Ruimao Zhang, Jiamin Ren, Zhanglin Peng are with SenseTime Research, Shenzhen, China. (Email: ruimao.zhang@ieee.org; renjiamin@sensetime.com; pengzhanglin@sensetime.com ). Ruimao Zhang is also with The Chinese University of Hong Kong. He is the corresponding author of this paper.
Jingyu Li is with Department of Electronic Engineering, The Chinese University of Hong Kong, Hong Kong, China. (Email: jingyuli@cuhk.edu.hk ).

"ideal case" more than BN and GN. For $(8, 32)$ as an example[1], ResNet50 trained with SN is able to achieve 76.9% top-1 accuracy, surpassing BN and GN by 0.5% and 1.0% respectively. In general, SN obtains better or comparable results than both BN and GN in all batch settings.

Overall, this work has three key **contributions**. (1) We introduce Switchable Normalization (SN), which is applicable in both CNNs and RNNs/LSTMs, and improves the other normalization techniques on many challenging benchmarks and tasks including image recognition in ImageNet [8], object detection in COCO [9], scene parsing in Cityscapes [10] and ADE20K [11], face recognition in MegaFace [12], video recognition in Kinetics [13], artistic image stylization [4] and neural architecture search [14]. (2) The analyses of SN are presented where multiple normalizers can be compared and understood with geometric interpretation. In addition, through systematic experiments, we answer three critical questions: i) Is it useful to allow each normalization layer to select its own normalizer? ii) What impacts the choices of normalizers? iii) Do different tasks and datasets prefer different normalizers? The answer of these questions give the suggestions of how to use SN in modern deep neural networks. (3) By enabling *each normalization layer in a deep network to have its own operation*, SN helps ease the usage of normalizers, pushes the frontier of normalization in deep learning, as well as opens up new research direction. We believe that all existing models could be reexamined with this new perspective. We make the code of SN available and recommend it as an alternative of existing handcrafted approaches.

In the following sections, we first review the related work in Sec.2 and then present SN in Sec.3. The performance of SN is evaluated extensively in Sec.4. We conclude our work and summarize some advices that help training CNNs with SN in Sec. 5. The future research directions are also presented in this section.

## 2 RELATED WORK

**Normalization.** As one of the most significant component in deep neural networks, normalization technique has achieved much attention in literature [1], [2], [3], [5], [17], [18]. According to the normalized space, related methods can be divided into two groups: methods normalizing activated representation over feature space such as [1], [2], [3], [5], and methods normalizing weights over the parameter space like [17], [18].

Among former group, Batch Normalization (BN) [1] is one of the most well-known method, which is motivated by the fact that whitening input features (*i.e.* centering, decorrelating and scaling) [19], [20] of each hidden layer can mimic the fast convergence of natural gradient descent (NGD) [21] by using stochastic gradient descent (SGD). Even though BN only pursuit standardization (*i.e.* centering and scaling) instead of whitening transformation, it still stabilizes the training process and boost the performance of neural network. To better understand the effectiveness of BN, a lot of theory studies have been proposed [22], [23], [24], [25], [26], [27]. For example, Santurkar *et. al.* [23] argued that BN did not help covariate shift but contributed to smooth loss landscape. Hoffer *et. al.* [24] developed a framework to decouple weights' norm from the underlying optimized objective. Luo *et. al.* [25] decomposed BN into population normalization (PN) and

1. In this work, minibatch size refers to the number of samples per GPU, and batch size is '#GPUs' times '#samples per GPU'. A batch setting is denoted as a 2-tuple, (#GPUs, #samples per GPU).

gamma decay as an explicit regularization, and claimed that BN allowed bigger learning rate in the training process. The similar conclusion was also achieved by Bjorck*et. al.* [28]. In [26], Yang *et. al.* developed a mean field theory for batch normalization in fully-connected feedforward neural networks. Arora *et. al.* [27] found that BN had the ability to allow gradient descent to succeed with less tuning of learning rates.

Despite a series of theory analyses, a lot of work have been proposed in practice to address the issue that BN required reasonable mini-batch size to estimate the mean and variance. Ba *et. al.* [3] proposed Layer Normalization (LN) to calculate the mean and variance for each sample on a single layer. By leveraging BN and LN, Ren *et. al.* [29] proposed Division Normalization to explore spatial region. In [2], Instance Normalization (IN) is proposed as a channel-wise normalization method to filter out complex appearance variance [30]. Wu and He [5] further proposed Group Normalization (GN) by dividing the channels into groups and computing within each group the mean and variance for normalization. Luo *et. al.* [31] proposed Dynamic Normalization (DN) to learn arbitrary normalization for different convolutional layers. Such method achieved stable accuracy in a wide range of batch sizes. Other attempts to deal with instability of BN with small batch-size including Batch Renormalization (BRN) [15], Batch Kalman Normalization (BKN) [16] and Stream Normalization (StN) [32].

Some work also proposed to promote the performance of BN on other aspects. For example, Huang *et. al.* [33] proposed Decorrelated Batch Normalization (DBN) to whiten the activations of each layer within a mini-batch to improve optimization efficiency and generalization ability of BN. In [30], Pan *et. al.* showed that the hybrid of BN and IN in the neural networks can greatly strengthen the generalization ability of CNN on different domains. In [34], Perez *et. al.* adopted Conditional Batch Normalization (CBN) to affect the feature representation of each sample independently in BN layers, and achieved promising results on visual question answering task. In [35], Pan *et. al.* proposed Switchable Whitening (SW) to select different whitening methods as well as standardization methods in a single layer.

For the methods normalizing the weights, Salimans *et. al.* [17] firstly proposed Weight Normalization (WN) to normalize the weight for each neuron to decouple the length of weight vectors from their directions. Huang *et. al.*[36] further introduced Centered Weight Normalization (CWN) to further powered WN by centering the input weight. In recent studies on Generative Adversarial Network (GAN) [37], Miyato *et. al.* [18] proposed Spectral Normalization (SpN) to re-scale the weight vector by dividing its largest singular value to satisfy the Lipschitz constraint. This method stabilized the training process of the discriminator in GAN, and effectively improved the quality of generated images.

In Table 1, we compare SN with five popular normalization methods, *i.e.* BN, IN, LN, GN and WN, as well as three variants of BN including Batch Renormalization (BRN) and Batch Kalman Normalization (BKN) in details. In general, we see that SN possesses comparable numbers of parameters and computations, as well as rich statistics. **First**, although SN has richer statistics, the computational complexity to estimate them is comparable to previous methods, as shown in the second portion of Table 1. As introduced in Sec.3, IN, LN, and BN estimate the means and variances along axes $(H, W)$, $(C, H, W)$, and $(N, H, W)$ respectively, leading to $2CN$, $2N$, and $2C$ numbers of statistics. Therefore, SN has $2CN + 2N + 2C$ statistics by combining them.

| | Parameter | | | Statistical Estimation | | |
|---|---|---|---|---|---|---|
| | params | #params | hyper-params | statistics | computation complexity | #statistics |
| BN [1] | $\gamma, \beta$ | $2C$ | $p, \epsilon$ | $\mu, \sigma, \mu', \sigma'$ | $\mathcal{O}(NCHW)$ | $2C$ |
| IN [2] | $\gamma, \beta$ | $2C$ | $\epsilon$ | $\mu, \sigma$ | $\mathcal{O}(NCHW)$ | $2CN$ |
| LN [3] | $\gamma, \beta$ | $2C$ | $\epsilon$ | $\mu, \sigma$ | $\mathcal{O}(NCHW)$ | $2N$ |
| GN [5] | $\gamma, \beta$ | $2C$ | $g, \epsilon$ | $\mu, \sigma$ | $\mathcal{O}(NCHW)$ | $2gN$ |
| BRN [15] | $\gamma, \beta$ | $2C$ | $p, \epsilon, r, d$ | $\mu, \sigma, \mu', \sigma'$ | $\mathcal{O}(NCHW)$ | $2C$ |
| BKN [16] | $A$ | $C^2$ | $p, \epsilon$ | $\mu, \Sigma, \mu', \Sigma'$ | $\mathcal{O}(NC^2HW)$ | $C + C^2$ |
| WN [17] | $\gamma$ | $C$ | $-$ | $-$ | $-$ | $-$ |
| SN | $\gamma, \beta,$ $\{w_k\}_{k \in \Omega}$ | $2C + 6$ | $\epsilon$ | $\{\mu_k, \sigma_k\}_{k \in \Omega}$ | $\mathcal{O}(NCHW)$ | $2C + 2N$ $+2CN$ |

TABLE 1: **Comparisons of normalization methods**. First, we compare their types of parameters, numbers of parameters (#params), and hyper-parameters. Second, we compare types of statistics, computational complexity to estimate statistics, and numbers of statistics (#statistics). Specifically, $\gamma, \beta$ denote the scale and shift parameters. $\mu, \sigma, \Sigma$ are a vector of means, a vector of standard deviations, and a covariance matrix. $\mu'$ represents the moving average. Moreover, $p$ is the momentum of moving average, $g$ in GN is the number of groups, $\epsilon$ is a small value for numerical stability, and $r, d$ are used in BRN. In SN, $k \in \Omega$ indicates a set of different kinds of statistics, $\Omega = \{\text{in}, \text{ln}, \text{bn}\}$, and $w_k$ is an importance weight of each kind.

Although BKN has the largest number of $C + C^2$ statistics, it also has the highest computations because it estimates the covariance matrix other than the variance vector. Also, approximating the covariance matrix in a minibatch is nontrivial as discussed in [19], [20], [38]. BN, BRN, and BKN also compute moving averages. **Second**, SN is demonstrated in various networks, tasks, and datasets. Its applications are much wider than existing normalizers and it also has rich theoretical value that is worth exploring.

**Meta Learning.** Our work is also related with meta learning (ML) problem [39], [40], [14], which can be also used to learn the control parameters in SN. In general, ML is defined as $\min_\Theta \mathcal{L}(\Theta, \Phi^*) + \min_\Phi \varphi \mathcal{L}(\Theta^*, \Phi)$ where $\varphi$ is a constant multiplier. Unlike SN trained with a single stage, this loss function is minimized by performing two feed-forward stages iteratively until converged. First, by fixing the current estimated control parameters $\Phi^*$, the network parameters are optimized by $\Theta^* = \min_\Theta \mathcal{L}(\Theta, \Phi^*)$. Second, by fixing $\Theta^*$, the control parameters are found by $\Phi^* = \min_\Phi \varphi \mathcal{L}(\Theta^*, \Phi)$.

The above two stages are usually optimized by using two different sets of training data. For example, previous work [41], [14] used $\Phi$ to search network architectures from a set of modules with different numbers of parameters and computational complexities. They divided an entire training set into a training and a validation set without overlapping, where $\Phi$ is learned from the validation set while $\Theta$ is learned from the training set. This is because $\Phi$ would choose the module with large complexity to overfit training data, if both $\Theta$ and $\Phi$ are optimized in the same dataset.

The above 2-stage training increases computations and runtime. In contrast, $\Theta$ and $\Phi$ for SN can be generally optimized within a single stage in the same dataset, because $\Phi$ regularizes training by choosing different normalizers from $\Omega$ to prevent overfitting.

## 3 SWITCHABLE NORMALIZATION (SN)

We describe a general formulation of a normalization layer and then present SN in this section.

**A General Form.** We take CNN as an illustrative example. Let $h$ be the input data of an arbitrary normalization layer represented by a 4D tensor $(N, C, H, W)$, indicating number of samples, number of channels, height and width of a channel respectively.

Let $h_{ncij}$ and $\hat{h}_{ncij}$ be a pixel before and after normalization, where $n \in [1, N]$, $c \in [1, C]$, $i \in [1, H]$, and $j \in [1, W]$. Let $\mu$ and $\sigma$ be a mean and a standard deviation. We have

$$\hat{h}_{ncij} = \gamma \frac{h_{ncij} - \mu}{\sqrt{\sigma^2 + \epsilon}} + \beta, \tag{1}$$

where $\gamma$ and $\beta$ are a scale and a shift parameter respectively. $\epsilon$ is a small constant to preserve numerical stability. Eqn.(1) shows that each pixel is normalized by using $\mu$ and $\sigma$, and then re-scale and re-shift by $\gamma$ and $\beta$.

In practice, IN, LN, and BN share the formulation of Eqn.(1), but they use different sets of pixels to estimate $\mu$ and $\sigma$. In other words, the numbers of their estimated statistics are different. In general, we have

$$\mu_k = \frac{1}{|I_k|} \sum_{(n,c,i,j) \in I_k} h_{ncij},$$
$$\sigma_k^2 = \frac{1}{|I_k|} \sum_{(n,c,i,j) \in I_k} (h_{ncij} - \mu_k)^2, \tag{2}$$

where $k \in \{\text{in}, \text{ln}, \text{bn}\}$ is used to distinguish different methods. $I_k$ is a set pixels and $|I_k|$ denotes the number of pixels. Specifically, $I_{\text{in}}$, $I_{\text{ln}}$, and $I_{\text{bn}}$ are the sets of pixels used to compute statistics in different approaches.

IN was established in the task of artistic image style transfer [4], [42], [43]. In IN, we have $\mu_{\text{in}}, \sigma_{\text{in}}^2 \in \mathbb{R}^{N \times C}$ and $I_{\text{in}} = \{(i, j) | i \in [1, H], j \in [1, W]\}$, meaning that IN has $2NC$ elements of statistics, where each mean and variance value is computed along $(H, W)$ for each channel of each sample.

LN [3] was proposed to ease optimization of recurrent neural networks (RNNs). In LN, we have $\mu_{\text{ln}}, \sigma_{\text{ln}}^2 \in \mathbb{R}^{N \times 1}$ and $I_{\text{ln}} = \{(c, i, j) | c \in [1, C], i \in [1, H], j \in [1, W]\}$, implying that LN has $2N$ statistical values, where a mean value and a variance value are computed in $(C, H, W)$ for each one of the $N$ samples.

BN [1] was first demonstrated in the task of image classification [6], [44] by normalizing the hidden feature maps of CNNs. In BN, we have $\mu_{\text{bn}}, \sigma_{\text{bn}}^2 \in \mathbb{R}^{C \times 1}$ and $I_{\text{bn}} = \{(n, i, j) | n \in [1, N], i \in [1, H], j \in [1, W]\}$, in the sense that BN treats each channel independently like IN, but not only normalizes across $(H, W)$, but also the $N$ samples in a minibatch, leading to $2C$ elements of statistics.

## 3.1 Formulation of SN

SN has an intuitive expression

$$\hat{h}_{ncij} = \gamma \frac{h_{ncij} - \Sigma_{k \in \Omega} w_k \mu_k}{\sqrt{\Sigma_{k \in \Omega} w_k' \sigma_k^2 + \epsilon}} + \beta, \qquad (3)$$

where $\Omega$ is a set of statistics estimated in different ways. In this work, we define $\Omega = \{\text{in}, \text{ln}, \text{bn}\}$ the same as above where $\mu_k$ and $\sigma_k^2$ can be calculated by following Eqn.(2). However, this strategy leads to large redundant computations. In fact, the three kinds of statistics of SN depend on each other. Therefore we could reduce redundancy by reusing computations,

$$\mu_{\text{in}} = \frac{1}{HW} \sum_{i,j}^{H,W} h_{ncij}, \quad \sigma_{\text{in}}^2 = \frac{1}{HW} \sum_{i,j}^{H,W} (h_{ncij} - \mu_{\text{in}})^2,$$

$$\mu_{\text{ln}} = \frac{1}{C} \sum_{c=1}^{C} \mu_{\text{in}}, \quad \sigma_{\text{ln}}^2 = \frac{1}{C} \sum_{c=1}^{C} (\sigma_{\text{in}}^2 + \mu_{\text{in}}^2) - \mu_{\text{ln}}^2,$$

$$\mu_{\text{bn}} = \frac{1}{N} \sum_{n=1}^{N} \mu_{\text{in}}, \quad \sigma_{\text{bn}}^2 = \frac{1}{N} \sum_{n=1}^{N} (\sigma_{\text{in}}^2 + \mu_{\text{in}}^2) - \mu_{\text{bn}}^2, \quad (4)$$

showing that the means and variances of LN and BN can be computed based on IN. Using Eqn.(4), the computational complexity of SN is $\mathcal{O}(NCHW)$, which is comparable to previous work.

Furthermore, $w_k$ and $w_k'$ in Eqn.(3) are importance ratios used to weighted average the means and variances respectively. Each $w_k$ or $w_k'$ is a scalar variable, which is shared across all channels. There are $3 \times 2 = 6$ importance weights in SN. We have $\Sigma_{k \in \Omega} w_k = 1$, $\Sigma_{k \in \Omega} w_k' = 1$, and $\forall w_k, w_k' \in [0, 1]$, and define

$$w_k = \frac{e^{\lambda_k}}{\Sigma_{z \in \{\text{in}, \text{ln}, \text{bn}\}} e^{\lambda_z}} \quad \text{and} \quad k \in \{\text{in}, \text{ln}, \text{bn}\}. \qquad (5)$$

Here each $w_k$ is computed by using a softmax function with $\lambda_{\text{in}}$, $\lambda_{\text{ln}}$, and $\lambda_{\text{bn}}$ as the control parameters, which can be learned by back-propagation (BP). $w_k'$ are defined similarly by using another three control parameters $\lambda_{\text{in}}'$, $\lambda_{\text{ln}}'$, and $\lambda_{\text{bn}}'$.

**Training.** Let $\Theta$ be a set of network parameters (*e.g.* filters) and $\Phi$ be a set of control parameters that control the network architecture. In SN, we have $\Phi = \{\lambda_{\text{in}}, \lambda_{\text{ln}}, \lambda_{\text{bn}}, \lambda_{\text{in}}', \lambda_{\text{ln}}', \lambda_{\text{bn}}'\}$. Training a deep network with SN is to minimize a loss function $\min_{\{\Theta, \Phi\}} \frac{1}{N} \sum_{j=1}^{N} \mathcal{L}(\mathbf{y}_j, f(\mathbf{x}_j; \Theta, \Phi))$, where $\{\mathbf{x}_j, \mathbf{y}_j\}_{j=1}^{N}$ indicates a set of training samples and their labels. $f(\mathbf{x}_j; \Theta)$ is a function learned by the CNN to predict $\mathbf{y}_j$. $\Theta$ and $\Phi$ are the parameters of the model that can be optimized jointly by back-propagation (BP). This training procedure is different from previous meta-learning algorithms such as network architecture search [39], [41], [14]. In previous work, $\Phi$ represents as a set of network modules with different learning capacities, where $\Theta$ and $\Phi$ were optimized in two BP stages iteratively by using two training sets that are non-overlapped. For example, previous work divided an entire training set into a training and a validation set. However, if $\Theta$ and $\Phi$ in previous work are optimized in the same set of training data, $\Phi$ would choose the module with large complexity to overfit these data. In contrast, SN essentially prevents overfitting by choosing normalizers to improve both learning and generalization ability as discussed below.

**Implementation.** SN can be easily implemented in existing softwares such as PyTorch and TensorFlow. The backward computation of SN can be obtained by automatic differentiation (AD)

in these softwares. Without AD, we need to implement back-propagation (BP) of SN, where the errors are propagated through $\mu_k$ and $\sigma_k^2$. We provide the derivations of BP in Appendix.

## 3.2 Analyses of SN

**Geometric View of SN.** To understand SN, we theoretically compare SN with BN, IN, and LN by representing them using weight normalization (WN) [17] that is independent of mean and variance. Specifically, the computation of WN is defined by $\hat{h}_{\text{wn}} = \frac{h}{\|\mathbf{w}_i\|_2} = \frac{\mathbf{w}_i^\top \mathbf{x}}{\|\mathbf{w}_i\|_2}$, where $\mathbf{w}$ and $\mathbf{x}$ represent a filter and an image patch. We use $i$ to indicate a filter of the $i$-th channel. As shown in Fig.2 (a), WN normalizes the norm of each filter to a unit sphere with length '1', and then rescales the length to $\gamma$ that is a learnable scale parameter. To simplify our discussions, we suppose this scale parameter is shared among all channels. In other words, WN decomposes the optimization of a filter into its direction and length.

By assuming $\mu = 0$ and $\sigma = 1$ for all the normalizers, we see that they can be also represented by using WN. For instance, IN turns into $\hat{h}_{\text{in}} = \frac{h - \mathbf{w}_i^\top \mathbb{E}[\mathbf{x}]}{\sqrt{\mathbf{w}_i^\top (\mathbb{E}[\mathbf{x}\mathbf{x}^\top] - \mathbb{E}[\mathbf{x}]^\top \mathbb{E}[\mathbf{x}]) \mathbf{w}_i}} = \frac{\mathbf{w}_i^\top \mathbf{x}}{\|\mathbf{w}_i\|_2}$ that is identical to WN as shown in Fig.2 (b). For LN, each channel is normalized by all the channels and thus $\hat{h}_{\text{ln}} = \frac{\mathbf{w}_i^\top \mathbf{x}}{\sqrt{\frac{1}{C} \sum_{i=1}^{C} \|\mathbf{w}_i\|_2^2}}$ where $C$ is the number of channels. The filter norm in LN is less constrained than WN and IN as visualized in Fig.2 (c), where the filter norm can be either longer or shorter than $\gamma$, in the sense that LN increases learning capacity of the networks by diminishing regularization.

Furthermore, BN can be represented by WN with a similar formula as IN, where $\mathbf{x}$ indicates the image patch sampled from all of the instances in a mini-batch. In [25], Luo *et al.* showed that BN imposes regularization on norms of all the filters and reduces correlations between filters. Its geometry interpretation can be viewed in Fig.2 (d), where the filter norms would be shorter and angle between filters would be larger than the other normalizers. In conclusion, BN improves generalization [1]. In general, we would have the following relationships

$$\text{learning capacity} : LN > BN > IN,$$
$$\text{generalization capacity} : BN > IN > LN.$$

Finally, $\hat{h}_{\text{sn}}$ in SN inherits the benefits from all of them and enables balance between learning and generalization ability. For example, when the batch size is small, the random noise from the batch statistics of BN would be too strong. SN is able to maintain performance by decreasing $w_{\text{bn}}$ and increasing $w_{\text{ln}}$, such that the regularization from BN is reduced and the learning ability is enhanced by LN. This phenomenon is supported by our experiment.

**Variants of SN.** SN has many extensions. For instance, a pre-trained network with SN can be finetuned by applying the argmax function on its control parameters where each normalization layer selects only one normalizer, leading to sparse SN. For $(8, 32)$ as an example, SN with sparsity achieves top-1 accuracy of 77.0% in ImageNet with ResNet50, which is comparable to 76.9% of SN without sparsity. Moreover, when the channels are divided into groups, each group could select its own normalizer to increase representation power of SN. Our preliminary results suggest that group SN performs better than SN in some senses. For instance, group SN with only two groups boosts the top-1 accuracy of ResNet50 to 77.2% in ImageNet. The above two variants will be presented as future work due to the length of paper. This work
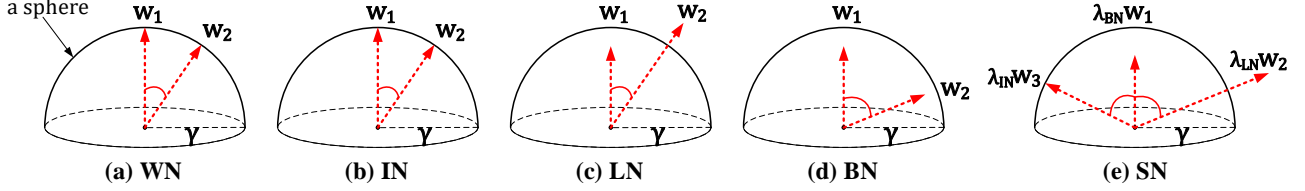
Fig. 2: Geometric view of directions and lengths of the filters in WN, IN, BN, LN, and SN. These normalizers are compared in an unify way by represent them by using WN that decomposes optimization of filters into their directions and lengths. In this way, IN is identical to WN that sets the filter norm to '1' (*i.e.* $\|\mathbf{w}_1\| = \|\mathbf{w}_2\| = 1$) and then rescales them to $\gamma$. LN is less constrained than IN and WN to increase learning ability. BN increases angle between filters and reduces filter length to improve generalization. SN inherits all their benefits by learning their importance ratios.

|  | (8,32) | (8,16) | (8,8) | (8,4) | (8,2) | (1,16) | (1,32) | (8,1) | (1,8) |
|---|---|---|---|---|---|---|---|---|---|
| BN | 76.4 | 76.3 | 75.2 | 72.7 | 65.3 | 76.2 | 76.5 | – | 75.4 |
| GN | 75.9 | 75.8 | 76.0 | 75.8 | **75.9** | 75.9 | 75.8 | **75.5** | 75.5 |
| SN | **76.9** | **76.7** | **76.7** | **75.9** | 75.6 | **76.3** | **76.6** | 75.0 | **75.9** |
| GN−BN | -0.5 | -0.5 | 0.8 | 3.1 | 10.6 | -0.3 | -0.7 | – | 0.1 |
| SN−BN | 0.5 | 0.4 | 1.5 | 3.2 | 10.3 | 0.1 | 0.1 | – | 0.5 |
| SN−GN | 1.0 | 0.9 | 0.7 | 0.1 | -0.3 | 0.4 | 0.8 | -0.5 | 0.4 |

TABLE 2: **Comparisons of top-1 accuracies** on the validation set of ImageNet, by using ResNet50 trained with SN, BN, and GN in different batch size settings. The bracket $(\cdot, \cdot)$ denotes (#GPUs, #samples per GPU). In the bottom part, 'GN-BN' indicates the difference between the accuracies of GN and BN. The '-' in $(8,1)$ indicates BN does not converge. The best-performing result of each setting is shown in bold.

focuses on SN where the importance weights are tied between channels.

**Inference.** When applying SN in test, the statistics of IN and LN are computed independently for each sample, while BN uses batch average *after* training without computing moving average in each iteration. Here batch average is performed in two steps. First, we freeze the parameters of the network and all the SN layers, and feed the network with a certain number of mini-batches randomly chosen from the training set. Second, we average the means and variances produced by all these mini-batches in each SN layer. The averaged statistics are used by BN in SN.

We find that the batch average can also achieve good performance by using a small amount of samples to calculate the statistics. For example, top-1 accuracies of ResNet50 on ImageNet by using batch average with 50k and all training samples are 76.90% and 76.92%. They are trained more stable in the early stage and slightly better than 76.89% of the moving average.

## 4 EXPERIMENTS

This section presents the main results of SN in multiple challenging problems and benchmarks, such as ImageNet [8], COCO [9], Cityscapes [10], ADE20K [11], MegaFace [12]and Kinetics [13], where the effectiveness of SN is demonstrated by comparing with existing normalization techniques.

### 4.1 Image Classification in ImageNet

SN is first compared with existing normalizers on the ImageNet classification dataset of 1k categories. All models in ImageNet are trained on 1.2M images and evaluated on 50K validation images. In the following, we first overview the experimental setting in Sec 4.1.1. The performance comparisons with other popular normalization methods are reported in Sec 4.1.2. More detailed analysis about different factors and learning dynamics of ratios are presented in Sec 4.1.3 and Sec 4.1.4.

### 4.1.1 Experimental Setting

All the models are trained by using SGD with different settings of batch sizes, which are denoted as a 2-tuple, (*number of GPUs*, *number of samples per GPU*). For each setting, the gradients are aggregated over all GPUs, and the means and variances of the normalization methods are computed in each GPU. The network parameters are initialized by following [6]. For all normalization methods, all $\gamma$'s are initialized as 1 and all $\beta$'s as 0. The parameters of SN ($\lambda_k$ and $\lambda'_k$) are initialized as 1. We use a weight decay of $10^{-4}$ for all parameters including $\gamma$ and $\beta$. Without special explanation, all the methods adopt ResNet50 as backbone network and are trained for 100 epoches with a initial learning rate of 0.1, which is deceased by $10\times$ after 30, 60, and 90 epoches. For different batch sizes, the initial learning rate is linearly scaled according to [45]. During training, we employ data augmentation the same as [6]. The top-1 classification accuracy on the $224\times224$ center crop is reported.

### 4.1.2 Performance Comparisons

The top-1 accuracy on the $224\times224$ center crop is reported for all models. SN is compared to BN and GN as shown in Table 2. In the first five columns, we see that the accuracy of BN reduces by 1.1% from $(8, 16)$ to $(8, 8)$ and declines to 65.3% of $(8, 2)$, implying that BN is unsuitable in small minibatch, where the random noise from the statistics is too heavy. GN obtains around 75.9% in all cases, while SN outperforms BN and GN in almost all cases, rendering its robustness to different batch sizes. Fig.3 plots the training and validation curves, where SN enables faster convergence while maintains higher or comparable accuracies than those of BN and GN.

The middle two columns of Table 2 average the gradients in a single GPU by using only 16 and 32 samples, such that their batch sizes are the same as $(8, 2)$ and $(8, 4)$. SN again performs best in these single-GPU settings, while BN outperforms GN. For example, unlike $(8, 4)$ that uses 8 GPUs, BN achieves 76.5% in $(1, 32)$, which is the best-performing result of BN, although the

(a) validation curves of SN in different batch sizes.  (b) train and validation curves of (8,32).  (c) train and validation curves of (8,2).
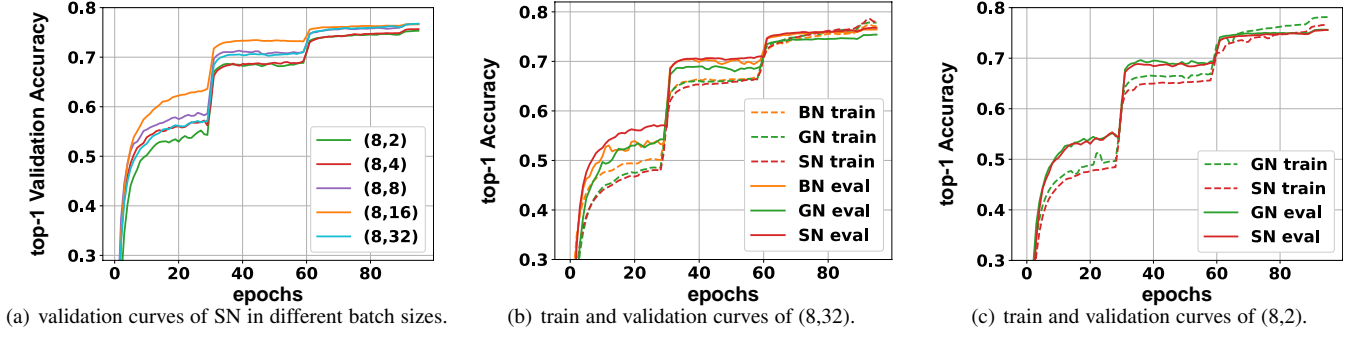
Fig. 3: **Comparisons of learning curves.** (a) visualizes the validation curves of SN with different settings of batch size. The bracket $(\cdot, \cdot)$ denotes (#GPUs, #samples per GPU). (b) compares the top-1 train and validation curves on ImageNet of SN, BN, and GN in the batch size of (8,32). (c) compares the train and validation curves of SN and GN in the batch size of (8,2).

batch size to compute the gradients is as small as 32. From the above results, we see that BN's performance are sensitive to the statistics more than the gradients, while SN are robust to both of them. The last two columns of Table 2 have the same batch size of 8, where $(1, 8)$ has a minibatch size of 8, while $(8, 1)$ is an extreme case with a single sample in a minibatch. For $(1, 8)$, SN performs best. For $(8, 1)$, SN consists of IN and LN but no BN, because IN and BN are the same in training when the minibatch size is 1. In this case, both SN and GN still perform reasonably well, while BN fails to converge.

Fig.3 (a) plots the validation curves of SN. Fig.3 (b) and (c) compare the training and validation curves of SN, BN and GN in $(8, 32)$ and $(8, 2)$ respectively. From all these curves, we see that SN enables faster convergence while maintains higher or comparable accuracies than those of BN and GN. On the other hand, SN prevents overfitting problem compared with GN when using both regular (*i.e.* 32 images) and small (*i.e.* 2 images) minibatch. It benefits from the combination of BN to introduce random noise to promote generalization ability.

In Fig.3, it is noteworthy that the evaluation accuracies are usually higher than training accuracies in the early stage. This phenomenon may be caused by two factors. (1) The training accuracy of each epoch is calculated by averaging the loss of each step within the corresponding epoch. In contrast, the evaluation accuracy is directly tested by using the model of the last step of each epoch. Since the values of accuracies are calculated in a slightly different way, in the early phase of the training, evaluation accuracies are usually higher than the training ones. (2) The data augmentation (e.g. random crop, multi-scale training, etc.) and dropout are adopted in the training phase but are omitted in the evaluation. In the early stage, these may also cause the difference between training and evaluation accuracies.

### 4.1.3 Ablation Study

Fig.1 (a) and Fig.4 plot histograms to compare the importance weights of SN with respect to different tasks and batch sizes. These histograms are computed by averaging the importance weights of all SN layers in a network. They show that SN adapts to various scenarios by changing its importance weights. For example, SN prefers BN when the minibatch is sufficiently large, while it selects LN instead when small minibatch is presented, as shown in the green and red bars of Fig.4. These results are in line with our analyses in Sec.3.1.
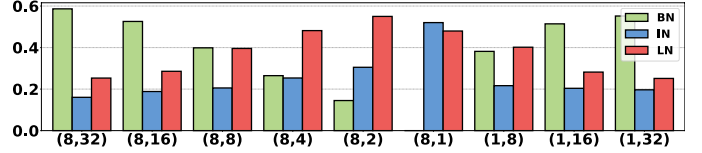


Fig. 4: Importance weights *v.s.* batch sizes. The bracket $(\cdot, \cdot)$ indicates (#GPUs, #samples per GPU). SN doesn't have BN in $(8, 1)$.

Furthermore, we repeat training of ResNet50 several times in ImageNet, to show that when the network, task, batch setting and data are fixed, the importance weights of SN are not sensitive to the change of training protocols such as solver, parameter initialization, and learning rate decay. As a result, we find that all trained models share similar importance weights.

The importance weights in each SN layer are visualized in Fig.5. We have several observations to answer what factors that impact the choices of normalizers. **First,** for the same batch size, the importance weights of $\mu$ and $\sigma$ could have notable differences, especially when comparing 'res1,4,5' of (a,b) and 'res2,4,5' of (c,d). For example, $\sigma$ of BN (green) in 'res5' in (b,d) are mostly reduced compared to $\mu$ of BN in (a,c). As discussed in [1], [17], this is because the variance estimated in a minibatch produces larger noise than the mean, making training instable. SN is able to restrain the noisy statistics and stabilize training. **Second,** the SN layers in different places of a network may select distinct operations. In other words, when comparing the adjacent SN layers after the $3 \times 3$ conv layer, shortcut, and the $1 \times 1$ conv layer, we see that they may choose different importance weights, *e.g.* 'res2,3'. The selectivity of operations in different places (normalization layers) of a deep network has not been observed in previous work. **Third,** deeper layers prefer LN and IN more than BN, as illustrated in 'res5', which tells us that putting BN in an appropriate place is crucial in the design of network architecture. Although the stochastic uncertainty in BN (*i.e.* the minibatch statistics) acts as a regularizer that might benefit generalization, using BN uniformly in all normalization layers may impede performance.

**Inference of SN.** In SN, BN employs batch average rather than moving average. We provide comparisons between them as shown in Fig.6, where SN is evaluated with both moving average and batch average to estimate the statistics used in test. They are used to train ResNet50 on ImageNet. The two settings of SN produce similar results of 76.9% when converged, which is better than 76.4% of BN. We see that SN with batch average converges more
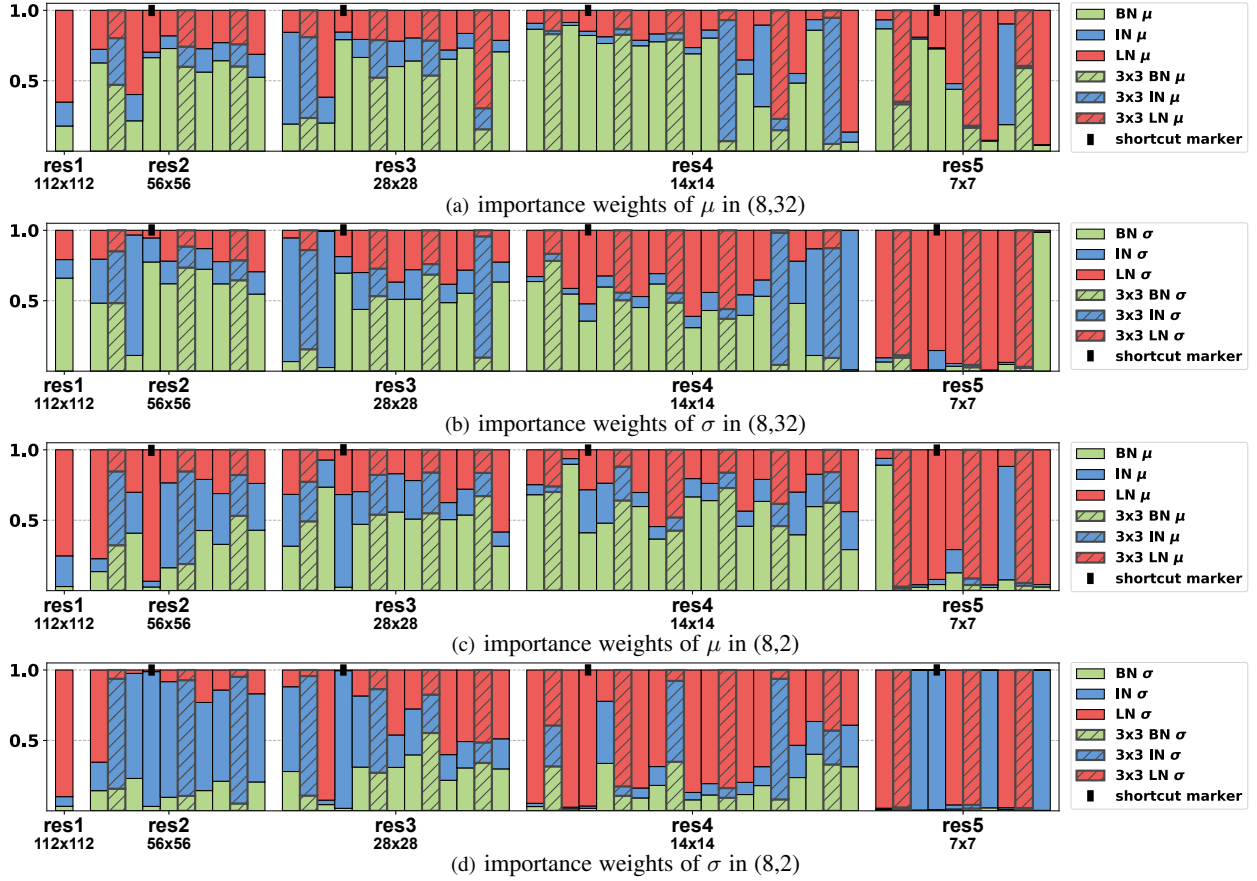
(a) importance weights of $\mu$ in $(8, 32)$

(b) importance weights of $\sigma$ in $(8, 32)$

(c) importance weights of $\mu$ in $(8, 2)$

(d) importance weights of $\sigma$ in $(8, 2)$

Fig. 5: **Selected operations of each SN layer in ResNet50.** There are 53 SN layers. (a,b) show the importance weights for $\mu$ and $\sigma$ of $(8, 32)$, while (c,d) show those of $(8, 2)$. The $y$-axis represents the importance weights that sum to 1, while the $x$-axis shows different residual blocks of ResNet50. The SN layers in different places are highlighted differently. For example, the SN layers follow the $3 \times 3$ conv layers are outlined by shaded color, those in the shortcuts are marked with '■', while those follow the $1 \times 1$ conv layers are in flat color. The first SN layer follows a $7 \times 7$ conv layer. We see that SN learns distinct importance weights for different normalization methods as well as $\mu$ and $\sigma$, adapting to different batch sizes, places, and depths of a deep network.

stably than BN and SN that use moving average. In this work, we also find that for all batch settings, SN with the batch average provides results better than the moving average.

**SN *v.s.* IN and LN.** IN and LN are not optimal in image classification as reported in [2] and [3]. With a regular setting of $(8, 32)$, ResNet50 trained with IN and LN achieve 71.6% and 74.7% respectively, which reduce 5.3% and 2.2% compared to 76.9% of SN.

**SN *v.s.* BRN and BKN.** BRN has two extra hyper-parameters, $r_{max}$ and $d_{max}$, which renormalize the means and variances. We choose their values as $r_{max} = 1.5$ and $d_{max} = 0.5$, which work best for ResNet50 in the setting of $(8, 4)$ following [15]. 73.7% of BRN surpasses 72.7% of BN by 1%, but it reduces 2.2% compared to 75.9% of SN.

BKN [16] estimated the statistics in the current layer by combining those computed in the preceding layers. It estimates the covariance matrix rather than the variance vector. In particular, how to connect the layers requires careful design for every specific network. For ResNet50 with $(8, 32)$, BKN achieved 76.8%, which is comparable to 76.9% of SN. However, for small minibatch, BKN reported 76.1% that was evaluated in a micro-batch setting where 256 samples are used to compute gradients and 4 samples to estimate the statistics. This setting is easier than $(8, 4)$ that uses 32 samples to compute gradients. Furthermore, it is unclear how to apply BRN and BKN in the other tasks such as object detection

| $(8, 32)$ | ResNet50 | ResNet101 |
|---|---|---|
| BN | 76.4 / 93.0 | 77.6 / 93.6 |
| SN | 76.9 / 93.2 | 78.6 / 94.1 |
| SN-BN | 0.5 / 0.2 | 1.0 / 0.5 |

| epoch | scratch | finetune |
|---|---|---|
| 30th | 76.5 / 93.0 | 77.4 / 93.4 |
| 60th | 76.1 / 93.0 | 77.1 / 93.4 |
| 90th | 76.2 / 93.0 | 76.9 / 93.3 |

TABLE 3: **Top:** comparisons of top1 / top5 accuracy of BN and SN trained with $(8, 32)$ in ImageNet. "SN-BN" is the difference between their results. **Bottom:** comparisons of top1/top5 accuracy between 'training from scratch' and 'finetuning' ResNet50+SN$(8, 32)$ with hard ratios until 100 epochs.

and segmentation.

### 4.1.4 Learning Dynamics of Ratios

**Soft ratios vary in training.** The values of soft ratios $\lambda_z^\mu$ and $\lambda_z^\sigma$ are between 0 and 1. Fig.7(a,b) plot their values for each normalization layer at every epoch. These values would have smooth fluctuation in training, implying that different epochs may have their own preference of normalizers. In general, we see that $\lambda_z^\mu$ mostly prefers BN, while $\lambda_z^\sigma$ prefers IN and BN when receptive field (RF) $<49$ and prefers LN when RF$>299$. Fig.7(c) shows the

discrepancy between (a) and (b) by computing a symmetry metric, that is, $\mathcal{D}(\lambda_z^\mu \| \lambda_z^\sigma) = \mathcal{KL}(\lambda_z^\mu \| \lambda_z^\sigma) + \mathcal{KL}(\lambda_z^\sigma \| \lambda_z^\mu)$ where $\mathcal{KL}(\cdot \| \cdot)$ is Kullback-Leibler divergence. Larger value of $\mathcal{D}(\lambda_z^\mu \| \lambda_z^\sigma)$ indicates larger discrepancy between the distributions of $\lambda_z^\mu$ and $\lambda_z^\sigma$.

For example, $\lambda_z^\mu$ and $\lambda_z^\sigma$ of the first layer choose different normalizers (see the first subfigure in (a,b)), making them had moderately large divergence $\mathcal{D}(\lambda_z^\mu \| \lambda_z^\sigma) \approx 1$ (see the first subfigure in (c)). Moreover, the ratios of the $2^{\text{nd}}$, $3^{\text{rd}}$, and $4^{\text{th}}$ layer are similar and they have small divergence $\mathcal{D}(\lambda_z^\mu \| \lambda_z^\sigma) < 0.5$. We also see that $\lambda_z^\mu$ and $\lambda_z^\sigma$ prefer different normalizers when RF is 199~299 and 299~427 where $\mathcal{D}(\lambda_z^\mu \| \lambda_z^\sigma) > 2$, as shown in the last row of (c). In conclusion, more than 50% number of layers have large divergence between $\lambda_z^\mu$ and $\lambda_z^\sigma$, confirming that different ratios should be learned for $\mu$ and $\sigma$.
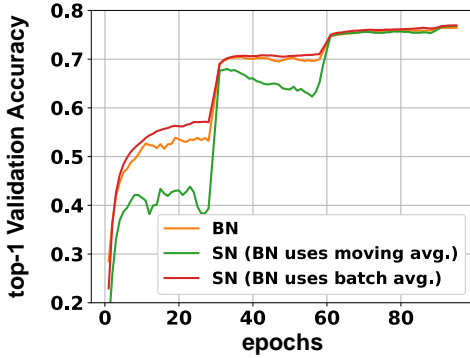


Fig. 6: Comparisons of 'BN', 'SN with moving average', and 'SN with batch average', when training ResNet50 on ImageNet in $(8, 32)$. We see that SN with batch average produces more stable convergence than the other methods.

**Hard ratios are relatively stable,** although the soft ratios are varying in training. A hard ratio is a sparse vector obtained by applying $\max$ function to $\lambda_z^\mu$ or $\lambda_z^\sigma$ such as $\max(\lambda_z^\sigma)$, that is, only one entry is '1' and the others are '0' to select only one normalizer. Fig.8 shows hard ratios for each layer in three snapshots, which are ResNet50+SN(8,32) trained after 30, 60, and 90 epochs respectively. For example, $\sigma$ and $\mu$ use LN and IN respectively in the first layer in Fig.8.

We have several observations. First, the size of filters ($3 \times 3$ or $1 \times 1$) seem to have no preference of specific normalizer, while the skip-connections prefer different normalizers for $\sigma$ and $\mu$ at the $90^{\text{th}}$ epoch. Second, around $50\%$ number of layers select two different normalizers for $\sigma$ and $\mu$ in these three snapshots, which are $49\%$ (26/53), $53\%$ (28/53), and $51\%$ (27/53) respectively. Third, the discrepancy between snapshots are small. For instance, 10 layers are different when comparing between $30^{\text{th}}$ and $60^{\text{th}}$ epoch, while only 2 layers are different between $60^{\text{th}}$ and $90^{\text{th}}$ epoch. Fourth, the layers that choose different normalizers are mainly presented when RF<40 and >200, rendering depth would be a major factor that affects the ratios.

**Performance of hard ratios.** We further examine performance of hard ratios. We finetune the above snapshots by replacing soft ratios with hard ratios. The right of Table 3 shows that these models achieve slightly better performance compared with their soft counterpart (76.9/93.3). An intuitive explanation is that sparseness can effectively prevent the model from overfitting. The similar results are also presented in the recent proposed Sparse Switchable Normalization (SSN) [46]. It implies that we could increase sparsity in ratios to reduce computations of multiple normalizers while maintaining good performance.

Furthermore, we train the above models from scratch by initializing the ratios as the hard ratios in the above snapshots, rather than using the default initial value $\frac{1}{3}$. However, this setting harms generalization as shown in Table 3. We conjecture that all normalizers in $\Omega$ are helpful to smooth loss landscape at the beginning of training. Therefore, the sparsity of ratios should be enhanced gently as training progresses. In other words, initializing ratios by $\frac{1}{3}$ is a good practice. Tuning this value is cumbersome and may also imped generalization ability.

## 4.2 Object Detection and Instance Segmentation

Next we evaluate SN in object detection and instance segmentation in COCO [9]. Unlike image classification, these two tasks benefit from large size of input images, making large memory footprint and therefore leading to small minibatch size, such as 2 samples per GPU [47], [48]. In this case, as BN is not applicable in small minibatch, previous work [47], [48], [49] often freeze BN and turns it into a constant linear transformation layer, which actually performs no normalization. Overall, SN selects different operations in different components of a detection system (see Fig.1), showing much more superiority than both BN and GN.

### 4.2.1 Experimental Settings

In practice, we implement object detection and instance segmentation on existing detection softwares of PyTorch and Caffe2-Detectron [50] respectively. We conduct 3 settings, including **setting-1:** Faster R-CNN [47] on PyTorch; **setting-2:** Faster R-CNN+FPN [48] on Caffe2; and **setting-3:** Mask R-CNN [49]+FPN on Caffe2. For all these settings, we choose ResNet50 as the backbone network. In each setting, the experimental configurations of all the models are the same, while only the normalization layers are replaced. All models of SN are finetuned from $(8, 2)$ in ImageNet.

For **setting-1**, we employ a fast implementation [51] of Faster R-CNN in PyTorch and follow its protocol. Specifically, we train all models on 4 GPUs and 3 images per GPU. Each image is re-scaled such that its shorter side is 600 pixels. All models are trained for 80k iterations with a learning rate of 0.01 and then for another 40k iterations with 0.001. For **setting-2** and **setting-3**, we employ the configurations of the Caffe2-Detectron [50]. We train all models on 8 GPUs and 2 images per GPU. Each image is re-scaled to its shorter side of 800 pixels. In particular, for setting-2, the learning rate (LR) is initialized as 0.02 and is decreased by a factor of 0.1 after 60k and 80k iterations and finally terminates at 90k iterations. This is referred as the 1x schedule in Detectron. In setting-3, the LR schedule is twice as long as the 1x schedule with the LR decay points scaled twofold proportionally, referred as 2x schedule. For all settings, we set weight decay to 0 for both $\gamma$ and $\beta$ following [5].

All of the models are trained in the *2017 train* set of COCO by using SGD with a momentum of 0.9 and a weight decay of $10^{-4}$ on the network parameters, and tested in the *2017 val* set. We report the standard metrics of COCO, including average precisions at IoU=0.5:0.05:0.75 (AP), IoU=0.5 (AP$_{.5}$), and IoU=0.75 (AP$_{.75}$) for both bounding box (AP$^{\text{b}}$) and segmentation mask (AP$^{\text{m}}$). Also, we report average precisions for small (AP$_s$), medium (AP$_m$), and large (AP$_l$) objects.

Fig. 7: Ratios of (a) $\lambda_z^\mu$ and (b) $\lambda_z^\sigma$ in ResNet50+SN(8,32) for each normalization layer for 100 epochs, as well as (c) their divergence $\mathcal{D}(\lambda_z^\mu \| \lambda_z^\sigma)$. Receptive field (RF) of each layer is given (53 normalization layers in total). The last 6 subfigures at the $4^{\text{th}}$, $8^{\text{th}}$, and $12^{\text{th}}$ row show results of different ranges of RF including 'RF<49', '49~99', '99~199', '199~299', '299~427', and 'ALL' (i.e. 7~427).
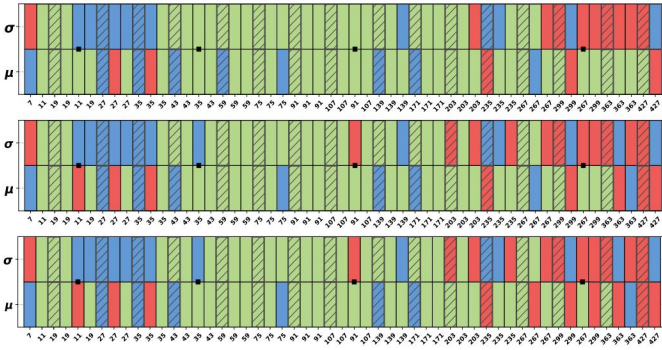


Fig. 8: **Hard ratios** for variance ($\sigma$) and mean ($\mu$) including BN (**green**), IN (**blue**), and LN (**red**). Snapshots of ResNet50 trained after 30 (**top**), 60 (**middle**), and 90 (**bottom**) epochs are shown. RF is given for each layer (53 normalization layers in total). A bar with slashes denotes SN after $3 \times 3$ conv layer (the others are $1 \times 1$ conv). A black square '■' indicates SN at the shortcut. It's better to zoom in 200%.

| backbone | head | AP | AP$_{.5}$ | AP$_{.75}$ | AP$_l$ | AP$_m$ | AP$_s$ |
|---|---|---|---|---|---|---|---|
| BN$^\dagger$ | BN$^\dagger$ | 29.6 | 47.8 | 31.9 | 45.5 | 33.0 | 11.5 |
| BN | BN | 19.3 | 33.0 | 20.0 | 32.3 | 21.3 | 7.4 |
| GN | GN | 32.7 | 52.4 | 35.1 | **49.1** | 36.1 | 14.9 |
| SN | SN | **33.0** | **52.9** | **35.7** | 48.7 | **37.2** | **15.6** |
| BN$^\ddagger$ | BN | 20.0 | 33.5 | 21.1 | 32.1 | 21.9 | 7.3 |
| GN$^\ddagger$ | GN | 28.3 | 46.3 | 30.1 | 41.2 | 30.0 | 12.7 |
| SN$^\ddagger$ | SN | **29.5** | **47.8** | **31.6** | **44.2** | **32.6** | **13.0** |

TABLE 4: **Faster R-CNN for detection in COCO** using ResNet50 and RPN. BN$^\dagger$ represents BN is frozen without finetuning. The superscript '$\ddagger$' indicates the backbones are trained from scratch without pretraining on ImageNet. The best results in the upper and lower parts are bold.

### 4.2.2 Performance Comparisons

**Results of Faster R-CNN.** As shown in Table 4, SN is compared with both BN and GN in the Faster R-CNN. In this setting, the layers up to conv4 of ResNet50 are used as *backbone* to extract features, and the layers of conv5 are used as the Region-of-Interest *head* for classification and regression. As the layers are inherited from the pretrained model, both the backbone and head involve normalization layers. Different results of Table 4 use different normalization methods in the backbone and head. Its upper part shows results of finetuning the ResNet50 models pretrained on ImageNet. The lower part compares training COCO from scratch without pretraining on ImageNet.

In the upper part of Table 4, the baseline is denoted as BN$^\dagger$, where the BN layers are frozen. We see that freezing BN performs significantly better than finetuning BN (29.6 v.s. 19.3). SN and GN enable finetuning the normalization layers, where SN obtains the best-performing AP of 33.0 in this setting. Fig.9 (a) compares their AP curves. As reported in the lower part of Table 4, SN and GN allow us to train COCO from scratch without pretraining on ImageNet, and they still achieve competitive results. For instance, 29.5 of SN$^\ddagger$ outperforms BN$^\ddagger$ by a large margin of 9.5 AP and GN$^\ddagger$ by 1.2 AP. Their learning curves are compared in Fig.9 (b).
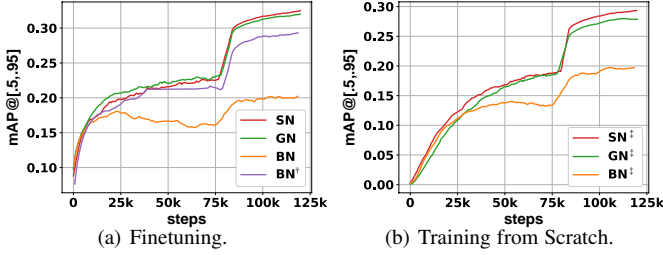
Fig. 9: Average precision (AP) curves of Faster R-CNN on the *2017 val* set of COCO. (a) plots the results of finetuning pretrained networks. (b) shows training the models from scratch.

| backbone | head | AP | AP$_{.5}$ | AP$_{.75}$ | AP$_l$ | AP$_m$ | AP$_s$ |
|---|---|---|---|---|---|---|---|
| BN$^\dagger$ | – | 36.7 | 58.4 | 39.6 | 48.1 | 39.8 | 21.1 |
| BN$^\dagger$ | GN | 37.2 | 58.0 | 40.4 | 48.6 | 40.3 | 21.6 |
| BN$^\dagger$ | SN | 38.0 | 59.4 | 41.5 | 48.9 | 41.3 | 22.7 |
| GN | GN | 38.2 | 58.7 | 41.3 | 49.6 | 41.0 | 22.4 |
| SN | SN | **39.3** | **60.9** | **42.8** | **50.3** | **42.7** | **23.5** |

TABLE 5: **Faster R-CNN+FPN** using ResNet50 and FPN with 1x LR schedule. BN$^\dagger$ represents BN is frozen. The best results are bold.

**Results of Faster R-CNN + FPN.** Table 5 reports results of Faster R-CNN by using ResNet50 and the Feature Pyramid Network (FPN) [48]. A baseline BN$^\dagger$ achieves an AP of 36.7 without using normalization in the detection head. When using SN and GN in the head and BN$^\dagger$ in the backbone, BN$^\dagger$+SN improves the AP of BN$^\dagger$+GN by 0.8 (from 37.2 to 38.0). We investigate using SN and GN in both the backbone and head. In this case, we find that GN improves BN$^\dagger$+SN by only a small margin of 0.2 AP (38.2 *v.s.* 38.0), although the backbone is pretrained and finetuned by using GN. When finetuning the SN backbone, SN obtains a significant improvement of 1.1 AP over GN (39.3 *v.s.* 38.2). Furthermore, the 39.3 AP of SN and 38.2 of GN both outperform 37.8 in [52], which synchronizes BN layers in the backbone (*i.e.* BN layers are not frozen).

**Results of Mask R-CNN + FPN.** Table 6 reports results of Mask R-CNN [49] with FPN. In the upper part, SN is compared to a head with no normalization and a head with GN, while the backbone is pretrained with BN, which is then frozen in finetuning (*i.e.* the ImageNet pretrained features are the same). We see that the baseline BN$^\dagger$ achieves a box AP of 38.6 and a mask AP of 34.2. SN improves GN by 0.5 box AP and 0.4 mask AP, when finetuning the same BN$^\dagger$ backbone.

More direct comparisons with GN are shown in the lower part of Table 6. We apply SN in the head and finetune the same backbone network pretrained with GN. In this case, SN outperforms GN by 0.2 and 0.3 box and mask APs respectively. Moreover, when finetuning the SN backbone, SN surpasses GN by a large margin of both box and mask AP (41.0 *v.s.* 40.2 and 36.5 *v.s.* 35.7). Note that the performance of SN even outperforms 40.9 and 36.4 of the 101-layered ResNet [50].

**Analysis of Dynamics of Ratios.** For object detection in COCO, we find that the ratio of BN gradually increases when 49<RF<299 (>0.5). This could be attributed to the two-stage pipeline of Mask R-CNN. To see this, Fig.10(b) and Fig.11 plot the ratios in COCO. We see that BN has larger impact in backbone and box stream than in FPN and mask stream. It demonstrates the regularization effect of BN could change depending on the

| backbone | head | AP$^b$ | AP$^b_{.5}$ | AP$^b_{.75}$ | AP$^m$ | AP$^m_{.5}$ | AP$^m_{.75}$ |
|---|---|---|---|---|---|---|---|
| BN$^\dagger$ | – | 38.6 | 59.5 | 41.9 | 34.2 | 56.2 | 36.1 |
| BN$^\dagger$ | GN | 39.5 | 60.0 | 43.2 | 34.4 | 56.4 | 36.3 |
| BN$^\dagger$ | SN | 40.0 | 61.0 | 43.3 | 34.8 | 57.3 | 36.3 |
| GN | GN | 40.2 | 60.9 | 43.8 | 35.7 | 57.8 | 38.0 |
| GN | SN | 40.4 | 61.4 | 44.2 | 36.0 | 58.4 | 38.1 |
| SN | SN | **41.0** | **62.3** | **45.1** | **36.5** | **58.9** | **38.7** |

TABLE 6: **Mask R-CNN** using ResNet50 and FPN with 2x LR schedule. BN$^\dagger$ represents BN is frozen without finetuning. The best results are bold.

tasks (*e.g.* bounding-box based or mask based) and network architectures design. Furthermore, ratios in the backbone have different dynamics in pretraining and finetuning by comparing two ResNet50 models in Fig.10(a,b), that is, the BN ratios decrease in pretraining for recognition, while increase in finetuning for detection even though the batch size is $(8, 2)$.

### 4.3 Semantic Image Segmentation

We investigate SN in semantic image segmentation in ADE20K [11] and Cityscapes [10]. Similar to object detection, semantic image segmentation also benefits from large input size, making the minibatch size is small during training. We use 2 samples per GPU for both of the datasets. We employ the open-source software in PyTorch[2] and only replace the normalization layers in CNNs with the other settings fixed.

#### 4.3.1 Experimental Settings

For both datasets, we use DeepLab [53] with ResNet50 as the backbone network, where $\text{output\_stride} = 8$ and the last two blocks in the original ResNet contains atrous convolution with $\text{rate} = 2$ and $\text{rate} = 4$ respectively. Following [54], we employ "poly" learning rate policy with $\text{power} = 0.9$ and use the auxiliary loss with the weight $0.4$ during training. The bilinear operation is adopted to upsmaple the score maps in the validation phase.

**ADE20K.** SyncBN and GN adopt the pretrained models on ImageNet. SyncBN collects the statistics from 8 GPUs. Thus the actual "batchsize" is 16 during training. To evaluate the performance of SN and SyncSN, we use SN $(8, 2)$, $(8, 4)$, $(8, 32)$ in ImageNet as the pretrained model, respectively. For all models, we resize each image to $450 \times 450$ and train for $100,000$ iterations. We performance multi-scale testing with $\text{input\_size} = \{300, 400, 500, 600\}$.

**Cityscapes.** For all models, we finetune from their pretrained ResNet50 models. Same as ADE20K, both SN and SyncSN finetune from $(8, 2)$, $(8, 4)$, $(8, 32)$ repsectively. For all models, the batchsize is 16 in finetuning. We use random crop with the size $713 \times 713$ and train for $400$ epoches. For multi-scale testing, the inference scales are $\{1.0, 1.25, 1.5, 1.75\}$.

#### 4.3.2 Performance Comparisons

Table 7 reports mIoU on the ADE20K validation set and Cityscapes test set, by using both single-scale and multi-scale testing. In SN, BN is not synchronized across GPUs, while it is synchronized in SyncSN. In ADE20K, the best model of SN outperforms SyncBN with a large margin in both testing schemes
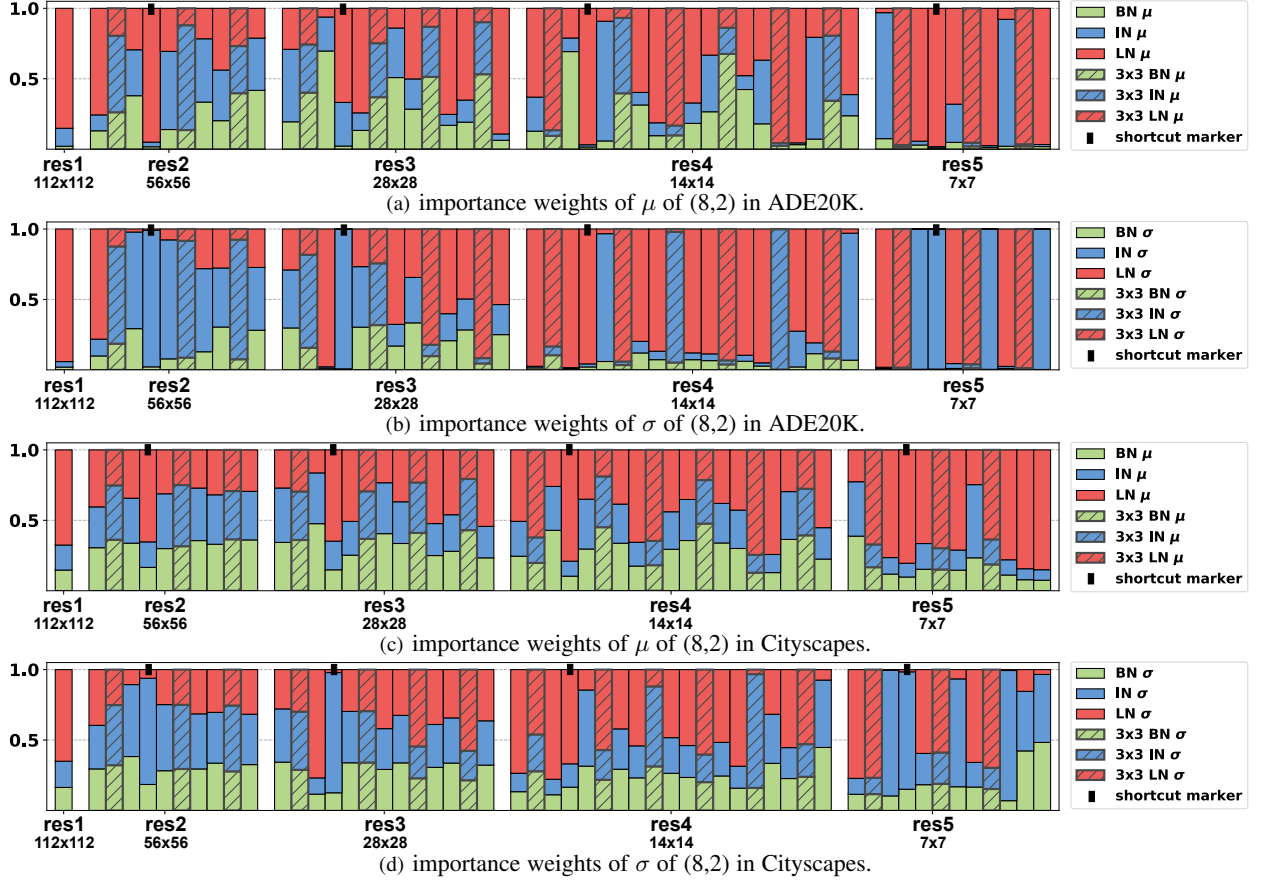
2. https://github.com/CSAILVision/semantic-segmentation-pytorch

Fig. 10: (a) shows ratios of ResNet50+SN(8, 2) in ImageNet. (b) shows ratios of Mask R-CNN+SN(8, 2) when finetuning in COCO including backbone (ResNet50), FPN, box stream, and mask stream.



Fig. 11: $\lambda_z^\mu$ (top) and $\lambda_z^\sigma$ (bottom) of Mask R-CNN+SN(8,2) are shown when training converged in COCO, including backbone (ResNet50), FPN, box stream, and mask stream.

(38.7 *v.s.* 36.4 and 39.2 *v.s.* 37.7), and improve GN ( *i.e.* channels in each group is 32 ) by 3.0 and 2.9. By using SyncSN, the best performance has been further improved to 40.0 for single-scale and 40.4 for multi-scale test. In Cityscapes, SN also performs better than SyncBN and GN. For example, the best model of SN surpasses SyncBN by 2.5 and 2.8 in both testing scales. SyncSN further improves the margins to 7.0 and 4.5. We see that GN performs worse than SyncBN in these two benchmarks. Note that SyncSN and SN share the same pre-trained models for this task. It means that we only use SyncSN to finetune the model on both of above datasets. Fig.12 compares the importance weights of SN (8, 2) in ResNet50 trained on both ADE20K and Cityscapes, showing that different datasets would choose different normalizers when the models and tasks are the same.

**Finetuning appropriate pretrained models.** For semantic image segmentation tasks, we observe that ratios pretrained with different batch sizes bring different impacts in finetuning. Using

models that are pretrained and finetuned with comparable batch size would be a best practice for good results. Otherwise, performance may degenerate.

In ADE20K, SN(8, 2) performs significantly better than the others. In this dataset, SN(8, 32) and SN(8, 4) may reduce performance, because BN has large ratio in these models, implying that finetuning SN pretrained with large batch to small batch would be unstable. Fig.14 shows the ratios when SN(8, 32) is used in pretraining but SN(8, 2) in finetuning. In line with expectation, the BN ratios are suppressed during finetuning, as the batch statistics become unstable. However, these ratios are still suboptimal until training converged, where the BN ratios finetuned from SN(8, 32) are still larger than those directly finetuned from SN(8, 2), reducing performance in ADE20K. When adopting SyncSN, the actual "batchsize" is 16 during training. SyncSN(8, 32) becomes more stable in the finetune phase and achieves best performance under the same setting.

(a) importance weights of $\mu$ of $(8,2)$ in ADE20K.

(b) importance weights of $\sigma$ of $(8,2)$ in ADE20K.

(c) importance weights of $\mu$ of $(8,2)$ in Cityscapes.

(d) importance weights of $\sigma$ of $(8,2)$ in Cityscapes.

Fig. 12: **Selected normalizers of each SN layer in ResNet50 for semantic image parsing in ADE20K and Cityscapes**. There are 53 SN layers. (a,b) show the importance weights for $\mu$ and $\sigma$ of $(8,2)$ in ADE20K, while (c,d) show those of $(8,2)$ in Cityscapes. The $y$-axis represents the importance weights that sum to 1, while the $x$-axis shows different residual blocks of ResNet50. The SN layers in different places are highlighted differently. For example, the SN layers follow the $3 \times 3$ conv layers are outlined by shaded color, those in the shortcuts are marked with '■', while those follow the $1 \times 1$ conv layers are in flat color.

|  | ADE20K | | Cityscapes | |
|---|---|---|---|---|
|  | mIoU$_{ss}$ | mIoU$_{ms}$ | mIoU$_{ss}$ | mIoU$_{ms}$ |
| SyncBN | 36.4 | 37.7 | 69.7 | 73.0 |
| GN | 35.7 | 36.3 | 68.4 | 73.1 |
| SN (8,2) | **38.7** | **39.2** | 71.6 | 75.4 |
| SN (8,4) | 38.6 | 39.0 | 72.1 | **75.8** |
| SN (8,32) | 37.7 | 38.4 | **72.2** | **75.8** |
| SyncSN (8,2) | 39.0 | 39.3 | 75.5 | 76.9 |
| SyncSN (8,4) | 39.8 | 39.9 | 75.8 | 77.0 |
| SyncSN (8,32) | **40.0** | **40.4** | **76.7** | **77.5** |

TABLE 7: **Results in ADE20K validation set and Cityscapes test set** by using ResNet50 with dilated convolutions. 'ss' and 'ms' indicate single-scale and multi-scale inference. SyncBN represents mutli-GPU synchronization of BN. SyncSN indicates the BN in SN is synchronized across mutli-GPU.

|  | SyncBN | GN | SN | | | SyncSN | | |
|---|---|---|---|---|---|---|---|---|
|  |  |  | (8,2) | (8,4) | (8,32) | (8,2) | (8,4) | (8,32) |
| mIoU$_{ss}$ | 76.3 | 72.6 | 75.9 | **77.1** | 76.2 | 76.0 | **77.1** | 76.7 |
| mIoU$_{ms}$ | 76.9 | 74.3 | 76.4 | 77.9 | 77.3 | 76.6 | **78.3** | 77.8 |

TABLE 8: **Results in Cityscapes validation set by using PSPNet** with ResNet50 as backbone network. 'ss' and 'ms' indicate single-scale and multi-scale inference, respectively.

Nevertheless, according to Table 7, SN(8, 32) and SN(8, 4) achieve better results than SN(8, 2) in Cityscapes. This could be attributed to large input image size 713×713 that diminishes noise in the batch statistics of BN (in SN). Same as ADE20K, SyncSN(8, 32) also outperforms the other synchronization models in this dataset.

**Analysis of Dynamics of Ratios.** Fig.13 visualizes ratios in finetuning for semantic image segmentation by using SN(8, 2),

which are more smooth than pretraining as compared to Fig.7. Intuitively, this is because a small learning rate is typically used in finetuning. In Fig.13, IN and LN ratios are generally larger than BN because of small batch size. The lower layers (RF<49) prefer IN more than the upper layers (RF>299) that choose LN. In addition, when comparing with detection in COCO, Fig.13(b,c) in segmentation have analogue dynamics where BN ratios are gently decreased (<0.3).

**Extension to PSPNet.** To further evaluate the performance of SN by using complex context modeling, we report the performance of PSPNet [54] adopting SN on Cityscapes validation set in Table 8. Both SN and SyncSN outperform GN with a margin. When compared with SyncBN, the mIoU of the best SN model (*i.e.* SN(8, 4)) is higher than SyncBN in both inference schemes (77.1 *v.s.* 76.3 and 77.9 *v.s.* 76.9). However, the ratio of improvement is less than using DeepLab as the backbone model.
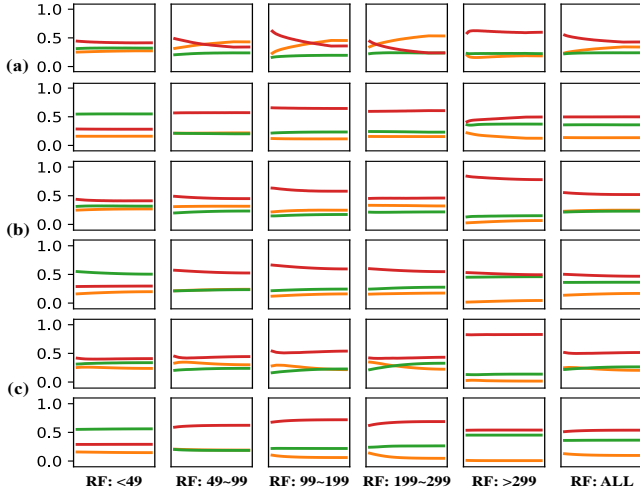
Fig. 13: **Ratios for detection and segmentation** including BN (**orange**), IN (**green**), and LN (**red**). We show $\lambda_z^\mu$ and $\lambda_z^\sigma$ in ResNet50+SN(8,2) finetuned to (a) COCO, (b) Cityscapes, and (c) ADE20K.

On one hand, this could be attributed to PSPNet providing a superior baseline compared with DeepLab. On the other hand, the spatial pyramid pooling (*i.e.* one type of multi-scale global pooling) may make IN and LN unstable after pooling operation. By using SyncSN, our best model (*i.e.* SyncSN$(8, 4)$) outperforms SyncBN by $0.8$ and $1.4$ in single and multiple scale test.

## 4.4 Face Recognition

We investigate face recognition task in MegaFace dataset [12]. We use MS1MV2 dataset [55] as our training data. For data processing, we follow [55], [56], [57] to generate the normalised face crop (*i.e.* $112$) by using five facial point. We adopt widely used CNN architecture, *i.e.* ResNet50, ResNet100 [6], as the backbone networks, and employ ArcFace [55] as the loss function. Following [55], the feature scale is set as $64$ and the angular margin of ArcFace is $0.5$. We set the batch size to $64$ for each GPU and train models on 8 GPUs. The learning rate is initialized as $0.1$ and divided by 10 at $[12, 15, 18]$ epoch. In training phase, we set momentum to $0.9$ and weight decay to $5 \times 10^{-4}$. For both SN and SyncSN, the pretrained model is trained on ImageNet with batch size 32 per GPU. During testing, we follow [55], [57] to keep the backbone network without the fully connected layer and extract the $512\text{-}d$ feature for each normalized face. We employ overlap list[3] to do dataset filtering to avoid the overlap ID in probe and gallery set, and then randomly select 1 million in remaining $1, 026, 351$ images to evaluate the performance of proposed SN.

Table 9 shows the performance of SN and SyncSN compared with BN and SyncBN. When using ResNet50 as the backbone architecture, proposed SN works better than BN no matter exploring synchronization version or not. For example, the verification result of SN is better than BN and SyncBN by $0.3$ and $0.4$. SyncSN further improves the margin to $0.5$ and $0.6$. By using ResNet101, SyncSN still achieves the best verification accuracy, and outperforms BN and SyncBN by $0.2$ and $0.3$, respectively.

---

3. https://github.com/deepinsight/insightface/tree/master/src/megaface

|  | BN | SyncBN | SN | SyncSN |
|---|---|---|---|---|
| ResNet50 | 95.3 | 95.2 | 95.6 | **95.8** |
| ResNet101 | 96.2 | 96.1 | 95.9 | **96.4** |

TABLE 9: **Verification Results of MegaFace dataset** by using different backbone architecture.

|  | batch=8, length=32 | | batch=4, length=32 | |
|---|---|---|---|---|
|  | top1 | top5 | top1 | top5 |
| BN | 73.2 | 90.9 | 72.1 | 90.0 |
| GN | 73.0 | 90.6 | 72.8 | 90.6 |
| SN | 73.5 | 91.2 | **73.3** | **91.2** |

TABLE 10: **Results of Kinetics dataset.** In training, the clip length of 32 frames is regularly sampled with a frame interval of 2. We study a batch size of 8 or 4 clips per GPU. BN is not synchronized across GPUs.

## 4.5 Video Recognition

We evaluate video recognition in Kinetics dataset [13], which has 400 action categories. We experiment with Inflated 3D (I3D) convolutional networks [58] and employ the ResNet50 I3D baseline as described in [5]. The models are pre-trained from ImageNet. For all normalizers, we extend the normalization from over $(H, W)$ to over $(T, H, W)$, where $T$ is the temporal axis. We train in the training set and evaluate in the validation set. The top1 and top5 classification accuracy are reported by using standard 10-clip testing that averages softmax scores from 10 clips sampled regularly.

Table 10 shows that SN works better than BN and GN in both batch sizes. For example, when batch size is 4, top1 accuracy of SN is better than BN and GN by 1.2% and 0.5%. It is seen that SN already surpasses BN and GN with batch size of 8.

## 4.6 Artistic Image Stylization

We also evaluate SN in the tasks of artistic image stylization [4]. We adopt a recent advanced approach [4], which jointly minimizes two loss functions. Specifically, one is a feature reconstruction loss that penalizes an output image when its content is deviated from a target image, and the other is a style reconstruction loss that penalizes differences in style (*e.g.* color, texture, exact boundary). [4], [42] show that IN works better than BN in this task.

We compare SN with IN and BN using VGG16 [59] as backbone network. All models are trained on the COCO dataset [9]. For each model in training, we resize each image to $256 \times 256$ and train for $40, 000$ iterations with a batch size setting of $(1, 4)$. We do not employ weight decay or dropout. The other training protocols are the same as [4]. In test, we evaluate the trained models on $512 \times 512$ images selected following [4]. Fig.16 (a) compares the style and feature reconstruction losses. We see that SN enables faster convergence than both IN and BN. As shown in Fig.1 (a), SN automatically selects IN in image stylization. Some stylization results are visualized in Fig.15.

## 4.7 Neural Architecture Search

We investigate SN in LSTM for efficient neural architecture search (ENAS) [14], which is designed to search the structures of convolutional cells. In ENAS, a convolutional neural network (CNN) is constructed by stacking multiple convolutional cells. It consists of two steps, training controllers and training child models. A

(a) $\lambda_z^\mu$ in **SN(8,2)** for ADE20K.

(b) $\lambda_z^\sigma$ in **SN(8,2)** for ADE20K.

(c) $\lambda_z^\mu$ in **SN(8,32)** for ADE20K.

(d) $\lambda_z^\sigma$ in **SN(8,32)** for ADE20K.

Fig. 14: Finetuning ResNet50+SN in ADE20K.

controller is a LSTM whose parameters are trained by using the REINFORCE [60] algorithm to sample a cell architecture, while a child model is a CNN that stacks many sampled cell architectures and its parameters are trained by back-propagation with SGD. In [14], the LSTM controller is learned to produce an architecture with high reward, which is the classification accuracy on the validation set of CIFAR-10 [61]. Higher accuracy indicates the controller produces better architecture.

We compare SN with LN and GN by using them in the LSTM controller to improve architecture search. As BN only achieves limited success in recurrent architecture [62] and IN is not applicable to LSTM when it is adopted to process non-image data (*i.e.* the statistics of IN are computed across height and width, which do not exist for non-image data), SN only combines LN and GN in this experiment. Fig.16 (b) shows the validation accuracy of CIFAR10. We see that SN obtains better accuracy than both LN and GN.

1. Initializing ratios of normalizers uniformly *e.g.* $1/3$. Carefully tuning the initial ratios may harm generalization.
2. Adding dropout with a small ratio (*e.g.* 0.1∼0.2) after each SN layer provides minor improvement of generalization in ImageNet. It reduces over-fitting.
3. Adding 0.5 dropout in the last fully-connected layer helps generalization in ImageNet.
4. A model in pretraining and finetuning should have comparable batch size.
5. Do not put SN after global pooling when feature map size is 1×1, because IN and LN are unstable after global pooling.
6. SN with hard ratios improves SN in ImageNet.
7. SN with hard ratios reduces computational runtime in inference compared to SN. 50% number of layers in sparse SN select BN for both $\mu$ and $\sigma$, meaning that these BN layers can be turned into linear transformation to reduce runtime in inference.
8. Synchronizing BN in SN improves generalization.

TABLE 11: Summary of practices that help training CNNs with SN.

## 5 DISCUSSIONS AND FUTURE WORK

This work presented Switchable Normalization (SN) to learn different operations in different normalization layers of a deep network. This novel perspective opens up new direction in many research fields that employ deep learning, such as CV, ML, NLP, Robotics, and Medical Imaging. This work has demonstrated SN in multiple tasks of CV such as recognition, detection, segmentation, image stylization, and neural architecture search, where SN outperforms previous normalizers without bells and whistles. Our analyses suggest that SN has an appealing characteristic to balance learning and generalization when training deep networks.
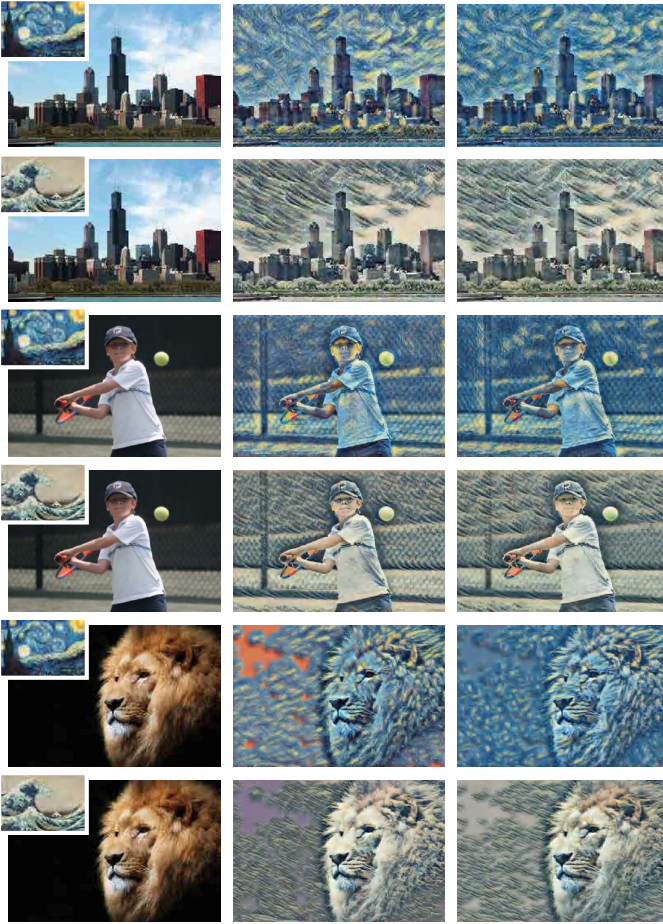
Fig. 15: **Results of Image Stylization.** The first column visualizes the content and the style images. The second and third columns are the results of IN and SN respectively. SN works comparably well with IN in this task.



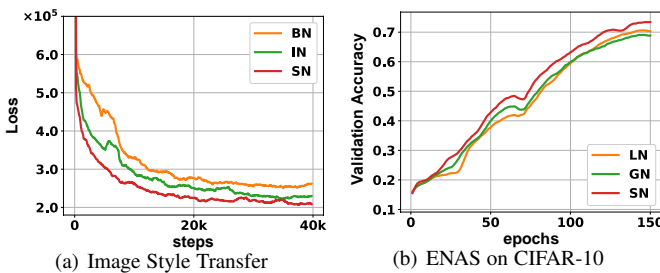(a) Image Style Transfer       (b) ENAS on CIFAR-10

Fig. 16: (a) shows the losses of BN, IN, and SN in the task of image stylization. SN converges faster than IN and BN. As shown in Fig.1 and the supplementary material, SN adapts its importance weight to IN while producing comparable stylization results. (b) plots the accuracy on the validation set of CIFAR-10 when searching network architectures.

Investigating SN facilitates the understanding of normalization approaches. To make better use of the proposed SN, we also summarize several important practices when learning to select normalizers in Table 11.

Future work will explore SN in more research fields as mentioned above. Moreover, constructing a theoretical framework to exploring the convergence and generalization ability of SN have

importance value and will be treated as future work either. As an extension of SN, it is valuable to design the algorithm to lean arbitrary normalization operations for different convolutional layers in a deep ConvNet.

## REFERENCES

[1] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *ICML*, 2015.

[2] D. Ulyanov, A. Vedaldi, and V. Lempitsky, "Instance normalization: The missing ingredient for fast stylization," *arXiv:1607.08022*, 2016.

[3] J. L. Ba, J. R. Kiros, and G. E. Hinton, "Layer normalization," *arXiv:1607.06450*, 2016.

[4] J. Johnson, A. Alahi, and L. Fei-Fei, "Perceptual losses for real-time style transfer and super-resolution," in *ECCV*, 2016.

[5] Y. Wu and K. He, "Group normalization," in *ECCV*, 2018.

[6] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016.

[7] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *CVPR*, 2009.

[8] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein *et al.*, "Imagenet large scale visual recognition challenge," *International Journal of Computer Vision*, vol. 115, no. 3, pp. 211–252, 2015.

[9] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *ECCV*, 2014.

[10] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele, "The cityscapes dataset for semantic urban scene understanding," in *CVPR*, 2016.

[11] B. Zhou, H. Zhao, X. Puig, S. Fidler, A. Barriuso, and A. Torralba, "Scene parsing through ADE20K dataset," in *CVPR*, 2017.

[12] I. Kemelmacher-Shlizerman, S. M. Seitz, D. Miller, and E. Brossard, "The megaface benchmark: 1 million faces for recognition at scale," in *CVPR*, 2016.

[13] W. Kay, J. Carreira, K. Simonyan, B. Zhang, C. Hillier, S. Vijayanarasimhan, F. Viola, T. Green, T. Back, P. Natsev *et al.*, "The kinetics human action video dataset," *arXiv:1705.06950*, 2017.

[14] H. Pham, M. Y. Guan, B. Zoph, Q. V. Le, and J. Dean, "Efficient neural architecture search via parameter sharing," *arXiv:1802.03268*, 2018.

[15] S. Ioffe, "Batch renormalization: Towards reducing minibatch dependence in batch-normalized models," in *NIPS*, 2017.

[16] G. Wang, J. Peng, P. Luo, X. Wang, and L. Lin, "Batch kalman normalization: Towards training deep neural networks with micro-batches," *NIPS*, 2018.

[17] T. Salimans and D. P. Kingma, "Weight normalization: A simple reparameterization to accelerate training of deep neural networks," in *NIPS*, 2016.

[18] T. Miyato, T. Kataoka, M. Koyama, and Y. Yoshida, "Spectral normalization for generative adversarial networks," in *ICLR*, 2018.

[19] G. Desjardins, K. Simonyan, R. Pascanu, and K. Kavukcuoglu, "Natural neural networks," *NIPS*, 2015.

[20] P. Luo, "Learning deep architectures via generalized whitened neural networks," *ICML*, 2017.

[21] S.-i. Amari and H. Nagaoka, *Methods of information geometry*. American Mathematical Soc., 2007, vol. 191.

[22] D. Balduzzi, M. Frean, L. Leary, J. Lewis, K. W.-D. Ma, and B. McWilliams, "The shattered gradients problem: If resnets are the answer, then what is the question?" *arXiv preprint arXiv:1702.08591*, 2017.

[23] S. Santurkar, D. Tsipras, A. Ilyas, and A. Madry, "How does batch normalization help optimization?(no, it is not about internal covariate shift)," in *NIPS*, 2018.

[24] E. Hoffer, R. Banner, I. Golan, and D. Soudry, "Norm matters: efficient and accurate normalization schemes in deep networks," in *NIPS*, 2018.

[25] P. Luo, X. Wang, W. Shao, and Z. Peng, "Towards understanding regularization in batch normalization," *ICLR*, 2019.

[26] V. R. J. S.-D. S. S. S. Greg Yang, Jeffrey Pennington, "A mean field theory of batch normalization," in *ICLR*, 2019.

[27] K. L. Sanjeev Arora, Zhiyuan Li, "Theoretical analysis of auto rate-tuning by batch normalization," in *ICLR*, 2019.

[28] B. S. K. Q. W. Johan Bjorck, Carla Gomes, "Understanding batch normalization," in *NIPS*, 2018.

[29] M. Ren, R. Liao, R. Urtasun, F. H. Sinz, and R. S. Zemel, "Normalizing the normalizers: Comparing and extending network normalization schemes," in *ICLR*, 2016.

[30] J. S. Xinggang Pan, Ping Luo and X. Tang, "Two at once: enhancing learning and generalization capacities via ibn-net," in *ECCV*, 2018.

[31] P. Luo, P. Zhanglin, S. Wenqi, Z. Ruimao, R. Jiamin, and W. Lingyun, "Differentiable dynamic normalization for learning deep representation," in *ICML*, 2019.

[32] Q. Liao, K. Kawaguchi, and T. Poggio, "Streaming normalization: Towards simpler and more biologically-plausible normalizations for online and recurrent learning," *arXiv preprint arXiv:1610.06160*, 2016.

[33] L. Huang, D. Yang, B. Lang, and J. Deng, "Decorrelated batch normalization," in *CVPR*, 2018.

[34] E. Perez, H. de Vries, and F. Strub, "Learning visual reasoning without strong priors," in *arXiv:1707.03017*, 2017.

[35] X. Pan, X. Zhan, J. Shi, X. Tang, and P. Luo, "Switchable whitening for deep representation learning," *arXiv*, 2019.

[36] L. Huang, X. Liu, Y. Liu, B. Lang, and D. Tao, "Centered weight normalization in accelerating training of deep neural networks," in *ICCV*, 2017.

[37] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. C. Courville, and Y. Bengio, "Generative adversarial nets," in *NIPS*, 2014.

[38] P. Luo, "Eigennet: Towards fast and structural learning of deep neural networks," *IJCAI*, 2017.

[39] B. Colson, P. Marcotte, and G. Savard, "An overview of bilevel optimization," *Annals of operations research*, 2007.

[40] D. Maclaurin, D. Duvenaud, and R. Adams, "Gradient-based hyperparameter optimization through reversible learning," *ICML*, 2015.

[41] H. Liu, K. Simonyan, and Y. Yang, "Darts: Differentiable architecture search," *arXiv:1806.09055*, 2018.

[42] X. Huang and S. Belongie, "Arbitrary style transfer in real-time with adaptive instance normalization," in *ICCV*, 2017.

[43] D. Ulyanov, A. Vedaldi, and V. Lempitsky, "Improved texture networks: Maximizing quality and diversity in feed-forward stylization and texture synthesis," in *CVPR*, 2017.

[44] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *NIPS*, 2012.

[45] P. Goyal, P. Dollr, R. Girshick, P. Noordhuis, L. Wesolowski, A. Kyrola, A. Tulloch, Y. Jia, and K. He, "Accurate, large minibatch sgd: Training imagenet in 1 hour," *arXiv:1706.02677*, 2017.

[46] W. Shao, T. Meng, J. Li, R. Zhang, Y. Li, X. Wang, and P. Luo, "Ssn: Learning sparse switchable normalization via sparsestmax," in *CVPR*, 2019.

[47] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *NIPS*, 2015.

[48] T.-Y. Lin, P. Dollra, R. Girshick, K. He, B. Hariharan, and S. Belongie, "Feature pyramid networks for object detection," *arXiv:1612.03144*, 2016.

[49] K. He, G. Gkioxari, P. Dollr, and R. Girshick, "Mask r-cnn," *ICCV*, 2017.

[50] R. Girshick, I. Radosavovic, G. Gkioxari, P. Dollár, and K. He, "Detectron," *https://github.com/facebookresearch/detectron*, 2018.

[51] J. Yang, J. Lu, D. Batra, and D. Parikh, "A faster pytorch implementation of faster r-cnn," *https://github.com/jwyang/faster-rcnn.pytorch*, 2017.

[52] C. Peng, T. Xiao, Z. Li, Y. Jiang, X. Zhang, K. Jia, G. Yu, and J. Sun, "Megdet: A large mini-batch object detector," in *CVPR*, 2018.

[53] L. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 40, no. 4, pp. 834–848, 2018.

[54] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia, "Pyramid scene parsing network," in *CVPR*, 2017.

[55] J. Deng, J. Guo, N. Xue, and S. Zafeiriou, "Arcface: Additive angular margin loss for deep face recognition," *arXiv preprint arXiv:1801.07698*, 2018.

[56] W. Liu, Y. Wen, Z. Yu, M. Li, B. Raj, and L. Song, "Sphereface: Deep hypersphere embedding for face recognition," in *CVPR*, 2017.

[57] H. Wang, Y. Wang, Z. Zhou, X. Ji, D. Gong, J. Zhou, Z. Li, and W. Liu, "Cosface: Large margin cosine loss for deep face recognition," in *CVPR*, 2018.

[58] J. Carreira and A. Zisserman, "Quo vadis, action recognition? a new model and the kinetics dataset," *arXiv:1705.07750*, 2017.

[59] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv:1409.1556*, 2014.

[60] R. J. Williams, "Simple statistical gradient-following algorithms for connectionist reinforcement learning," *Machine Learning*, 1992.

[61] A. Krizhevsky, "Learning multiple layers of features from tiny images," *Technical report*, 2009.

[62] C. Laurent, G. Pereyra, P. Brakel, Y. Zhang, and Y. Bengio, "Batch normalized recurrent neural networks," in *ICASSP*, 2016.

# APPENDIX A
# BACK-PROPAGATION OF SN

In practice, the back-propagation (BP) stage can be computed by auto differentiation (AD). For the software without AD, we provide the backward computations of SN for single GPU and multiple GPUs as below.

**Backward for Single GPU**. Let $\hat{h}$ be the output of the SN layer represented by a 4D tensor $(N, C, H, W)$ with index $n, c, i, j$. Let $\hat{h} = \gamma \tilde{h} + \beta$ and $\tilde{h} = \frac{h - \mu}{\sqrt{\sigma^2 + \epsilon}}$, where $\mu = w_{\mathrm{bn}} \mu_{\mathrm{bn}} + w_{\mathrm{in}} \mu_{\mathrm{in}} + w_{\mathrm{ln}} \mu_{\mathrm{ln}}$, $\sigma^2 = w_{\mathrm{bn}} \sigma_{\mathrm{bn}}^2 + w_{\mathrm{in}} \sigma_{\mathrm{in}}^2 + w_{\mathrm{ln}} \sigma_{\mathrm{ln}}^2$, and $w_{\mathrm{bn}} + w_{\mathrm{in}} + w_{\mathrm{ln}} = 1$. Note that the importance weights are shared among the means and variances for clarity of notations. Suppose that each one of $\{\mu, \mu_{\mathrm{bn}}, \mu_{\mathrm{in}}, \mu_{\mathrm{ln}}, \sigma^2, \sigma_{\mathrm{bn}}^2, \sigma_{\mathrm{in}}^2, \sigma_{\mathrm{ln}}^2\}$ is reshaped into a vector of $N \times C$ entries, which are the same as the dimension of IN's statistics. In the following, let $\mathcal{L}$ be the loss function and $\frac{\partial \mathcal{L}}{\partial \mu_n}$ be the gradient with respect to the $n$-th entry of $\mu$. We have,

$$\frac{\partial \mathcal{L}}{\partial \tilde{h}_{ncij}} = \frac{\partial \mathcal{L}}{\partial \hat{h}_{ncij}} \cdot \gamma, \tag{6}$$

$$\frac{\partial \mathcal{L}}{\partial \sigma^2} = -\frac{1}{2(\sigma^2 + \epsilon)} \sum_{i,j}^{H,W} \frac{\partial \mathcal{L}}{\partial \tilde{h}_{ncij}} \cdot \tilde{h}_{ncij}, \tag{7}$$

$$\frac{\partial \mathcal{L}}{\partial \mu} = -\frac{1}{\sqrt{\sigma^2 + \epsilon}} \sum_{i,j}^{H,W} \frac{\partial \mathcal{L}}{\partial \tilde{h}_{ncij}}, \tag{8}$$

$$
\begin{aligned}
\frac{\partial \mathcal{L}}{\partial h_{ncij}} &= \frac{\partial \mathcal{L}}{\partial \tilde{h}_{ncij}} \cdot \frac{1}{\sqrt{\sigma^2 + \epsilon}} \\
&+ \left[ \frac{2 w_{\mathrm{in}} (h_{ncij} - \mu_{\mathrm{in}})}{HW} \frac{\partial \mathcal{L}}{\partial \sigma^2} \right. \\
&+ \frac{2 w_{\mathrm{ln}} (h_{ncij} - \mu_{\mathrm{ln}})}{CHW} \sum_{c=1}^{C} \frac{\partial \mathcal{L}}{\partial \sigma_c^2} \\
&+ \left. \frac{2 w_{\mathrm{bn}} (h_{ncij} - \mu_{\mathrm{bn}})}{NHW} \sum_{n=1}^{N} \frac{\partial \mathcal{L}}{\partial \sigma_n^2} \right] \\
&+ \left[ \frac{w_{\mathrm{in}}}{HW} \frac{\partial \mathcal{L}}{\partial \mu} + \frac{w_{\mathrm{ln}}}{CHW} \sum_{c=1}^{C} \frac{\partial \mathcal{L}}{\partial \mu_c} \right. \\
&+ \left. \frac{w_{\mathrm{bn}}}{NHW} \sum_{n=1}^{N} \frac{\partial \mathcal{L}}{\partial \mu_n} \right],
\end{aligned}
\tag{9}
$$

The gradients for $\gamma$ and $\beta$ are

$$\frac{\partial \mathcal{L}}{\partial \gamma} = \sum_{n,c,i,j}^{N,C,H,W} \frac{\partial \mathcal{L}}{\partial \hat{h}_{ncij}} \cdot \tilde{h}_{ncij}, \tag{10}$$

$$\frac{\partial \mathcal{L}}{\partial \beta} = \sum_{n,c,i,j}^{N,C,H,W} \frac{\partial \mathcal{L}}{\partial \hat{h}_{ncij}}, \tag{11}$$

and the gradients for $\lambda_{\mathrm{in}}, \lambda_{\mathrm{ln}}$, and $\lambda_{\mathrm{bn}}$ are

$$
\begin{aligned}
\frac{\partial \mathcal{L}}{\partial \lambda_{\mathrm{in}}} &= w_{\mathrm{in}} (1 - w_{\mathrm{in}}) \sum_{n,c}^{N,C} \left( \frac{\partial \mathcal{L}}{\partial \mu_{nc}} \mu_{\mathrm{in}} + \frac{\partial \mathcal{L}}{\partial \sigma_{nc}^2} \sigma_{\mathrm{in}}^2 \right) \\
&- w_{\mathrm{in}} w_{\mathrm{ln}} \sum_{n,c}^{N,C} \left( \frac{\partial \mathcal{L}}{\partial \mu_{nc}} \mu_{\mathrm{ln}} + \frac{\partial \mathcal{L}}{\partial \sigma_{nc}^2} \sigma_{\mathrm{ln}}^2 \right) \\
&- w_{\mathrm{in}} w_{\mathrm{bn}} \sum_{n,c}^{N,C} \left( \frac{\partial \mathcal{L}}{\partial \mu_{nc}} \mu_{\mathrm{bn}} + \frac{\partial \mathcal{L}}{\partial \sigma_{nc}^2} \sigma_{\mathrm{bn}}^2 \right),
\end{aligned}
\tag{12}
$$

$$\frac{\partial \mathcal{L}}{\partial \lambda_{\ln}} = w_{\ln}(1 - w_{\ln}) \sum_{n,c}^{N,C} \left( \frac{\partial \mathcal{L}}{\partial \mu_{nc}} \mu_{\ln} + \frac{\partial \mathcal{L}}{\partial \sigma_{nc}^2} \sigma_{\ln}^2 \right)$$
$$- w_{\text{in}} w_{\ln} \sum_{n,c}^{N,C} \left( \frac{\partial \mathcal{L}}{\partial \mu_{nc}} \mu_{\text{in}} + \frac{\partial \mathcal{L}}{\partial \sigma_{nc}^2} \sigma_{\text{in}}^2 \right)$$
$$- w_{\ln} w_{\text{bn}} \sum_{n,c}^{N,C} \left( \frac{\partial \mathcal{L}}{\partial \mu_{nc}} \mu_{\text{bn}} + \frac{\partial \mathcal{L}}{\partial \sigma_{nc}^2} \sigma_{\text{bn}}^2 \right), \quad (13)$$

$$\frac{\partial \mathcal{L}}{\partial \lambda_{\text{bn}}} = w_{\text{bn}}(1 - w_{\text{bn}}) \sum_{n,c}^{N,C} \left( \frac{\partial \mathcal{L}}{\partial \mu_{nc}} \mu_{\text{bn}} + \frac{\partial \mathcal{L}}{\partial \sigma_{nc}^2} \sigma_{\text{bn}}^2 \right)$$
$$- w_{\text{in}} w_{\text{bn}} \sum_{n,c}^{N,C} \left( \frac{\partial \mathcal{L}}{\partial \mu_{nc}} \mu_{\text{in}} + \frac{\partial \mathcal{L}}{\partial \sigma_{nc}^2} \sigma_{\text{in}}^2 \right)$$
$$- w_{\ln} w_{\text{bn}} \sum_{n,c}^{N,C} \left( \frac{\partial \mathcal{L}}{\partial \mu_{nc}} \mu_{\ln} + \frac{\partial \mathcal{L}}{\partial \sigma_{nc}^2} \sigma_{\ln}^2 \right). \quad (14)$$

**Forward and Backward for Multiple GPUs**. Considering multi-GPU synchronization, let $p \in \{1, 2, ..., P\}$ denote the index of GPU, the mean and variance on $p$-th GPU can be denoted as $\mu_p = w_{\text{bn}} \mu_{\text{sybn}} + w_{\text{in}} \mu_{\text{in},p} + w_{\ln} \mu_{\ln,p}$ and $\sigma_p^2 = w_{\text{bn}} \sigma_{\text{sybn}}^2 + w_{\text{in}} \sigma_{\text{in},p}^2 + w_{\ln} \sigma_{\ln,p}^2$, where $\mu_{\text{sybn}}$ and $\sigma_{\text{sybn}}^2$ indicate the synchronized statistics of batch mean and batch variance on multiple GPUs. Let $\hat{h}_{p,ncij}$ indicate the output of the SN layer on the $p$-th GPU, we have

$$\hat{h}_{p,ncij} = \gamma \tilde{h}_{p,ncij} + \beta$$
$$\tilde{h}_{p,ncij} = \frac{h_{p,ncij} - \mu_p}{\sqrt{\sigma_p^2 + \epsilon}} \quad (15)$$

Let $\mathcal{L}$ be the loss function, the gradients for $\gamma$ and $\beta$ are,

$$\frac{\partial \mathcal{L}}{\partial \gamma} = \sum_{p,n,i,j}^{P,N,C,H,W} \frac{\partial \mathcal{L}}{\partial \hat{h}_{p,ncij}} \cdot \tilde{h}_{p,ncij}, \quad (16)$$

$$\frac{\partial \mathcal{L}}{\partial \beta} = \sum_{p,n,i,j}^{P,N,H,W} \frac{\partial \mathcal{L}}{\partial \hat{h}_{p,ncij}}, \quad (17)$$

and the gradients for $\lambda_{\text{in}}$, $\lambda_{\ln}$, and $\lambda_{\text{bn}}$ are

$$\frac{\partial \mathcal{L}}{\partial \lambda_{\text{in}}} = w_{\text{in}}(1 - w_{\text{in}}) \sum_{p,n,c}^{P,N,C} \left( \frac{\partial \mathcal{L}}{\partial \mu_{p,nc}} \mu_{\text{in},p} + \left( \frac{\partial \mathcal{L}}{\partial \sigma_{p,nc}^2} \sigma_{\text{in},p}^2 \right) \right)$$
$$- w_{\text{in}} w_{\ln} \sum_{p,n,c}^{P,N,C} \left( \frac{\partial \mathcal{L}}{\partial \mu_{p,nc}} \mu_{\ln,p} + \frac{\partial \mathcal{L}}{\partial \sigma_{p,nc}^2} \sigma_{\ln,p}^2 \right)$$
$$- w_{\text{in}} w_{\text{bn}} \sum_{p,n,c}^{P,N,C} \left( \frac{\partial \mathcal{L}}{\partial \mu_{p,nc}} \mu_{\text{sybn}} + \frac{\partial \mathcal{L}}{\partial \sigma_{p,nc}^2} \sigma_{\text{sybn}}^2 \right), \quad (18)$$

$$\frac{\partial \mathcal{L}}{\partial \lambda_{\ln}} = w_{\ln}(1 - w_{\ln}) \sum_{p,n,c}^{P,N,C} \left( \left( \frac{\partial \mathcal{L}}{\partial \mu_{p,nc}} \mu_{\ln,p} + \frac{\partial \mathcal{L}}{\partial \sigma_{p,nc}^2} \sigma_{\ln,p}^2 \right) \right)$$
$$- w_{\text{in}} w_{\ln} \sum_{p,n,c}^{P,N,C} \left( \frac{\partial \mathcal{L}}{\partial \mu_{p,nc}} \mu_{\text{in},p} + \frac{\partial \mathcal{L}}{\partial \sigma_{p,nc}^2} \sigma_{\text{in},p}^2 \right)$$
$$- w_{\ln} w_{\text{bn}} \sum_{p,n,c}^{P,N,C} \left( \left( \frac{\partial \mathcal{L}}{\partial \mu_{p,nc}} \mu_{\text{sybn}} + \frac{\partial \mathcal{L}}{\partial \sigma_{p,nc}^2} \sigma_{\text{sybn}}^2 \right) \right), \quad (19)$$

$$\frac{\partial \mathcal{L}}{\partial \lambda_{\text{bn}}} = w_{\text{bn}}(1 - w_{\text{bn}}) \sum_{p,n,c}^{P,N,C} \left( \frac{\partial \mathcal{L}}{\partial \mu_{p,nc}} \mu_{\text{sybn}} + \frac{\partial \mathcal{L}}{\partial \sigma_{p,nc}^2} \sigma_{\text{sybn}}^2 \right)$$
$$- w_{\text{in}} w_{\text{bn}} \sum_{p,n,c}^{P,N,C} \left( \frac{\partial \mathcal{L}}{\partial \mu_{p,nc}} \mu_{\text{in},p} + \left( \frac{\partial \mathcal{L}}{\partial \sigma_{p,nc}^2} \sigma_{\text{in},p}^2 \right) \right)$$
$$- w_{\ln} w_{\text{bn}} \sum_{p,n,c}^{P,N,C} \left( \frac{\partial \mathcal{L}}{\partial \mu_{p,nc}} \mu_{\ln,p} + \left( \frac{\partial \mathcal{L}}{\partial \sigma_{p,nc}^2} \sigma_{\ln,p}^2 \right) \right). \quad (20)$$

and the gradients for mean and variance on a specific GPU are

$$\frac{\partial \mathcal{L}}{\partial \sigma_{p,nc}^2} = -\frac{1}{2\sqrt{\sigma_{p,nc}^2 + \epsilon}} \sum_{i,j}^{H,W} \frac{\partial \mathcal{L}}{\partial \tilde{h}_{p,ncij}} \cdot \tilde{h}_{p,ncij}$$
$$- \frac{\gamma}{2\sqrt{\sigma_{p,nc}^2 + \epsilon}} \sum_{i,j}^{H,W} \frac{\partial \mathcal{L}}{\partial \hat{h}_{p,ncij}} \cdot \tilde{h}_{p,ncij},$$

$$\frac{\partial \mathcal{L}}{\partial \mu_{p,nc}} = -\frac{1}{\sqrt{\sigma_{p,nc}^2 + \epsilon}} \sum_{i,j}^{H,W} \frac{\partial \mathcal{L}}{\partial \tilde{h}_{p,ncij}}$$
$$- \frac{\gamma}{\sqrt{\sigma_{p,nc}^2 + \epsilon}} \sum_{i,j}^{H,W} \frac{\partial \mathcal{L}}{\partial \hat{h}_{p,ncij}}, \quad (21)$$

and the gradients respect to the input can be calculated by

$$\frac{\partial \mathcal{L}}{\partial h_{p,ncij}} = \underbrace{\frac{\partial \mathcal{L}}{\partial \tilde{h}_{p,ncij}} \cdot \frac{\partial \tilde{h}_{p,ncij}}{\partial h_{p,ncij}}}_{\text{Term1}} + \underbrace{\sum_{\ddot{p},\ddot{n},\ddot{c}}^{P,N,C} \frac{\partial \mathcal{L}}{\partial \sigma_{\ddot{p},\ddot{n}\ddot{c}}^2} \cdot \frac{\partial \sigma_{\ddot{p},\ddot{n}\ddot{c}}^2}{\partial h_{p,ncij}}}_{\text{Term2}}$$
$$+ \underbrace{\sum_{\ddot{p},\ddot{n},\ddot{c}}^{P,N,C} \frac{\partial \mathcal{L}}{\partial \mu_{\ddot{p},\ddot{n}\ddot{c}}} \cdot \frac{\partial \mu_{\ddot{p},\ddot{n}\ddot{c}}}{\partial h_{p,ncij}}}_{\text{Term3}}, \quad (22)$$

and,

$$\text{Term1} = \frac{\partial \mathcal{L}}{\partial \hat{h}_{p,ncij}} \cdot \frac{\gamma}{\sqrt{\sigma_{p,nc}^2 + \epsilon}}, \quad (23)$$

$$\text{Term2} = \sum_{\ddot{p},\ddot{n},\ddot{c}}^{P,N,C} \frac{\partial \mathcal{L}}{\partial \sigma_{\ddot{p},\ddot{n}\ddot{c}}^2} \cdot \left( w_{\text{in}} \frac{\partial \sigma_{\text{in},p,nc}^2}{\partial h_{p,ncij}} \delta_{p\ddot{p},n\ddot{n},c\ddot{c}} \right.$$
$$\left. + w_{\ln} \frac{\partial \sigma_{\ln,p,n}^2}{\partial h_{p,ncij}} \delta_{p\ddot{p},n\ddot{n}} + w_{\text{bn}} \frac{\partial \sigma_{\text{sybn},c}^2}{\partial h_{p,ncij}} \delta_{c\ddot{c}} \right)$$
$$= w_{\text{in}} \frac{\partial \sigma_{\text{in},p,nc}^2}{h_{p,ncij}} \frac{\partial \mathcal{L}}{\partial \sigma_{p,nc}^2} + w_{\ln} \frac{\partial \sigma_{\ln,p,n}^2}{h_{p,ncij}} \sum_{\ddot{c}}^{C} \frac{\partial \mathcal{L}}{\partial \sigma_{p,n\ddot{c}}^2}$$
$$+ w_{\text{bn}} \frac{\partial \sigma_{\text{sybn},c}^2}{h_{p,ncij}} \sum_{\ddot{p},\ddot{n}}^{P,N} \frac{\partial \mathcal{L}}{\partial \sigma_{\ddot{p},\ddot{n}c}^2}$$
$$= w_{\text{in}} \frac{2(h_{p,ncij} - \mu_{\text{in},p,nc})}{HW} \frac{\partial \mathcal{L}}{\partial \sigma_{p,nc}^2}$$
$$+ w_{\ln} \frac{2(h_{p,ncij} - \mu_{\ln,p,n})}{CHW} \sum_{\ddot{c}}^{C} \frac{\partial \mathcal{L}}{\partial \sigma_{p,n\ddot{c}}^2}$$
$$+ w_{\text{bn}} \frac{2(h_{p,ncij} - \mu_{\text{sybn},c})}{PNHW} \sum_{\ddot{p},\ddot{n}}^{P,N} \frac{\partial \mathcal{L}}{\partial \sigma_{\ddot{p},\ddot{n}c}^2} \quad (24)$$

$$
\begin{aligned}
\text{Term3} \;=\;& \sum_{\ddot{p},\ddot{n},\ddot{c}}^{P,N,C} \frac{\partial \mathcal{L}}{\partial \mu_{\ddot{p},\ddot{n}\ddot{c}}} \cdot \Big( w_{\text{in}} \frac{\partial \mu_{\text{in},p,nc}}{\partial h_{p,ncij}} \delta_{p\ddot{p},n\ddot{n},c\ddot{c}} \\
& + w_{\text{ln}} \frac{\partial \mu_{\text{ln},p,n}}{\partial h_{p,ncij}} \delta_{p\ddot{p},n\ddot{n}} + w_{\text{bn}} \frac{\partial \mu_{\text{sybn},c}}{\partial h_{p,ncij}} \delta_{c\ddot{c}} \Big) \\
=\;& w_{\text{in}} \frac{\partial \mu_{\text{in},p,nc}}{h_{p,ncij}} \frac{\partial \mathcal{L}}{\partial \mu_{p,nc}} + w_{\text{ln}} \frac{\partial \mu_{\text{ln},p,n}}{h_{p,ncij}} \sum_{\ddot{c}}^{C} \frac{\partial \mathcal{L}}{\partial \mu_{p,n\ddot{c}}} \\
& + w_{\text{bn}} \frac{\partial \mu_{\text{sybn},c}}{h_{p,ncij}} \sum_{\ddot{p},\ddot{n}}^{P,N} \frac{\partial \mathcal{L}}{\partial \mu_{\ddot{p},\ddot{n}c}} \\
=\;& w_{\text{in}} \frac{1}{HW} \frac{\partial \mathcal{L}}{\partial \mu_{p,nc}} + w_{\text{ln}} \frac{1}{CHW} \sum_{\ddot{c}}^{C} \frac{\partial \mathcal{L}}{\partial \mu_{p,n\ddot{c}}} \\
& + w_{\text{bn}} \frac{1}{PNHW} \sum_{\ddot{p},\ddot{n}}^{P,N} \frac{\partial \mathcal{L}}{\partial \mu_{\ddot{p},\ddot{n}c}} \qquad (25)
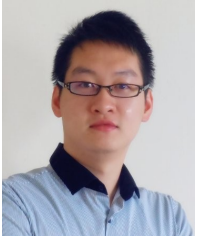\end{aligned}
$$

where $\delta$ indicates a Dirac Delta function that $\delta_{t\ddot{t}} = 1$ if $t = \ddot{t}$ and $\delta_{t\ddot{t}} = 0$ if $t \neq \ddot{t}$. Since we adopt synchronized statistics of batch mean and batch variance on multiple GPU, the information from the other GPUs also effect the output of specific GPU in the above formulas, we introduce $\ddot{p}, \ddot{n}$ and $\ddot{c}$ to help to describe the such interaction between different GPUs.

**Jiamin Ren** is currently a Researcher in Sense-Time Research. He received the B.E. and M.S. degrees from Harbin Institute of Technology(HIT), Harbin, China in 2014 and 2017, respectively. His research interests include computer vision, machine learning and robotics.

**Zhanglin Peng** is currently a Researcher in SenseTime Research. She received the B.E. and M.S. degrees from Sun Yat-sen University (SYSU), Guangzhou, China in 2013 and 2016, respectively. Her research interests include computer vision and deep learning.

**Ping Luo** is an Assistant Professor in the department of computer science, The University of Hong Kong (HKU). He received his PhD degree in 2014 from Information Engineering, the Chinese University of Hong Kong (CUHK), supervised by Prof. Xiaoou Tang and Prof. Xiaogang Wang. He was a Postdoctoral Fellow in CUHK from 2014 to 2016. He joined SenseTime Research as a Principal Research Scientist from 2017 to 2018. His research interests are machine learning and computer vision. He has published 70+ peer-reviewed articles in top-tier conferences and journals such as TPAMI, IJCV, ICML, ICLR, CVPR, and NIPS. His work has high impact with 7,000 citations according to Google Scholar. He has won a number of competitions and awards such as the first runner up in 2014 ImageNet ILSVRC Challenge, the first place in 2017 DAVIS Challenge on Video Object Segmentation, Gold medal in 2017 Youtube 8M Video Classification Challenge, the first place in 2018 Drivable Area Segmentation Challenge for Autonomous Driving, 2011 HK PhD Fellow Award, and 2013 Microsoft Research Fellow Award (ten PhDs in Asia).

**Jingyu Li** received the B.S. degree from The Chinese University of Hong Kong(CUHK), Hong Kong, China. He currently is a Ph.D. student in the Department of Electronic Engineering. In 2017, he joined Sensetime Research as intern where he worked on deep learning for video analysis and normalization algorithm. His research interests include computer vision , deep learning and multimedia data processing. He is a recipient of Hong Kong Postgraduate Fellowship.

**Ruimao Zhang** is currently a Senior Researcher in SenseTime Research. He received the B.E. and Ph.D. degrees from Sun Yat-sen University (SYSU), Guangzhou, China in 2011 and 2016, respectively. From 2017 to 2019, he was a Postdoctoral Research Fellow in the Department of Electronic Engineering, The Chinese University of Hong Kong (CUHK), Hong Kong, China. His research interests include computer vision, deep learning and related multimedia applications. He currently serves as a reviewer of numerous academic journals and conferences, including IJCV, T-NNLS, T-IP, T-CSVT, T-MM, CVPR, ICCV and IJCAI. He is a member of IEEE.