

Conditional Group Normalisation

BMVC 2020 Submission # 857

Abstract

Batch normalisation has been widely used to improve optimisation in deep neural networks. While the uncertainty in batch statistics can act as a regulariser, using these dataset statistics specific to the training set impairs generalisation in certain tasks. Recently, alternative methods for normalizing feature activations in neural networks have been proposed. Among them, group normalisation has been shown to yield competitive performance to batch normalisation. All these methods utilise a learned affine transformation after the normalisation operation to increase representational power. Methods used in conditional computation define the parameters of these transformations as learnable functions of conditioning information. In this work, we study how the conditional formulation of group normalisation can improve generalisation compared to conditional batch normalisation. We evaluate performances on the tasks of visual question answering and conditional image generation. Our experiments indicate that conditional group normalisation is a reasonable replacement for conditional batch normalisation, and can achieve improved performance in certain tasks.

1 Introduction

In machine learning, the parameters of a model are typically optimised using a fixed training set. The model is then evaluated on a separate partition of the data to estimate its generalisation capability. In practice, even under the i.i.d. assumption¹, the distribution of these two finite sets can *appear* quite different to the learning algorithm, making it challenging to achieve strong and robust generalisation. This difference is often the result of the fact that a training set of limited size cannot adequately cover the cross-product of all relevant factors of variation. In other cases, the i.i.d. assumption is dropped on purpose to study whether a model can capture regularity that generalises *out-of-distribution*. Several benchmarks for evaluating task-specific models for their generalisation capacity [6, 22, 23] have been proposed recently. The problem of out-of-distribution generalisation can in some cases be addressed by making strong assumptions that simplify discovering a family of patterns from limited data. Bahdanau et al. [5], for example, show that their proposed synthetic relational reasoning task can be solved by a Neural Module Network (NMN) [2] with fixed tree structure while models without this structural prior fail. At the other end of the spectrum are more “generic” models, as Feature-wise Linear Modulation (FiLM) [15], which are able to uncover some of the compositionality in the tasks they are trained for.

In this paper, we focus on such generic models that leverage conditional normalisation methods for applications in visual question answering (VQA) and generative modeling of images. Despite our focus on these tasks, any improvement in this area can also benefit other

© 2020. The copyright of this document resides with its authors.
It may be distributed unchanged freely in print or electronic forms.
¹All data samples are assumed to be drawn independently from an identical distribution (i.i.d.).

domains such as deep reinforcement learning or metalearning, where conditional normalisation layers are also being used [8, 21, 32]. We study strong deep neural network models for these tasks that employ Conditional Batch Normalisation (CBN) [12] for modulating normalised activations with contextual information.

Since Batch Normalisation (BN) normalises activations with statistics computed across multiple training samples, one has to precompute activation statistics over the training set to be used during inference. Due to this reliance on dataset statistics, it seems that BN [20] (and thus also CBN) may be vulnerable to significant domain shifts between training and test data. To train models with BN one has to use a sufficiently large mini-batch size to limit the noise of activation statistics. Further potential issues with BN include limited diversity of samples generated by Generative Adversarial Networks (GANs) involving BN [12] and vulnerability to adversarial examples [15].

Group Normalisation (GN) [10] normalises across groups of feature maps instead of across samples in mini-batches. Here, we explore whether a conditional formulation of GN is a viable alternative for CBN. GN is conceptually simpler than BN, as its function is the same during training and inference. Further, GN can be used with small batch sizes, which may help in applications with particularly large feature maps, such as medical imaging or video processing, in which the available memory can be a constraint.

Our contribution is an extensive empirical study of two conditional normalisation techniques over multiple tasks and benchmarks. Our experiments show that Conditional Group Normalisation (CGN) has advantages for out-of-distribution generalisation.

2 Background

2.1 Normalisation Layers

Several normalisation methods have been proposed to stabilise and speed-up the training of deep neural networks [9, 20, 39, 40]. To stabilise the range of variation of network activations x_i , methods such as BN [20] first normalise the activations by subtracting mean μ_i and dividing by standard deviation σ_i :

$$\hat{x}_i = \frac{1}{\sigma_i} (x_i - \mu_i) \quad (1)$$

The distinction between different methods lies in how exactly these statistics are computed. Wu and He [10] aptly summarise several methods using the following notation. Let $i = (i_N, i_C, i_H, i_W)$ be a four-dimensional vector, whose elements index the features along the batch, channel, height and width axes, respectively. The computation of the statistics can then be written as

$$\mu_i = \frac{1}{m} \sum_{k \in \mathcal{S}_i} x_k, \quad \sigma_i = \sqrt{\frac{1}{m} \sum_{k \in \mathcal{S}_i} (x_k - \mu_i)^2 + \varepsilon}, \quad (2)$$

where the set \mathcal{S}_i of size m is defined differently for each method and ε is a small constant for numerical stability. BN, for instance, corresponds to:

$$\mathcal{S}_i = \{k | k_C = i_C\}, \quad (3)$$

i.e. \mathcal{S}_i is the set of all pixels sharing the same channel axis, resulting in μ_i and σ_i being computed along the (N, H, W) axes.

As Ba et al. [4] point out, the performance of BN is highly affected by the batch size hyperparameter. This insight led to the introduction of several alternative normalisation schemes, that normalise per sample, i.e. not along batch axis N . For instance, Layer Normalisation (LN) [4], which normalises activations within each layer, corresponds to the following set definition:

$$\mathcal{S}_i = \{k | k_N = i_N\}. \quad (4)$$

Ulyanov et al. [19] introduce Instance Normalisation (IN) in the context of image stylisation. IN is motivated by the insight that BN with fixed training set statistics fails to remove instance-specific contrast information at test time. IN normalises separately for each sample and each channel along the spatial dimensions:

$$\mathcal{S}_i = \{k | k_N = i_N, k_C = i_C\}. \quad (5)$$

Drawing inspiration from classical features such as Histogram of Oriented Gradients (HOG) [10], Wu and He [44] proposed GN. It normalises features per sample, separately within each of G groups, along the channel axis:

$$\mathcal{S}_i = \{k | k_N = i_N, \lfloor \frac{k_C}{C/G} \rfloor = \lfloor \frac{i_C}{C/G} \rfloor\} \quad (6)$$

GN can be seen as a way to interpolate between the two extremes of LN (corresponding to $G = 1$, i.e. all channels are in a single group) and IN (corresponding to $G = C$, i.e. each channel is in its own group).

After normalisation, all above mentioned methods insert a scaling and shifting operation using learnable per-channel parameters γ and β :

$$y_i = \gamma \hat{x}_i + \beta \quad (7)$$

This “de-normalisation” is done to restore the representational power of the normalised network layer [40].

CBN [14, 35] is a conditional variant of BN, in which the learnable parameters γ and β in Equation 7 are replaced by learnable functions

$$\gamma(c_k) = W_\gamma c_k + b_\gamma, \quad \beta(c_k) = W_\beta c_k + b_\beta \quad (8)$$

of some per-sample conditioning input c_k to the network with parameters W_γ , W_β , b_γ , b_β . In a VQA model, c_k would for instance be an embedding of the question [35]. Before CBN was introduced, Dumoulin et al. [44] proposed Conditional Instance Normalisation (CIN), a conditional variant of IN very similar to CBN, using IN instead of BN. In our experiments, we explore a conditional variant of GN, i.e. CGN.

2.2 Visual Question Answering

In VQA [8, 76], the task is to answer a question about an image. This task is usually approached by feeding both image and question to a parametric model, which is trained to predict the correct answer, for instance via classification among all possible answers in the dataset. One recent successful model for VQA is the FiLM architecture [35], which employs CBN to modulate visual features based on an embedding of the question.

2.3 Conditional Image Generation

Some of the most successful models for generating images are GANs [16]. This approach involves training one neural network (Generator) to generate an image, while the only supervisory signal is that from another neural network (Discriminator) which indicates whether the image looks real or not. Several variants of GANs [27, 31] have been proposed to condition the image generation process on a class label. More recently, the generators that work best stack multiple ResNet-style [8] architectural blocks, involving two CBN-ReLU-Conv operations and an upsampling operation [8, 43]. These blocks are followed by a BN-ReLU-Conv operation to transform the last features into the shape of an image. Such models can be trained as Wasserstein GANs using gradient penalty (WGAN-GP) as proposed by Gulrajani et al. [17], which gives mathematically sound arguments for an optimisation framework.

More recently, Spectral Norm GAN (SNGAN) [28] uses the aforementioned architecture with spectral normalisation on the weights to stabilise training at each iteration. Two noteworthy GAN models that use architectures based on SNGAN are Self-Attention GAN (SAGAN) [44] and BigGAN [8]. SAGAN inserts a self-attention mechanism [9, 33, 40] to attend over important parts of features during the generation process. The architecture of BigGAN is the same as for SAGAN, with the exception of an increase in batch size and channel widths, as well as some architectural changes to improve memory and computational efficiency. Both of these models have been used successfully in generating high quality natural images. In our experiments, we use SNGAN and compare performance metrics of two types of normalisation — CBN and CGN.

3 Experiments

3.1 Visual Question Answering

We study whether CGN in the VQA architecture FiLM [35] yields performance improvements over CBN. We run experiments on the following recently proposed benchmarks for compositional generalisation:

CLEVR-CoGenT CLEVR Compositional Generalisation Test (CLEVR-CoGenT) [21] is a variant of the popular Compositional Language and Elementary Visual Reasoning (CLEVR) dataset [21] that tests for compositional generalisation. See Figure 1 (a) for an example from this dataset. The images consist of rendered three-dimensional scenes containing several shapes (small and large cubes, spheres and cylinders) of differing material properties (*metal* or *rubber*) and colors. Questions involve *queries* for object attributes, *comparisons*, *counting* of sets and combinations thereof. In contrast to the regular CLEVR dataset, the training set of CLEVR-CoGenT explicitly combines some shapes only with different subsets of four out of eight colors, and provides two validation sets: one with the same combinations (*valA*) and one in which the shape-color assignments are swapped (*valB*). To perform well on *valB*, the model has to generalise to unseen combinations of shapes and colors, i.e. it needs to capture the compositionality of the task.

SQOOP Spatial Queries On Object Pairs (SQOOP) [6] is a recently introduced dataset that tests for systematic generalisation. Figure 1 (b) shows an example from the training set. It consists of images containing five randomly chosen and arranged objects (digits and

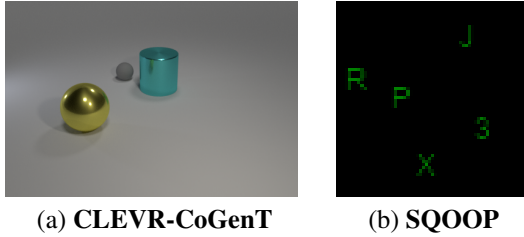


Figure 1: Examples of the VQA datasets used in our experiments. **Corresponding question-answer pairs:** (a) *Are there any gray things made of the same material as the big cyan cylinder?* - No. (b) *X right_of J?* - No

characters). Questions concern the four spatial relations *LEFT OF*, *RIGHT OF*, *ABOVE* and *BELOW* and the queries are all of the format "X R Y?", where X and Y are left-hand and right-hand objects and R is a relationship between them, e.g. "nine *LEFT OF* a?". To test for systematic generalisation, only a limited number of combinations of each left-hand object with different right-hand objects Y are shown during training. In the hardest version of the task (1 rhs/lhs), only a single right-hand side object is combined with each left-hand side object. For instance, the training set of this version may contain images with the query "a *RIGHT OF* b?", but no images with queries about relations of left-hand object *a* with any other object than *b*. The test set contains images and questions about all combinations, i.e. it evaluates generalisation to relations between novel object combinations.

3.1.1 Model

We experiment with three small variations of the FiLM architecture [5]. The original architecture in Perez et al. [5] consists of an unconditional *stem* network, followed by a core of four ResNet [18] blocks with CBN [10], and finally a classifier. The stem network is either a sequence of residual blocks trained from scratch or a fixed pre-trained feature extractor followed by a learnable layer of 3×3 convolutions. The scaling and shifting parameters of the core layers are affine transforms of a question embedding provided by a Gated Recurrent Unit (GRU) [11]. The output of the last residual block is fed to the classifier, which consists of a layer of $512 \ 1 \times 1$ convolutions, global max-pooling, followed by a fully-connected ReLU [6] layer using (unconditional) BN and a softmax layer, which outputs the probability of each possible answer. We train the following three variants that include CGN²:

1. "all GN": all conditional and regular BN layers are replaced with corresponding conditional or regular GN layers.
2. "BN stem": all CBN layers are replaced with CGN, regular BN layers in the stem and classifier are left unchanged, except those in the fully-connected hidden layer in the classifier, for which we remove normalisation.
3. "BN stem & cls": all CBN layers in the core ResNet blocks are replaced with CGN, regular BN in the stem and classifier are left unchanged.

²We always set the number of groups to 4, as the authors of Wu and He [10] showed that this hyperparameter does not have a large influence on the performance. This number was selected using uniform sampling from the set $\{2, 4, 8, 16\}$.

Table 1: Classification accuracy on CLEVR-CoGenT *valB*. Mean and standard deviation of three runs with early stopping on *valA* are reported for the models we trained.

Model	Accuracy (%)
CBN (FiLM [5])	75.60
CBN (FiLM, our results)	75.54 ± 0.67
CGN (all GN)	75.76 ± 0.36
CGN (BN stem)	75.70 ± 0.57
CGN (BN stem & cls)	75.81 ± 0.51

Table 2: Test accuracies on several versions of SQOOP. Mean and standard deviation of three runs after early stopping on the validation set are reported for the models we trained. Here, FiLM refers to the model specified in [6], whereas “FiLM, ours” indicates our run of the same.

Model	Accuracies (%)			
	1 rhs/lhs	2 rhs/lhs	4 rhs/lhs	35 rhs/lhs
CBN (FiLM)	65.27 ± 4.61	80.20 ± 4.32	90.42 ± 1.00	99.80 ± 0.22
CBN (FiLM, ours)	72.37 ± 0.53	84.97 ± 4.17	97.04 ± 1.96	99.84 ± 0.04
CGN (all GN)	74.02 ± 2.81	86.69 ± 6.31	91.40 ± 0.32	99.76 ± 0.03
CGN (BN stem)	73.82 ± 0.33	83.11 ± 0.38	91.60 ± 1.94	99.82 ± 0.12
CGN (BN stem & cls)	74.93 ± 3.89	85.86 ± 5.32	99.47 ± 0.25	99.78 ± 0.16

Besides the described changes in the normalisation layers, the architecture and hyperparameters are the same as used in Perez et al. [5] for all experiments, except for SQOOP where they are the same as in Bahdanau et al. [6]. The only difference is that we set the constant ϵ of the Adam optimiser [24] to $1e-5$ to improve training stability³. For SQOOP, the input to the residual network are the raw image pixels. For CLEVR-CoGenT, we instead feed features extracted from layer *conv4* of a ResNet-101 [18], pre-trained on ImageNet [56], following Perez et al. [5].

3.1.2 Results

Tables 1 and 2 show the results of training FiLM with CBN and CGN on the two considered datasets. In the experiments on CLEVR-CoGenT, all three CGN variants of FiLM achieve a slightly higher average accuracy. Note that for the SQOOP dataset we rerun the original CBN experiments by Bahdanau et al. [6] and observe significantly higher accuracy on all versions of the task. In the hardest SQOOP variant with only one right-hand side object per left-hand side object (*1 rhs/lhs*), all three variants of CGN achieve a higher performance than the CBN experiments (both original and ours). For the SQOOP variant with four right-hand side objects per left-hand side object, CGN did not converge in some cases. It is possible that additional regularisation is required to guarantee convergence. Note that learning curves of models successfully trained on SQOOP all seem to follow the same pattern: For a relatively large number of gradient updates there is no significant improvement. Then, at some point, almost instantly the model achieves 100% training accuracy.

³The authors of Perez et al. [5] confirmed occasional gradient explosions with the original setting of $1e-8$.

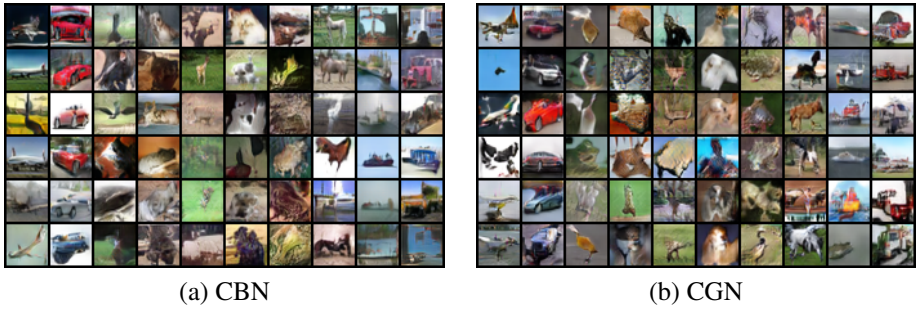


Figure 2: Samples from models trained with different normalisation techniques. The images in each column belong to the same class, ordered as ‘airplane’, ‘automobile’, ‘bird’, ‘cat’, ‘deer’, ‘dog’, ‘frog’, ‘horse’, ‘ship’, ‘truck’. Samples are not cherry-picked.

3.2 Conditional Image Generation

Here we compare CBN and CGN on the task of generating images conditioned on their class label using the SNGAN [28] architecture. We use the CIFAR-10 [25] dataset containing 60000 32×32 images, 6000 for each of 10 classes. The dataset is split into 50000 training and 10000 test samples.

3.2.1 Model

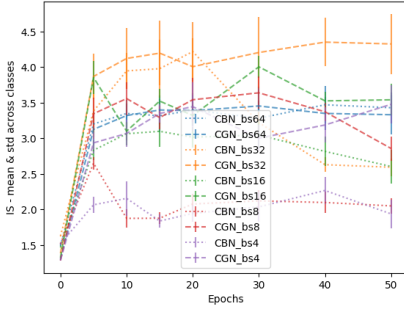
The SNGAN model architecture consists of a series of residual blocks followed by bn-relu-conv layers. Each residual block contains two bn-relu-conv modules, with an optional up-sampling layer. Since the architectures of more recent models such as SAGAN [43], BigGAN [8], BigBiGAN[13] are very similar to that of the one we used, it is likely that the conclusions we draw from the SNGAN experiments transfer to them.

We use the official implementation of SNGAN [29]. We replace the BN modules in the residual blocks with CBN and CGN for the respective cases, with number of groups set to 4 in case of CGN. We retain the optimisation setup of a learning rate of $2e-4$ for both generator and discriminator, five discriminator updates per generator update using the Adam optimiser [24]. We train using a single GPU (NVIDIA P100) and a batch size of 64. We also perform experiments where we use smaller batch sizes, to show the effects of CBN and CGN in helping generalisation.

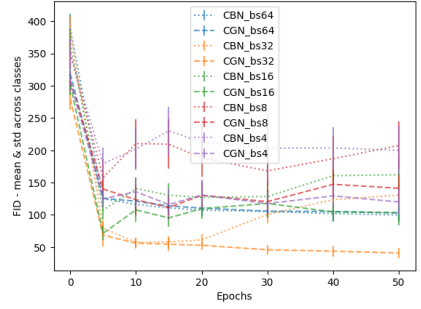
3.2.2 Results

Figure 2 shows samples from conditional SNGAN trained using each of the two normalisation methods.

For both normalisation methods, in addition to a qualitative check of the generated samples, we calculate two scores that are widely used in the community to evaluate image generation: Inception Score (IS) [7, 37] and Fréchet Inception Distance (FID) [19]. Since we are using publicly available PyTorch [64] implementations to compute these scores, the values for real data may differ slightly from scores computed using the original TensorFlow [1] implementation. However, we compare scores using the same implementation of these metrics for true and generated data.



(a) IS



(b) FID

Figure 3: (a) Inception Score (IS) [37] (higher is better), (b) Fréchet Inception Distance (FID) [19] (lower is better), each averaged across all classes for CBN- and CGN-based models for different batch sizes — 64, 32, 16, 8, 4

IS is meant to measure the natural-ness of an image by checking the embedding of the generated images on a pre-trained Inception network [38]. Although the suitability of IS for this purpose has been rightfully put into question [7], it continues to be used frequently. FID measures how similar two sets of images are, by computing the Fréchet distance between two multivariate Gaussians fitted to the embeddings of the images from the two sets. The embeddings are obtained from a pre-trained InceptionV3 network [38]. In this case, we measure the distance between the real and generated CIFAR-10 images. This is a better metric than IS, since there is no constraint on the images being natural, and it is able to quantify not only their similarity to the real images, but also diversity in the generated images.

We trained SNGAN models to generate CIFAR-10 images conditioned on the class for different batch sizes, viz. 64, 32, 16, 8, 4. In each case, we trained one model with CBN and another with CGN to compare them. We calculated IS and FID on these models at different stages of training, as can be seen in Figure 3 (a) and (b).

In the case of batch size 64, CGN performs similarly to CBN, and in all other cases using smaller batch sizes, CGN clearly outperforms CBN. This indicates the heavy dependence of CBN on batch size, which prevents it from generalizing well. In addition, CBN requires multi-GPU synchronization of batch statistics during optimisation, which can heavily hamper the training time, while CGN does not face this issue. Thus, we believe CGN is preferable as a module to use in a deep neural network-based generative models than CBN.

4 Conclusion

Because the performance of CBN heavily depends on the batch size and on how well training and test statistics match, we investigate the use of CGN as a potential alternative for CBN.

We experimentally show that the effect of this substitution is task-dependent, with performance increases in some VQA tasks that focus on systematic generalisation. In conditional image generation, we show that CGN can be trained with significantly smaller batch sizes than CBN, sometimes even with increased performance as measured by the IS and FID metrics. CGN’s simpler implementation, its consistent behaviour during training and inference time, as well as its independence from batch sizes, are all good reasons to explore its adop-

tion instead of CBN in tasks that require out-of-distribution generalisation. That being said, further analysis is required to be able to confidently suggest one method over the other. For instance, a hyperparameter search for each of the normalisation methods would be required to provide a more detailed performance comparison. As shown in our conditional image generation experiments, CGN is more amenable to training with smaller batch sizes. This suggests investigating applications in domains where efficient large-batch training is non-trivial, such as medical imaging or video processing.

Lastly, since some of the success of BN (and consequently CBN) can be attributed to the regularisation effect introduced by noisy batch statistics, it seems worthwhile to explore combinations of CGN with regularisation as suggested for GN by Wu and He [44]. This is also motivated by recent successful attempts at replacing (unconditional) BN with careful network initialisation [45], which relies on regularisation [44] to match generalisation performance.

References

- [1] Martín Abadi et al. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. URL <http://tensorflow.org/>.
- [2] Jacob Andreas, Marcus Rohrbach, Trevor Darrell, and Dan Klein. Neural module networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [3] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. Vqa: Visual question answering. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2015.
- [4] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- [5] Dzmitry Bahdanau, Felix Hill, Jan Leike, Edward Hughes, Pushmeet Kohli, and Edward Grefenstette. Learning to follow language instructions with adversarial reward induction. *arXiv preprint arXiv:1806.01946*, 2018.
- [6] Dzmitry Bahdanau, Shikhar Murty, Michael Noukhovitch, Thien Huu Nguyen, Harm de Vries, and Aaron Courville. Systematic generalization: What is required and can it be learned? *arXiv preprint arXiv:1811.12889*, 2018.
- [7] Shane Barratt and Rishi Kant Sharma. A note on the inception score. *CoRR*, abs/1801.01973, 2018.
- [8] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural image synthesis. *International Conference on Learning Representations*, 2019.
- [9] Jianpeng Cheng, Li Dong, and Mirella Lapata. Long short-term memory-networks for machine reading. *arXiv preprint arXiv:1601.06733*, 2016.
- [10] Kyunghyun Cho, Bart Van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. On the properties of neural machine translation: Encoder-decoder approaches. *arXiv preprint arXiv:1409.1259*, 2014.

- [11] Navneet Dalal and Bill Triggs. Histograms of oriented gradients for human detection. In *international Conference on computer vision & Pattern Recognition (CVPR'05)*, volume 1. IEEE Computer Society, 2005.
- [12] Harm De Vries, Florian Strub, Jérémie Mary, Hugo Larochelle, Olivier Pietquin, and Aaron C Courville. Modulating early visual processing by language. In *Advances in Neural Information Processing Systems*, 2017.
- [13] Jeff Donahue and Karen Simonyan. Large scale adversarial representation learning. In *NeurIPS*, 2019.
- [14] Vincent Dumoulin, Jonathon Shlens, and Manjunath Kudlur. A learned representation for artistic style. *International Conference on Learning Representations (ICLR)*, 2017.
- [15] Angus Galloway, Anna Golubeva, Thomas Tanay, Medhat Moussa, and Graham W Taylor. Batch normalization is a cause of adversarial vulnerability. *arXiv preprint arXiv:1905.02161*, 2019.
- [16] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C. Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems*, 2014.
- [17] Ishaan Gulrajani, Faruk Ahmed, Martín Arjovsky, Vincent Dumoulin, and Aaron C. Courville. Improved training of wasserstein gans. In *Advances in Neural Information Processing Systems*, 2017.
- [18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*, 2016.
- [19] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Advances in Neural Information Processing Systems*, 2017.
- [20] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning (ICML)*, 2015.
- [21] Xiang Jiang, Mohammad Havaei, Farshid Varno, Gabriel Chartrand, Nicolas Chapados, and Stan Matwin. Learning to learn with conditional class dependencies. 2018.
- [22] Justin Johnson, Bharath Hariharan, Laurens van der Maaten, Li Fei-Fei, C Lawrence Zitnick, and Ross Girshick. Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [23] Samira Ebrahimi Kahou, Vincent Michalski, Adam Atkinson, Ákos Kádár, Adam Trischler, and Yoshua Bengio. Figureqa: An annotated figure dataset for visual reasoning. *Workshop in the International Conference on Learning Representations*, 2017.
- [24] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *International Conference on Learning Representations (ICLR)*, 2014.

- [25] Alex Krizhevsky. Learning multiple layers of features from tiny images. 2009.
- [26] Mateusz Malinowski and Mario Fritz. A multi-world approach to question answering about real-world scenes based on uncertain input. In *Advances in neural information processing systems*, 2014.
- [27] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *CoRR*, abs/1411.1784, 2014.
- [28] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative adversarial networks. *ArXiv*, abs/1802.05957, 2018.
- [29] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative adversarial networks. *International Conference on Learning Representations (ICLR)*, 2018.
- [30] Vinod Nair and Geoffrey E Hinton. Rectified linear units improve restricted boltzmann machines. In *Proceedings of the 27th international conference on machine learning (ICML-10)*, 2010.
- [31] Augustus Odena, Christopher Olah, and Jonathon Shlens. Conditional image synthesis with auxiliary classifier gans. In *International Conference on Machine Learning (ICML)*, 2017.
- [32] Boris Oreshkin, Pau Rodríguez López, and Alexandre Lacoste. Tadam: Task dependent adaptive metric for improved few-shot learning. In *Advances in Neural Information Processing Systems*, 2018.
- [33] Ankur P Parikh, Oscar Täckström, Dipanjan Das, and Jakob Uszkoreit. A decomposable attention model for natural language inference. *arXiv preprint arXiv:1606.01933*, 2016.
- [34] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in PyTorch. In *NeurIPS Autodiff Workshop*, 2017.
- [35] Ethan Perez, Florian Strub, Harm De Vries, Vincent Dumoulin, and Aaron Courville. Film: Visual reasoning with a general conditioning layer. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [36] Olga Russakovsky et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3), 2015.
- [37] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, Xi Chen, and Xi Chen. Improved techniques for training gans. In *Advances in Neural Information Processing Systems*, 2016.
- [38] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016.
- [39] Dmitry Ulyanov, Andrea Vedaldi, and Victor S. Lempitsky. Instance normalization: The missing ingredient for fast stylization. *CoRR*, abs/1607.08022, 2016.

- [40] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, 2017.
- [41] Yuxin Wu and Kaiming He. Group normalization. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018.
- [42] Sitao Xiang and Hao Li. On the effects of batch and weight normalization in generative adversarial networks. *arXiv preprint arXiv:1704.03971*, 2017.
- [43] Han Zhang, Ian J. Goodfellow, Dimitris N. Metaxas, and Augustus Odena. Self-attention generative adversarial networks. *International Conference on Learning Representations (ICLR)*, 2018.
- [44] Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=r1Ddp1-Rb>.
- [45] Hongyi Zhang, Yann N. Dauphin, and Tengyu Ma. Fixup initialization: Residual learning without normalization. In *International Conference on Learning Representations*, 2019.