

MIT Open Access Articles

GAN Compression: Efficient Architectures for Interactive Conditional GANs

The MIT Faculty has made this article openly available. *Please share* how this access benefits you. Your story matters.

Citation: Li, Muyang et al. "GAN Compression: Efficient Architectures for Interactive Conditional GANs." Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2020 (June 2020) © 2020 The Author(s)

As Published: 10.1109/CVPR42600.2020.00533

Publisher: Institute of Electrical and Electronics Engineers (IEEE)

Persistent URL: <https://hdl.handle.net/1721.1/129446>

Version: Author's final manuscript: final author's manuscript post peer review, without publisher's formatting or copy editing

Terms of use: Creative Commons Attribution-Noncommercial-Share Alike



GAN Compression: Efficient Architectures for Interactive Conditional GANs

Muyang Li^{1,3} Ji Lin¹ Yaoyao Ding^{1,3} Zhijian Liu¹ Jun-Yan Zhu² Song Han¹
 lmxyy@mit.edu jilin@mit.edu yyding@mit.edu zhijian@mit.edu junzhu@adobe.com songhan@mit.edu
¹Massachusetts Institute of Technology ²Adobe Research ³Shanghai Jiao Tong University

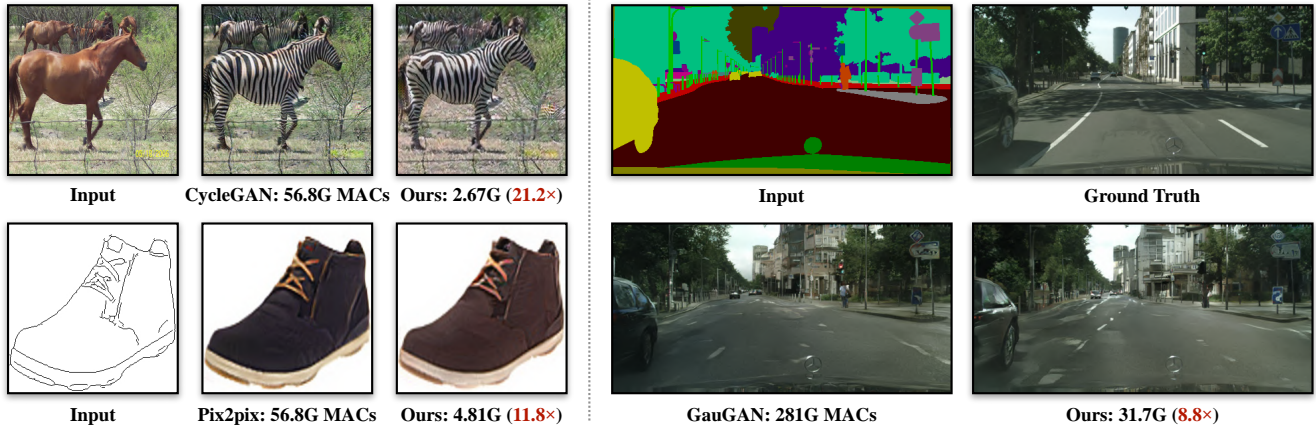


Figure 1: We introduce *GAN Compression*, a general-purpose method for compressing conditional GANs. Our method reduces the computation of widely-used conditional GAN models including pix2pix, CycleGAN, and GauGAN by 9-21 \times while preserving the visual fidelity. Our method is effective for a wide range of generator architectures, learning objectives, and both paired and unpaired settings.

Abstract

Conditional Generative Adversarial Networks (cGANs) have enabled controllable image synthesis for many computer vision and graphics applications. However, recent cGANs are 1-2 orders of magnitude more computationally-intensive than modern recognition CNNs. For example, GauGAN consumes 281G MACs per image, compared to 0.44G MACs for MobileNet-v3, making it difficult for interactive deployment. In this work, we propose a general-purpose compression framework for reducing the inference time and model size of the generator in cGANs. Directly applying existing CNNs compression methods yields poor performance due to the difficulty of GAN training and the differences in generator architectures. We address these challenges in two ways. First, to stabilize the GAN training, we transfer knowledge of multiple intermediate representations of the original model to its compressed model, and unify unpaired and paired learning. Second, instead of reusing existing CNN designs, our method automatically finds efficient architectures via neural architecture search (NAS). To accelerate the search process, we decouple the model training and architecture search via weight sharing. Experiments demon-

strate the effectiveness of our method across different supervision settings (paired and unpaired), model architectures, and learning methods (e.g., pix2pix, GauGAN, CycleGAN). Without losing image quality, we reduce the computation of CycleGAN by more than 20 \times and GauGAN by 9 \times , paving the way for interactive image synthesis. The *code* and *demo* are publicly available.

1. Introduction

Generative Adversarial Networks (GANs) [16] excel at synthesizing photo-realistic images. Their conditional extension, conditional GANs [50, 29, 80], allows controllable image synthesis and enables many computer vision and graphics applications such as interactively creating an image from a user drawing [52], transferring the motion of a dancing video stream to a different person [66, 9, 1], or creating VR facial animation for remote social interaction [68]. All of these applications require models to interact with humans and therefore demand low-latency on-device performance for better user experience. However, edge devices (mobile phones, tablets, VR headsets) are tightly constrained by hardware resources such as memory and battery. This com-

putational bottleneck prevents conditional GANs from being deployed on edge devices.

Different from image recognition CNNs [34, 61, 21, 27], image-conditional GANs are notoriously computationally intensive. For example, the widely-used CycleGAN model [80] requires more than 50G MACs*, 100× more than MobileNet [27]. A more recent model GauGAN [52], though generating photo-realistic high-resolution images, requires more than 250G MACs, 500× more than MobileNet [27, 57, 26].

In this work, we present *GAN Compression*, a general-purpose compression method for reducing the inference time and computational cost for conditional GANs. We observe that compressing generative models faces two fundamental difficulties: GANs are unstable to train, especially under the unpaired setting; generators also differ from recognition CNNs, making it hard to reuse existing CNN designs. To address these issues, we first transfer the knowledge from the intermediate representations of the original teacher generator to its corresponding layers of its compressed student generator. We also find it beneficial to create pseudo pairs using the teacher model’s output for unpaired training. This transforms the unpaired learning to a paired learning. Second, we use neural architecture search (NAS) to automatically find an efficient network with significantly fewer computation costs and parameters. To reduce the training cost, we decouple the model training from architecture search by training a “once-for-all network” that contains all possible channel number configurations. The once-for-all network can generate many sub-networks by weight sharing and enable us to evaluate the performance of each sub-network without retraining. Our method can be applied to various conditional GAN models regardless of model architectures, learning algorithms, and supervision settings (paired or unpaired).

Through extensive experiments, we show that our method can reduce the computation of three widely-used conditional GAN models including pix2pix [29], CycleGAN [80], and GauGAN [52] by 9× to 21× regarding MACs, without loss of the visual fidelity of generated images (see Figure 1 for several examples). Finally, we deploy our compressed pix2pix model on a mobile device (Jetson Nano) and demonstrate an interactive edges2shoes application (demo).

2. Related Work

Conditional GANs. Generative Adversarial Networks (GANs) [16] excel at synthesizing photo-realistic results [31, 5]. Its conditional form, conditional GANs [50, 29] further enables controllable image synthesis, allowing a user to synthesize images given various conditional inputs such

*We use the number of Multiply-Accumulate Operations (MAC) to quantify the computation cost. Modern computer architectures use fused multiply-add (FMA) instructions for tensor operations. These instructions compute $a = a + b \times c$ as one operation. 1 MAC=2 FLOPs.

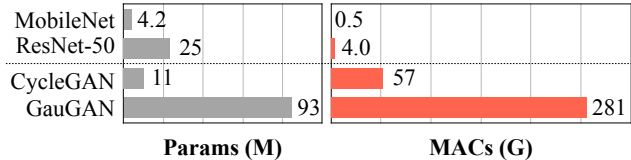


Figure 2: Conditional GANs require two orders of magnitude more computation than image classification CNNs, making it prohibitive to be deployed on edge devices.

as user sketches [29, 58], class labels [50, 5], or textual descriptions [55, 77]. Subsequent works further increased the resolution and realism of the results [67, 52]. Later, several algorithms were proposed to learn conditional GANs without paired data [63, 59, 80, 32, 71, 43, 13, 28, 35].

The high-resolution, photo-realistic synthesized results come at the cost of intensive computation. As shown in Figure 2, although the model size is of the same magnitude as the size of image recognition CNNs [21], conditional GANs require two orders of magnitudes more computations. This makes it challenging to deploy these models on edge devices given limited computational resources. In this work, we focus on efficient image-conditional GANs architectures for interactive applications.

Model acceleration. Extensive attention has been paid to hardware-efficient deep learning for various real-world applications [20, 19, 79, 65, 18]. To reduce redundancy in network weights, researchers proposed to prune the connections between layers [20, 19, 69]. However, the pruned networks require specialized hardware to achieve its full speedup. Several subsequent works proposed to prune entire convolution filters [23, 38, 44] to improve the regularity of computation. AutoML for Model Compression (AMC) [22] leverages reinforcement learning to determine the pruning ratio of each layer automatically. Liu *et al.* [45] later replaced the reinforcement learning by an evolutionary search algorithm. Recently, Shu *et al.* [60] proposed co-evolutionary pruning for CycleGAN by modifying the original CycleGAN algorithm. This method is tailored for a particular algorithm. The compressed model significantly increases FID under a moderate compression ratio (4.2×). In contrast, our model-agnostic method can be applied to conditional GANs with different learning algorithms, architectures, and both paired and unpaired settings. We assume no knowledge of the original cGAN learning algorithm. Experiments show that our general-purpose method achieves 21.1× compression ratio (5× better than CycleGAN-specific method [60]) while retaining the FID of original models.

Knowledge distillation. Hinton *et al.* [25] introduced the knowledge distillation for transferring the knowledge in a larger teacher network to a smaller student network. The student network is trained to mimic the behavior of the teacher

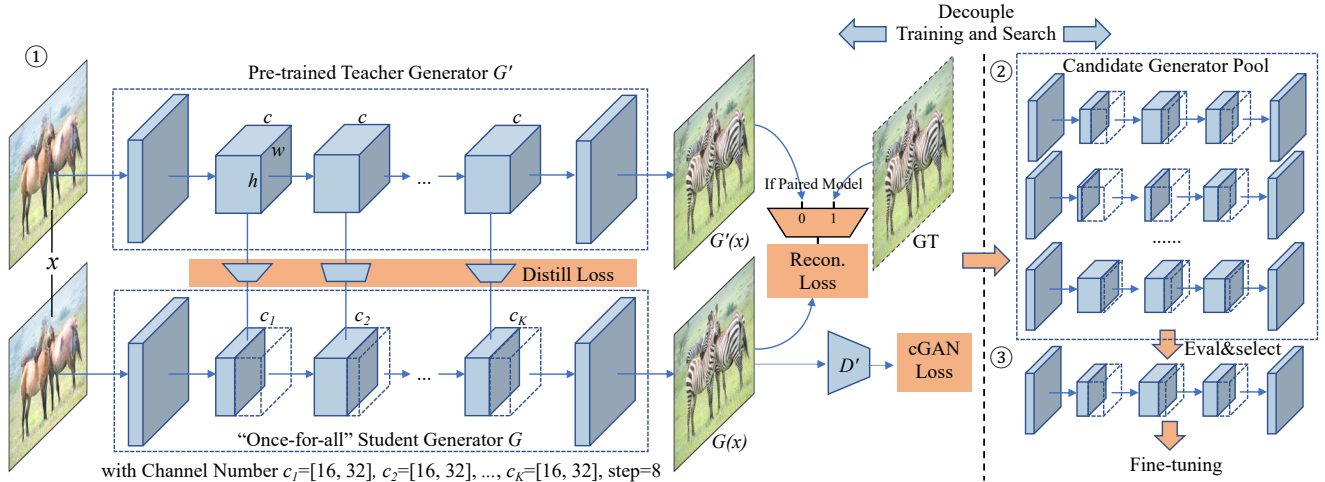


Figure 3: GAN Compression framework: ① Given a pre-trained teacher generator G' , we distill a smaller “once-for-all” student generator G that contains all possible channel numbers through weight sharing. We choose different channel numbers $\{c_k\}_{k=1}^K$ for the student generator G at each training step. ② We then extract many sub-generators from the “once-for-all” generator and evaluate their performance. No retraining is needed, which is the advantage of the “once-for-all” generator. ③ Finally, we choose the best sub-generator given the compression ratio target and performance target (FID or mIoU), perform fine-tuning, and obtain the final compressed model.

network. Several methods leverage knowledge distillation for compressing recognition models [48, 10, 36]. Recently, Aguinaldo *et al.* [2] adopts this method to accelerate unconditional GANs. Different from them, we focus on conditional GANs. We experimented with several distillation methods [2, 72] on conditional GANs and only observed marginal improvement, insufficient for interactive applications. Please refer to Appendix 6.2 for more details.

Neural architecture search. Neural Architecture Search (NAS) has successfully designed neural network architectures that outperform hand-crafted ones for large-scale image classification tasks [82, 40, 41]. To effectively reduce the search cost, researchers recently proposed one-shot neural architecture search [42, 8, 70, 17, 26, 4, 7] in which different candidate sub-networks can share the same set of weights. While all of these approaches focus on image classification models, we study efficient conditional GANs architectures using NAS.

3. Method

Compressing conditional generative models for interactive applications is challenging due to two reasons. Firstly, the training dynamic of GANs is highly unstable by nature. Secondly, the large architectural differences between recognition and generative models make it hard to apply existing CNN compression algorithms directly. To address the above issues, we propose a training protocol tailored for efficient generative models (Section 3.1) and further increase the compression ratio with neural architecture search

(NAS) (Section 3.2). The overall framework is illustrated in Figure 3. Here, we use the ResNet generator [30, 80] as an example. However, the same framework can be applied to different generator architectures and learning objectives.

3.1. Training Objective

Unifying unpaired and paired learning. Conditional GANs aim to learn a mapping function G between a source domain X and a target domain Y . They can be trained using either *paired* data ($\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^N$ where $\mathbf{x}_i \in X$ and $\mathbf{y}_i \in Y$) or *unpaired* data (source dataset $\{\mathbf{x}_i\}_{i=1}^N$ to target dataset $\{\mathbf{y}_j\}_{j=1}^M$). Here, N and M denote the number of training images. For simplicity, we omit the subscript i and j . Several learning objectives have been proposed to handle both paired and unpaired settings (e.g., [29, 52, 67, 80, 43, 28]). The wide range of training objectives makes it difficult to build a general-purpose compression framework. To address this, we unify the unpaired and paired learning in the model compression setting, regardless of how the teacher model is originally trained. Given the original teacher generator G' , we can transform the unpaired training setting to the paired setting. In particular, for the unpaired setting, we can view the original generator’s output as our ground-truth and train our compressed generator G with a paired learning objective. Our learning objective can be summarized as follows:

$$\mathcal{L}_{\text{recon}} = \begin{cases} \mathbb{E}_{\mathbf{x}, \mathbf{y}} \|G(\mathbf{x}) - \mathbf{y}\|_1 & \text{if paired cGANs,} \\ \mathbb{E}_{\mathbf{x}} \|G(\mathbf{x}) - G'(\mathbf{x})\|_1 & \text{if unpaired cGANs.} \end{cases} \quad (1)$$

Here we denote $\mathbb{E}_{\mathbf{x}} \triangleq \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})}$ and $\mathbb{E}_{\mathbf{x}, \mathbf{y}} \triangleq \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim p_{\text{data}}(\mathbf{x}, \mathbf{y})}$ for simplicity. $\|\cdot\|_1$ denotes L1 norm.

With such modifications, we can apply the same compression framework to different types of cGANs. Furthermore, As shown in Section 4.4, learning using the above pseudo pairs makes training more stable and yields much better results, compared to the original unpaired training setting.

As the unpaired training has been transformed into paired training, we will discuss the following sections in the paired training setting unless otherwise specified.

Inheriting the teacher discriminator. Although we aim to compress the generator, a discriminator D stores useful knowledge of a learned GAN as D learns to spot the weakness of the current generator [3]. Therefore, we adopt the same discriminator architecture, use the pre-trained weights from the teacher, and fine-tune the discriminator together with our compressed generator. In our experiments, we observe that a pre-trained discriminator could guide the training of our student generator. Using a randomly initialized discriminator often leads to severe training instability and the degradation of image quality. The GAN objective is formalized as:

$$\mathcal{L}_{\text{cGAN}} = \mathbb{E}_{\mathbf{x}, \mathbf{y}} [\log D(\mathbf{x}, \mathbf{y})] + \mathbb{E}_{\mathbf{x}} [\log(1 - D(\mathbf{x}, G(\mathbf{x})))] \quad (2)$$

where we initialize the student discriminator D using the weights from teacher discriminator D' . G and D are trained using a standard minimax optimization [16].

Intermediate feature distillation. A widely-used method for CNN model compression is knowledge distillation [25, 48, 10, 72, 36, 53, 12]. By matching the distribution of the output layer’s logits, we can transfer the dark knowledge from a teacher model to a student model, improving the performance of the student. However, conditional GANs [29, 80] usually output a deterministic image, rather than a probabilistic distribution. Therefore, it is difficult to distill the dark knowledge from the teacher’s output pixels. Especially for paired training setting, output images generated by the teacher model essentially contains no additional information compared to ground-truth target images. Experiments in Appendix 6.2 show that for paired training, naively mimicking the teacher model’s output brings no improvement.

To address the above issue, we match the intermediate representations of the teacher generator instead, as explored in prior work [36, 75, 10]. The intermediate layers contain more channels, provide richer information, and allow the student model to acquire more information in addition to outputs. The distillation objective can be formalized as

$$\mathcal{L}_{\text{distill}} = \sum_{t=1}^T \|f_t(G_t(\mathbf{x})) - G'_t(\mathbf{x})\|_2 \quad (3)$$

where $G_t(\mathbf{x})$ and $G'_t(\mathbf{x})$ are the intermediate feature activations of the t -th chosen layer in the student and teacher models, and T denotes the number of layers. A 1×1 learnable convolution layer f_t maps the features from the student model to the same number of channels in the features of the teacher model. We jointly optimize G_t and f_t to minimize the distillation loss $\mathcal{L}_{\text{distill}}$. Appendix 6.1 details which layers we choose in practice.

Full objective. Our final objective is written as follows:

$$\mathcal{L} = \mathcal{L}_{\text{cGAN}} + \lambda_{\text{recon}} \mathcal{L}_{\text{recon}} + \lambda_{\text{distill}} \mathcal{L}_{\text{distill}}, \quad (4)$$

where hyper-parameters λ_{recon} and λ_{distill} control the importance of each term. Please refer to Appendix 6.1 for more details.

3.2. Efficient Generator Design Space

Choosing a well-designed student architecture is essential for the final performance of knowledge distillation. We find that naively shrinking the channel numbers of the teacher model fails to produce a compact student model: the performance starts to degrade significantly above $4\times$ computation reduction. One of the possible reasons is that existing generator architectures are often adopted from image recognition models [46, 21, 56, 46], and may not be the optimal choice for image synthesis tasks. Below, we show how we derive a better architecture design space from an existing cGAN generator and perform neural architecture search (NAS) within the space.

Convolution decomposition and layer sensitivity. Existing generators usually adopt vanilla convolutions to follow the design of classification and segmentation CNNs. Recent efficient CNN designs widely adopt a decomposed version of convolutions (depthwise + pointwise) [27], which proves to have a better performance-computation trade-off. We find that using the decomposed convolution also benefits the generator design in cGANs.

Unfortunately, our early experiments have shown that directly applying decomposition to all the convolution layers (as in classifiers) will significantly degrade the image quality. Decomposing some of the layers will immediately hurt the performance, while other layers are more robust. Furthermore, this layer sensitivity pattern is not the same as recognition models. For example, in ResNet generator [21, 30], the resBlock layers consume the majority of the model parameters and computation cost while is almost immune to decomposition. On the contrary, the upsampling layers have much fewer parameters, but are fairly sensitive to model compression: moderate compression can lead to a large FID degradation. Therefore, we only decompose the resBlock layers. We conduct a comprehensive study regarding the sensitivity of layers in Section 4.4.

Automated channel reduction with NAS. Existing generators use a hand-crafted (and mostly uniform) channel numbers across all the layers, which contains redundancy and is far from optimal. To further improve the compression ratio, we automatically select the channel width in the generators using channel pruning [23, 22, 44, 81, 47] to remove the redundancy, which can reduce the computation quadratically. We support fine-grained choices regarding the numbers of channels. For each convolution layers, the number of channels can be chosen from multiples of 8, which balances MACs and hardware parallelism [22].

Given the possible channel configurations $\{c_1, c_2, \dots, c_K\}$, where K is the number of layers to prune, our goal is to find the best channel configuration $\{c_1^*, c_2^*, \dots, c_K^*\} = \arg \min_{c_1, c_2, \dots, c_K} \mathcal{L}, \quad s.t. \text{ MACs} < F_t$ using neural architecture search, where F_t is the computation constraint. A straight-forward approach is to traverse all the possible channel configuration, train it to convergence, evaluate, and pick the generator with the best performance. However, as K increases, the number of possible configurations increases exponentially, and each configuration might require different hyper-parameters regarding the learning rates and weights for each term. This trial and error process is far too time-consuming.

3.3. Decouple Training and Search

To address the problem, we decouple model training from architecture search, following recent work in one-shot neural architecture search methods [8, 7, 17]. We first train a “once-for-all” network [7] that supports different channel numbers. Each sub-network with different numbers of channels are equally trained and can operate independently. Sub-networks share the weights with the “once-for-all” network. Figure 3 illustrates the overall framework. We assume that the original teacher generator has $\{c_k^0\}_{k=1}^K$ channels. For a given channel number configuration $\{c_k\}_{k=1}^K, c_k \leq c_k^0$, we obtain the weight of the sub-network by extracting the first $\{c_k\}_{k=1}^K$ channels from the corresponding weight tensors of “once-for-all” network, following Guo et al. [17]. At each training step, we randomly sample a sub-network with a certain channel number configuration, compute the output and gradients, and update the extracted weights using our learning objective (Equation 4). Since the weights at the first several channels are updated more frequently, they play a more critical role among all the weights.

After the “once-for-all” network is trained, we find the best-performing sub-network by *directly* evaluating the performance of each candidate sub-network on the validation set through two types of search methods, as mentioned below:

- Brute Force: we directly evaluate the performance of each candidate sub-network under a certain computation budget. Though this method will definitely return

the best-performing one. However, this search process is rather time-consuming.

- Evolution: we include this search method since arXiv v3, which resorts to the evolution algorithm [54] to search for the best-performing sub-network. This method is much more efficient (about $20\times$ faster than the brute force) and could support a much larger search space.

Since the “once-for-all” network is thoroughly trained with weight sharing, no fine-tuning is needed. This approximates the model performance when it is trained from scratch. In this manner, we can decouple the training and search of the generator architecture: we only need to train once, but we can evaluate all the possible channel configurations without further training, and pick the best one as the search result. Optionally, we fine-tune the selected architecture to further improve the performance. We report both variants in Section 4.4.

4. Experiments

4.1. Setups

Models. We conduct experiments on three conditional GAN models to demonstrate the generality of our method.

- CycleGAN [80], an unpaired image-to-image translation model, uses a ResNet-based generator [21, 30] to transform an image from a source domain to a target domain, without using pairs.
- Pix2pix [29] is a conditional-GAN based paired image-to-image translation model. For this model, we replace the original U-Net generator [56] by the ResNet-based generator [30] as we observe that the ResNet-based generator achieves better results with less computation cost, given the same learning objective. See Appendix 6.2 for a detailed U-Net vs. ResNet benchmark.
- GauGAN [52] is a state-of-the-art paired image-to-image translation model. It can generate a high-fidelity image given a semantic label map.

We re-trained the pix2pix and CycleGAN using the official PyTorch repo with the above modifications. Our re-trained pix2pix and CycleGAN models (available at our [repo](#)) slightly outperform the official pre-trained models. We use these re-trained models as original models. For GauGAN, we use the pre-trained model from the authors. See Appendix 6.2 for more details.

Datasets. We use the following five datasets:

- Edges→shoes. We use 50,025 images from UT Zappos50K dataset [73]. We split the dataset randomly so

Model	Dataset	Method	#Parameters		MACs		Metric			
							FID (\downarrow)	mIoU (\uparrow)		
CycleGAN	horse \rightarrow zebra	Original	11.4M	–	56.8G	–	61.53	–	–	
		Shu <i>et al.</i> [60]	–	–	13.4G	(4.2 \times)	96.15	(34.6 \odot)	–	
		Ours (w/o fine-tuning)	0.34M	(33.3\times)	2.67G	(21.2 \times)	64.95	(3.42\odot)	–	
		Ours	0.34M	(33.3 \times)	2.67G	(21.2 \times)	71.81	(10.3 \odot)	–	
		Ours (Fast)	0.36M	(32.1 \times)	2.64G	(21.5\times)	65.19	(3.66 \odot)	–	
	edges \rightarrow shoes	Original	11.4M	–	56.8G	–	24.18	–	–	
		Ours (w/o fine-tuning)	0.70M	(16.3 \times)	4.81G	(11.8 \times)	31.30	(7.12 \odot)	–	
		Ours	0.70M	(16.3\times)	4.81G	(11.8\times)	26.60	(2.42 \odot)	–	
		Ours (Fast)	0.70M	(16.3 \times)	4.87G	(11.7 \times)	25.76	(1.58\odot)	–	
Pix2pix	Cityscapes	Original	11.4M	–	56.8G	–	–	42.06	–	
		Ours (w/o fine-tuning)	0.71M	(16.0 \times)	5.66G	(10.0 \times)	–	33.35	(8.71 \odot)	
		Ours	0.71M	(16.0\times)	5.66G	(10.0 \times)	–	40.77	(1.29 \odot)	
		Ours (Fast)	0.89M	(12.8 \times)	5.45G	(10.4\times)	–	41.71	(0.35\odot)	
	map \rightarrow aerial photo	Original	11.4M	–	56.8G	–	47.76	–	–	
		Ours (w/o fine-tuning)	0.75M	(15.1 \times)	4.68G	(12.1 \times)	71.82	(24.1 \odot)	–	
		Ours	0.75M	(15.1 \times)	4.68G	(12.1 \times)	48.02	(0.26\odot)	–	
		Ours (Fast)	0.71M	(16.1\times)	4.57G	(12.4\times)	48.67	(0.91 \odot)	–	
GauGAN	Cityscapes	Original	93.0M	–	281G	–	–	62.18	–	
		Ours (w/o fine-tuning)	20.4M	(4.6 \times)	31.7G	(8.8 \times)	–	59.44	(2.74 \odot)	
		Ours	20.4M	(4.6 \times)	31.7G	(8.8 \times)	–	61.22	(0.96\odot)	
		Ours (Fast)	20.2M	(4.6\times)	31.2G	(9.0\times)	–	61.17	(1.01 \odot)	
	COCO-Stuff	Original	97.5M	–	191G	–	21.38	–	38.78	–
Ours (Fast, w/o fine-tuning)		26.0M	(3.8 \times)	35.4G	(5.4 \times)	28.78	(7.40 \odot)	31.77	(7.01 \odot)	
		Ours (Fast)	26.0M	(3.8\times)	35.4G	(5.4\times)	25.06	(3.68\odot)	35.34	(3.45\odot)

Table 1: Quantitative evaluation of GAN Compression: we use the mIoU metric (the higher the better) for the Cityscapes and COCO-Stuff datasets, and FID (the lower the better) for other datasets. **Ours** denotes GAN Compression and **Ours (Fast)** denotes Fast GAN Compression as described in Section 4.1. Our method can compress state-of-the-art conditional GANs by **9-21 \times** in MACs and **5-33 \times** in model size, with only minor performance degradation. For CycleGAN compression, our general-purpose approach outperforms previous CycleGAN-specific Co-Evolution method [60] by a large margin.

that the validation set has 2,048 images for a stable evaluation of Fréchet Inception Distance (FID) (see Section 4.2). We evaluate the pix2pix model on this dataset.

- Cityscapes. The dataset [14] contains the images of German street scenes. The training set and the validation set consists of 2975 and 500 images, respectively. We evaluate both the pix2pix and GauGAN model on this dataset.
- Horse \leftrightarrow zebra. The dataset consists of 1,187 horse images and 1,474 zebra images originally from ImageNet [15] and is used in CycleGAN [80]. The validation set contains 120 horse images and 140 zebra images. We evaluate the CycleGAN model on this dataset.
- Map \leftrightarrow aerial photo. The dataset contains 2194 images scraped from Google Maps and used in pix2pix [29]. The training set and the validation set contains 1096 and 1098

images, respectively. We evaluate the pix2pix model on this dataset.

- COCO-Stuff. COCO-Stuff [6] dataset is derived from the COCO dataset [39]. It has 118,000 training images and 5,000 validation images. We evaluate the GauGAN model on this dataset.

Pipelines. We propose two pipelines for compressing conditional GAN generators.

- **GAN Compression.** For each model and dataset, we first train a MobileNet [27] style network from scratch, and then use the network as a teacher model to distill a student network. Initialized by the distilled student network, we train a “once-for-all” network. We then evaluate all

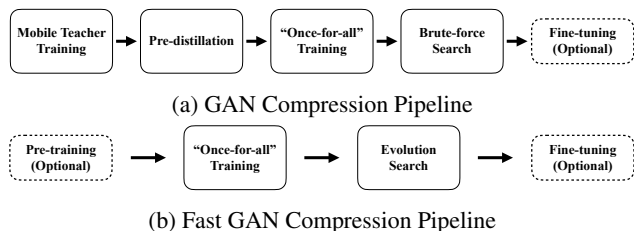


Figure 4: Detailed pipelines of the GAN Compression and Fast GAN Compression. Fast GAN Compression does not need “Mobile Teacher Training” and “Pre-distillation”, and switches the “Brute-force Search” to “Evolution Search”. The steps embraced by the dashed line are optional. If the original model has been pre-trained, the step “Pre-training” could be skipped. The “Fine-tuning” is also optional as reported in Section 3.3.

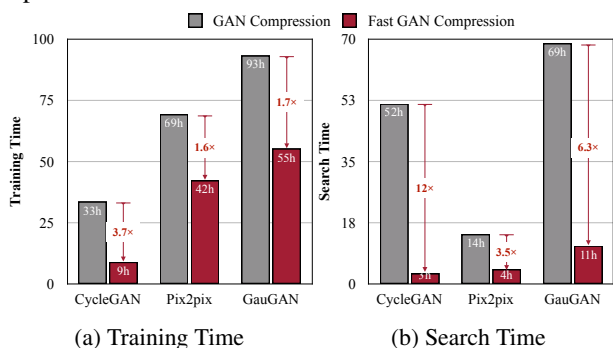


Figure 5: Training and search time comparison of GAN Compression and Fast GAN Compression. Fast GAN Compression could save 1.7 ~ 3.7× training time and 3.5 ~ 12× search time. CycleGAN, pix2pix and GauGAN models are measured on horse→zebra, edges→shoes and Cityscapes datasets. The training time of GauGAN is measured on 8 2080Ti GPUs, while others are all on a single 2080Ti.

sub-networks under a certain computation budget. After evaluation, we choose the best-performing sub-network within the “once-for-all” network and fine-tune it to obtain our final compressed model. The detailed pipeline is shown in Figure 4a. The sizes of the from-scratch MobileNet style teacher and the distilled student for each task are listed in Table 5. If not specified, our models are compressed with this pipeline.

- **Fast GAN Compression.** To further simplify and speed up the compression procedure, since arXiv v3, we propose an improved pipeline, Fast GAN Compression, which could produce comparable results as GAN Compression but with a much simpler and faster pipeline. In the training stage, we no longer need to train a MobileNet [27] style teacher network and run the pre-distillation. Instead, we directly train a MobileNet style “once-for-all” [7] network from scratch using the original full network as a teacher. In the search stage, instead of evaluating all sub-networks,

we adopt the evolution algorithm [54] to search for the best-performing sub-network under a certain computation budget within the “once-for-all” network and fine-tune it with the pre-trained discriminator to obtain our final compressed model. The detailed differences between the GAN Compression and Fast GAN Compression are shown in Figure 4b. With this pipeline, we could save up to 70% training time and 90% search time of GAN Compression as shown in Figure 5. Since the search space of Fast GAN Compression is larger than GAN Compression thanks to removing pre-distillation (the “once-for-all” network is larger) and a more efficient search method (see Appendix 6.1), its performance is on par with GAN Compression, as shown in Table 1. Please refer to our code for more details.

Implementation details. For the CycleGAN and pix2pix model, we use a learning rate of 0.0002 for both generator and discriminator during training in all the experiments. The batch sizes on dataset horse→zebra, edges→shoes, map→aerial photo, and cityscapes are 1, 4, 1, and 1, respectively. For the GauGAN model, we followed the setting in the original paper [52], except that the batch size is 16 instead of 32. We find that we can achieve a better result with a smaller batch size. See Appendix 6.1 and our code for more implementation details.

4.2. Evaluation Metrics

We introduce the metrics for assessing the equality of synthesized images.

Fréchet Inception Distance (FID) [24]. The FID score aims to calculate the distance between the distribution of feature vectors extracted from real and generated images using an InceptionV3 [62] network. The score measures the similarity between the distributions of real and generated images. A *lower* score indicates a *better* quality of generated images. We use an open-sourced FID evaluation code[†]. For paired image-to-image translation (pix2pix and GauGAN), we calculate the FID between translated test images to real test images. For unpaired image-to-image translations (CycleGAN), we calculate the FID between translated test images to real training+test images. This allows us to use more images for a stable FID evaluation, as done in previous unconditional GANs [31]. The FID of our compressed CycleGAN model slightly increases when we use real test images instead of real training+test images.

Semantic Segmentation Metrics. Following prior work [29, 80, 52], we adopt a semantic segmentation metric to evaluate the generated images on the Cityscapes dataset. We run a semantic segmentation model on the

[†]<https://github.com/mseitzer/pytorch-fid>

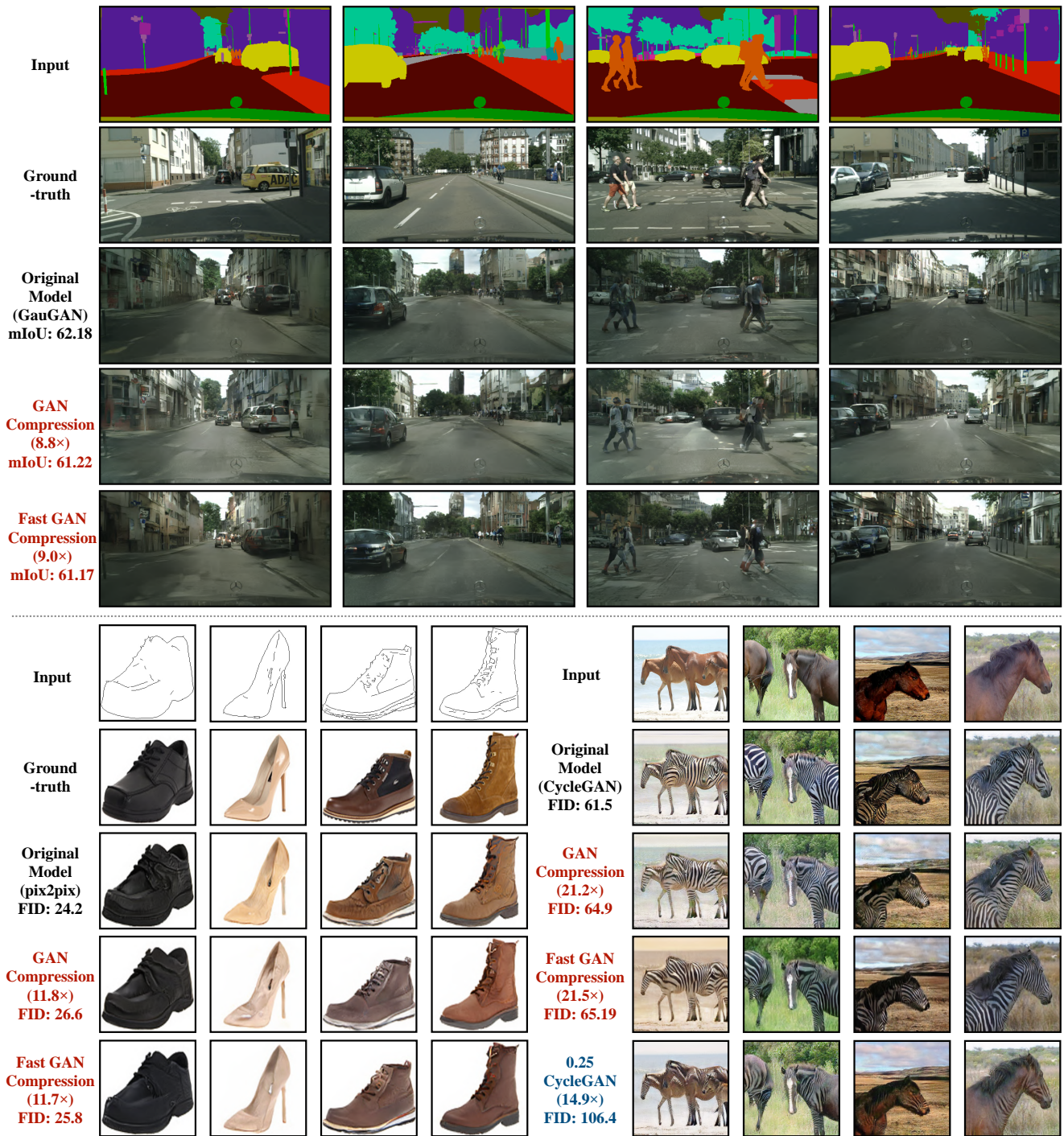


Figure 6: Qualitative compression results on Cityscapes, Edges→Shoes and Horse→Zebra. Our methods (GAN Compression and Fast GAN Compression) preserves the fidelity while significantly reducing the computation. In contrast, directly training a smaller model (e.g., 0.25 CycleGAN, which linearly scales each layer to 25% channels) yields poor performance.

generated images and compare how well the segmentation model performs. We choose the mean Intersection over Union (mIoU) as the segmentation metric, and we use DRN-D-105 [74] as our segmentation model for Cityscapes and DeepLabV2 [11] for COCO-Stuff. *Higher* mIoUs imply

that the generated images look more *realistic* and better reflect the input label map. For Cityscapes, we upsample the DRN-D-105’s output semantic map to 2048×1024 , which is the resolution of the ground truth images. For COCO-Stuff, we resize the generated images to the resolution of the

ground truth images. Please refer to our [code](#) for more evaluation details.

4.3. Results

Model	CycleGAN	Pix2pix	GauGAN	
Metric	FID (\downarrow)	61.5 \rightarrow 65.0	24.2 \rightarrow 26.6	–
	mIoU (\uparrow)	–	–	62.2 \rightarrow 61.2
MAC Reduction	21.2 \times	11.8 \times	8.8 \times	
Memory Reduction	2.0 \times	1.7 \times	1.8 \times	
Xavier	CPU	1.65s (18.5 \times)	3.07s (9.9 \times)	21.2s (7.9 \times)
Speedup	GPU	0.026s (3.1 \times)	0.035s (2.4 \times)	0.10s (3.2 \times)
Nano	CPU	6.30s (14.0 \times)	8.57s (10.3 \times)	65.3s (8.6 \times)
Speedup	GPU	0.16s (4.0 \times)	0.26s (2.5 \times)	0.81s (3.3 \times)
1080Ti Speedup		0.005s (2.5 \times)	0.007s (1.8 \times)	0.034s (1.7 \times)
Xeon Silver 4114	CPU Speedup	0.11s (3.4 \times)	0.15s (2.6 \times)	0.74s (2.8 \times)

Table 2: Measured memory reduction and latency speedup on NVIDIA Jetson AGX Xavier, NVIDIA Jetson Nano, 1080Ti GPU and Xeon CPU. CycleGAN, pix2pix, and GauGAN models are trained on horse \rightarrow zebra, edges \rightarrow shoes and Cityscapes datasets.

Quantitative Results We report the quantitative results of compressing CycleGAN, pix2pix, and GauGAN on five datasets in Table 1. By using the best performing sub-network from the “once-for-all” network, our method *GAN Compression* achieves large compression ratios. It can compress state-of-the-art conditional GANs by **5-21 \times** , and reduce the model size by **4-33 \times** , with only negligible degradation in the model performance. Specifically, our proposed method shows a clear advantage of CycleGAN compression compared to the previous Co-Evolution method [60]. We can reduce the computation of CycleGAN generator by 21.2 \times , which is 5 \times better compared to the previous CycleGAN-specific method [60] while achieving a better FID by more than 30 \ddagger .

Performance vs. Computation Trade-off Apart from the large compression ratio we can obtain, we verify that our method can consistently improve the performance at different model sizes. Taking the pix2pix model as an example, we plot the performance vs. computation trade-off on Cityscapes and Edges \rightarrow shoes dataset in Figure 8.

\ddagger In CycleGAN setting, for our model, the original model, and baselines, we report the FID between translated test images and real training+test images, while Shu et al. [60]’s FID is calculated between translated test images and real test images. The FID difference between the two protocols is small. The FIDs for the original model, Shu et al. [60], and our compressed model are 65.48, 96.15, and 69.54 using their protocol.

First, in the large model size regime, prune + distill (without NAS) outperforms training from scratch, showing the effectiveness of intermediate layer distillation. Unfortunately, with the channels continuing shrinking down uniformly, the capacity gap between the student and the teacher becomes too large. As a result, the knowledge from the teacher may be too recondite for the student, in which case the distillation may even have negative effects on the student model.

On the contrary, our training strategy allows us to automatically find a sub-network with a smaller gap between the student and teacher model, which makes learning easier. Our method consistently outperforms the baselines by a large margin.

Qualitative Results Figure 6 shows several example results. We provide the input, its ground-truth (except for unpaired setting), the output of the original model, and the output of our compressed models. Our compression methods well preserve the visual fidelity of the output image even under a large compression ratio. For CycleGAN, we also provide the output of a baseline model (0.25 CycleGAN: 14.9 \times). The baseline model 0.25 CycleGAN contains $\frac{1}{4}$ channels and has been trained from scratch. Our advantage is distinct: the baseline model can hardly create a zebra pattern on the output image, given a much smaller compression ratio. There might be some cases where compressed models show a small degradation (*e.g.*, the leg of the second zebra in Figure 6), but compressed models sometimes surpass the original one in other cases (*e.g.*, the first and last shoe images have a better leather texture). Generally, GAN models compressed by our methods perform comparatively compared to the original model, as shown by quantitative results.

Accelerate Inference on Hardware For real-world interactive applications, inference acceleration on hardware is more critical than the reduction of computation. To verify the practical effectiveness of our method, we measure the inference speed of our compressed models on several devices with different computational powers. To simulate interactive applications, we use a batch size of 1. We first perform 100 warm-up runs and measure the average timing of the next 100 runs. The results are shown in Table 2. The inference speed of compressed CycleGAN generator on edge GPU of Jetson Xavier can achieve about **40 FPS**, meeting the demand of interactive applications. We notice that the acceleration on GPU is less significant compared to CPU, mainly due to the large degree of parallelism. Nevertheless, we focus on making generative models more accessible on edge devices where powerful GPUs might not be available, so that more people can use interactive cGAN applications.

4.4. Ablation Study

Below we perform several ablation studies regarding our individual system components and design choices.

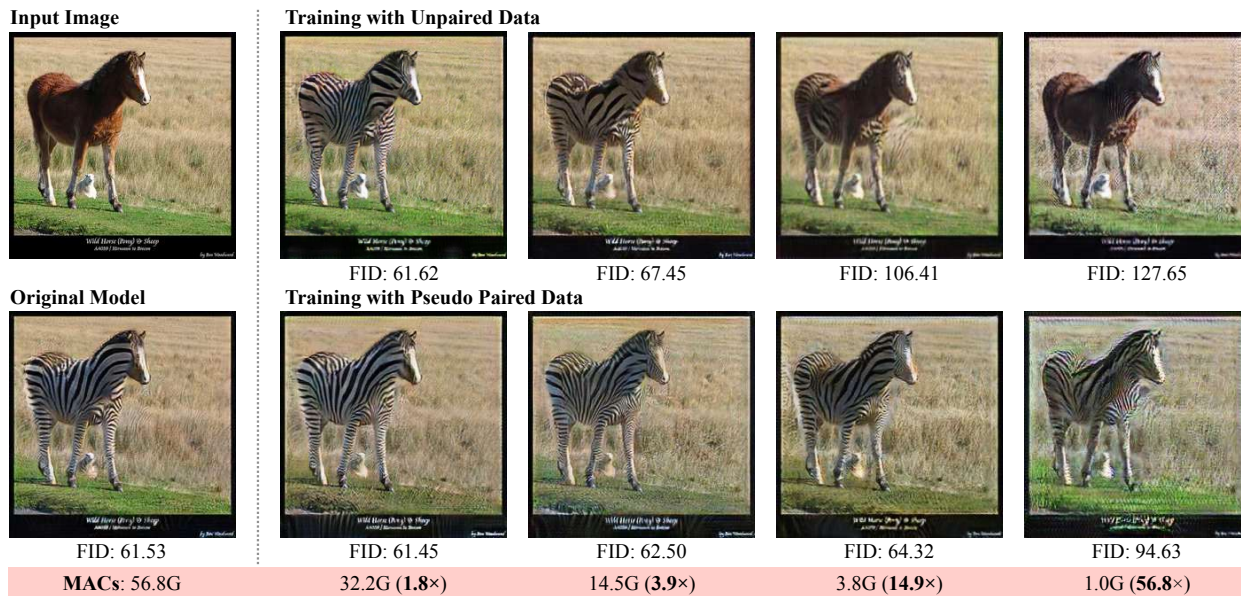


Figure 7: The comparison between training with unpaired data (naive) and training with pseudo paired data (proposed). The latter consistently outperforms the former, especially for small models. The generator’s computation can be compressed by 14.9× without hurting the fidelity using the proposed pseudo pair method. In this comparison, both methods do not use automated channel reduction and convolution decomposition.

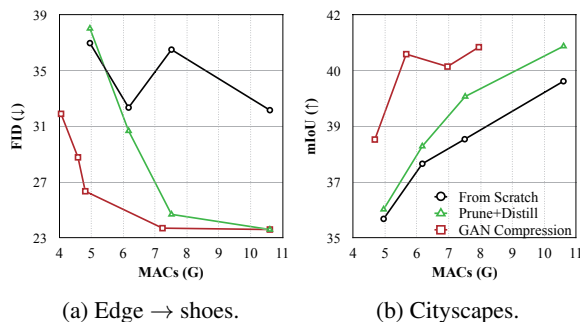


Figure 8: Trade off curve of pix2pix on Cityscapes and Edges→Shoes dataset. Pruning + distill method outperforms training from scratch for larger models, but works poorly when the model is aggressively shrunk. Our GAN Compression method can consistently improve the performance vs. computation trade-off at various scales.

Advantage of unpaired-to-paired transform. We first analyze the advantage of transforming unpaired conditional GANs into a pseudo paired training setting using the teacher model’s output.

Figure 9a shows the comparison of performance between the original unpaired training and our pseudo paired training. As our computation budget reduces, the quality of images generated by the unpaired training method degrades dramatically, while our pseudo paired training method remains relatively stable. The unpaired training requires the model to be strong enough to capture the complicated and ambiguous

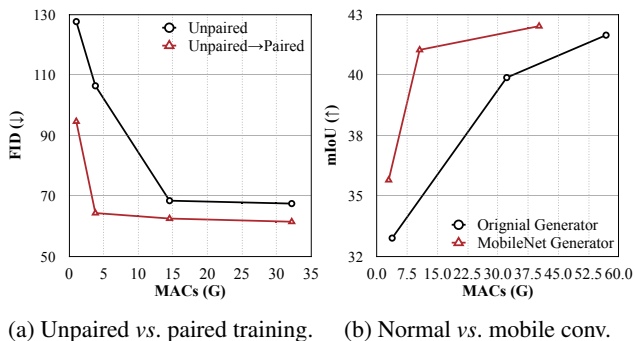


Figure 9: Ablation study: (a) Transforming the unpaired training into a paired training (using the pseudo pairs generated by the teacher model) significantly improves the performance of efficient models. (b) Decomposing the convolutions in the original ResNet-based generator into a channel-wise and depth-wise convolutional filters improves the performance vs. computation trade-off. We call our modified network MobileNet generator.

mapping between the source domain and the target domain. Once the mapping is learned, our student model can learn it from the teacher model directly. Additionally, the student model can still learn extra information on the real target images from the inherited discriminator.

The effectiveness of intermediate distillation and inheriting the teacher discriminator. Table 3 demonstrates the

Dataset	Setting	Training Technique			Metric	
		Pr.	Dstl.	Keep D.	FID (\downarrow)	mIoU (\uparrow)
Edges → shoes	ngf=64				24.91	–
					27.91	–
				✓	28.60	–
	ngf=48	✓		✓	27.25	–
			✓	✓	26.32	–
		✓	✓		46.24	–
	✓	✓	✓	24.45	–	
Cityscapes	ngf=96				–	42.47
					–	40.49
				✓	–	38.64
	ngf=64	✓		✓	–	40.98
			✓	✓	–	41.49
		✓	✓		–	40.66
		✓	✓	✓	–	42.11

Table 3: Ablation study. **Pr.**: Pruning; **Dstl.**: Distillation; **Keep D.**: in this setting, we inherit the discriminator’s weights from the teacher discriminator. Pruning combined with distillation achieves the best performance on both datasets.

Model	ngf	FID	MACs	#Parameters
Original	64	61.75	56.8G	11.38M
Only change downsample	64	68.72	55.5G	11.13M
Only change resBlocks	64	62.95	18.3G	1.98M
Only change upsample	64	61.04	48.3G	11.05M
Only change downsample	16	74.77	3.6G	0.70M
Only change resBlocks	16	79.49	1.4G	0.14M
Only change upsample	16	95.54	3.3G	0.70M

Table 4: We report the performance after applying convolution decomposition in each of the three parts (Downsample, ResBlocks, and Upsample) of the ResNet generator respectively on the horse→zebra dataset. ngf denotes the number of the generator’s filters. Both computation and model size are proportional to ngf^2 . We evaluate two settings $ngf=64$ and $ngf=16$. We observe that modifying ResBlock blocks shows a significantly better performance *vs.* computation trade-off, compared to modifying other parts of the network.

effectiveness of intermediate distillation and inheriting the teacher discriminator on the pix2pix model. Solely pruning and distilling intermediate feature cannot render a significantly better result than the baseline from-scratch training. We also explore the role of the discriminator in the pruned

ing. As a pre-trained discriminator stores useful information of the original generator, it can guide the pruned generator to learn faster and better. If the student discriminator is reset, the knowledge of the pruned student generator will be spoiled by the randomly initialized discriminator, which sometimes yields even worse results than the from-scratch training baseline.

Effectiveness of convolution decomposition. We systematically analyze the sensitivity of conditional GANs regarding the convolution decomposition transform. We take the ResNet-based generator from CycleGAN to test its effectiveness. We divide the structure of ResNet generator into three parts according to its network structure: Downsample (3 convolutions), ResBlocks (9 residual blocks), and Upsample (the final two deconvolutions). To validate the sensitivity of each stage, we replace all the conventional convolutions in each stage into separable convolutions [27]. The performance drop is reported in Table. 4. The ResBlock part takes a fair amount of computation cost, so decomposing the convolutions in the ResBlock can notably reduce computation costs. By testing both the architectures with $ngf=64$ and $ngf=16$, the ResBlock-modified architecture shows better computation costs *vs.* performance trade-off. We further explore the computation costs *vs.* performance trade-off of the ResBlock-modified architecture on Cityscapes dataset. Figure. 9b illustrates that such Mobilenet-style architecture is consistently more efficient than the original one, which has already reduced about half of the computation cost.

5. Conclusion

In this work, we proposed a general-purpose compression framework for reducing the computational cost and model size of generators in conditional GANs. We have used knowledge distillation and neural architecture search to alleviate training instability and to increase the model efficiency. Extensive experiments have shown that our method can compress several conditional GAN models while preserving the visual quality. Future work includes reducing the latency of models and efficient architectures for generative video models [66, 64].

Acknowledgments

We thank NSF Career Award #1943349, MIT-IBM Watson AI Lab, Adobe, Intel, Samsung and AWS machine learning research award for supporting this research. We thank Ning Xu, Zhuang Liu, Richard Zhang, and Antonio Torralba for helpful comments. We thank NVIDIA for donating the Jetson AGX Xavier that runs our [demo](#).

References

- [1] Kfir Aberman, Rundi Wu, Dani Lischinski, Baoquan Chen, and Daniel Cohen-Or. Learning character-agnostic motion for motion retargeting in 2d. In *SIGGRAPH*, 2019. [1](#)
- [2] Angeline Aguineldo, Ping-Yeh Chiang, Alex Gain, Ameya Patil, Kolten Pearson, and Soheil Feizi. Compressing gans using knowledge distillation. *arXiv preprint arXiv:1902.00159*, 2019. [3](#), [15](#), [16](#)
- [3] Samaneh Azadi, Catherine Olsson, Trevor Darrell, Ian Goodfellow, and Augustus Odena. Discriminator rejection sampling. In *ICLR*, 2019. [4](#)
- [4] Gabriel Bender, Pieter-Jan Kindermans, Barret Zoph, Vijay Vasudevan, and Quoc Le. Understanding and simplifying one-shot architecture search. In *ICML*, 2019. [3](#)
- [5] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural image synthesis. In *ICLR*, 2019. [2](#)
- [6] Holger Caesar, Jasper Uijlings, and Vittorio Ferrari. Coco-stuff: Thing and stuff classes in context. In *CVPR*, 2018. [6](#)
- [7] Han Cai, Chuang Gan, and Song Han. Once for all: Train one network and specialize it for efficient deployment. *ICLR*, 2020. [3](#), [5](#), [7](#), [15](#)
- [8] Han Cai, Ligeng Zhu, and Song Han. Proxylessnas: Direct neural architecture search on target task and hardware. In *ICLR*, 2019. [3](#), [5](#)
- [9] Caroline Chan, Shiry Ginosar, Tinghui Zhou, and Alexei A Efros. Everybody dance now. In *ICCV*, 2019. [1](#)
- [10] Guobin Chen, Wongun Choi, Xiang Yu, Tony Han, and Manmohan Chandraker. Learning efficient object detection models with knowledge distillation. In *NeurIPS*, 2017. [3](#), [4](#)
- [11] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence*, 40(4):834–848, 2017. [8](#)
- [12] Yuntao Chen, Naiyan Wang, and Zhaoxiang Zhang. Darkcrank: Accelerating deep metric learning via cross sample similarities transfer. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018. [4](#)
- [13] Yunjey Choi, Minje Choi, Munyoung Kim, Jung-Woo Ha, Sunghun Kim, and Jaegul Choo. Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. In *CVPR*, 2018. [2](#)
- [14] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *CVPR*, 2016. [6](#), [15](#)
- [15] 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. [6](#)
- [16] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *NeurIPS*, 2014. [1](#), [2](#), [4](#), [15](#)
- [17] Zichao Guo, Xiangyu Zhang, Haoyuan Mu, Wen Heng, Zechun Liu, Yichen Wei, and Jian Sun. Single path one-shot neural architecture search with uniform sampling. *arXiv preprint arXiv:1904.00420*, 2019. [3](#), [5](#)
- [18] Song Han, Han Cai, Ligeng Zhu, Ji Lin, Kuan Wang, Zhijian Liu, and Yujun Lin. Design automation for efficient deep learning computing. *arXiv preprint arXiv:1904.10616*, 2019. [2](#)
- [19] Song Han, Huizi Mao, and William J Dally. Deep compression: Compressing deep neural networks with pruning, trained quantization and Huffman coding. In *ICLR*, 2015. [2](#)
- [20] Song Han, Jeff Pool, John Tran, and William Dally. Learning both weights and connections for efficient neural network. In *NeurIPS*, 2015. [2](#)
- [21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. [2](#), [4](#), [5](#)
- [22] Yihui He, Ji Lin, Zhijian Liu, Hanrui Wang, Li-Jia Li, and Song Han. Amc: Automl for model compression and acceleration on mobile devices. In *ECCV*, 2018. [2](#), [5](#)
- [23] Yihui He, Xiangyu Zhang, and Jian Sun. Channel pruning for accelerating very deep neural networks. In *ICCV*, 2017. [2](#), [5](#)
- [24] 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 *NeurIPS*, 2017. [7](#)
- [25] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. In *NeurIPS Workshop*, 2015. [2](#), [4](#)
- [26] Andrew Howard, Mark Sandler, Grace Chu, Liang-Chieh Chen, Bo Chen, Mingxing Tan, Weijun Wang, Yukun Zhu, Ruoming Pang, Vijay Vasudevan, et al. Searching for mobilenetv3. *arXiv preprint arXiv:1905.02244*, 2019. [2](#), [3](#)
- [27] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017. [2](#), [4](#), [6](#), [7](#), [11](#)
- [28] Xun Huang, Ming-Yu Liu, Serge Belongie, and Jan Kautz. Multimodal unsupervised image-to-image translation. *ECCV*, 2018. [2](#), [3](#)
- [29] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *CVPR*, 2017. [1](#), [2](#), [3](#), [4](#), [5](#), [6](#), [7](#), [15](#), [16](#)
- [30] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *ECCV*, 2016. [3](#), [4](#), [5](#), [15](#)
- [31] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *CVPR*, 2019. [2](#), [7](#)
- [32] Taeksoo Kim, Moonsu Cha, Hyunsoo Kim, Jungkwon Lee, and Jiwon Kim. Learning to discover cross-domain relations with generative adversarial networks. In *ICML*, 2017. [2](#)
- [33] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015. [15](#)
- [34] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *NeurIPS*, 2012. [2](#)

- [35] Hsin-Ying Lee, Hung-Yu Tseng, Jia-Bin Huang, Maneesh Singh, and Ming-Hsuan Yang. Diverse image-to-image translation via disentangled representations. In *ECCV*, 2018. 2
- [36] Tianhong Li, Jianguo Li, Zhuang Liu, and Changshui Zhang. Knowledge distillation from few samples. *arXiv preprint arXiv:1812.01839*, 2018. 3, 4
- [37] Jae Hyun Lim and Jong Chul Ye. Geometric gan. *arXiv preprint arXiv:1705.02894*, 2017. 15
- [38] Ji Lin, Yongming Rao, Jiwen Lu, and Jie Zhou. Runtime neural pruning. In *NeurIPS*, 2017. 2
- [39] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014. 6
- [40] Chenxi Liu, Barret Zoph, Maxim Neumann, Jonathon Shlens, Wei Hua, Li-Jia Li, Li Fei-Fei, Alan Yuille, Jonathan Huang, and Kevin Murphy. Progressive neural architecture search. In *ECCV*, 2018. 3
- [41] Hanxiao Liu, Karen Simonyan, Oriol Vinyals, Chrisantha Fernando, and Koray Kavukcuoglu. Hierarchical representations for efficient architecture search. In *ICLR*, 2018. 3
- [42] Hanxiao Liu, Karen Simonyan, and Yiming Yang. Darts: Differentiable architecture search. In *ICLR*, 2019. 3
- [43] Ming-Yu Liu, Thomas Breuel, and Jan Kautz. Unsupervised image-to-image translation networks. In *NeurIPS*, 2017. 2, 3
- [44] Zhuang Liu, Jianguo Li, Zhiqiang Shen, Gao Huang, Shoumeng Yan, and Changshui Zhang. Learning efficient convolutional networks through network slimming. In *ICCV*, 2017. 2, 5
- [45] Zechun Liu, Haoyuan Mu, Xiangyu Zhang, Zichao Guo, Xin Yang, Tim Kwang-Ting Cheng, and Jian Sun. Metapruning: Meta learning for automatic neural network channel pruning. In *ICCV*, 2019. 2
- [46] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, pages 3431–3440, 2015. 4
- [47] Jian-Hao Luo, Jianxin Wu, and Weiyao Lin. Thinet: A filter level pruning method for deep neural network compression. In *Proceedings of the IEEE international conference on computer vision*, pages 5058–5066, 2017. 5
- [48] Ping Luo, Zhenyao Zhu, Ziwei Liu, Xiaogang Wang, and Xiaoou Tang. Face model compression by distilling knowledge from neurons. In *AAAI*, 2016. 3, 4
- [49] Xudong Mao, Qing Li, Haoran Xie, YK Raymond Lau, Zhen Wang, and Stephen Paul Smolley. Least squares generative adversarial networks. In *ICCV*, 2017. 15
- [50] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014. 1, 2
- [51] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative adversarial networks. In *ICLR*, 2018. 15
- [52] Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu. Semantic image synthesis with spatially-adaptive normalization. In *CVPR*, pages 2337–2346, 2019. 1, 2, 3, 5, 7, 15
- [53] Antonio Polino, Razvan Pascanu, and Dan Alistarh. Model compression via distillation and quantization. *arXiv preprint arXiv:1802.05668*, 2018. 4
- [54] Esteban Real, Alok Aggarwal, Yanping Huang, and Quoc V Le. Regularized evolution for image classifier architecture search. In *Proceedings of the aaai conference on artificial intelligence*, volume 33, pages 4780–4789, 2019. 5, 7
- [55] Scott Reed, Zeynep Akata, Xinchun Yan, Lajanugen Logeswaran, Bernt Schiele, and Honglak Lee. Generative adversarial text to image synthesis. In *ICML*, 2016. 2
- [56] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*, pages 234–241, 2015. 4, 5, 15
- [57] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4510–4520, 2018. 2
- [58] Patsorn Sangkloy, Jingwan Lu, Chen Fang, Fisher Yu, and James Hays. Scribbler: Controlling deep image synthesis with sketch and color. In *CVPR*, pages 5400–5409, 2017. 2
- [59] Ashish Shrivastava, Tomas Pfister, Oncel Tuzel, Josh Susskind, Wenda Wang, and Russ Webb. Learning from simulated and unsupervised images through adversarial training. In *CVPR*, 2017. 2
- [60] Han Shu, Yunhe Wang, Xu Jia, Kai Han, Hanting Chen, Chunjing Xu, Qi Tian, and Chang Xu. Co-evolutionary compression for unpaired image translation. In *ICCV*, 2019. 2, 6, 9
- [61] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015. 2
- [62] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *CVPR*, 2016. 7
- [63] Yaniv Taigman, Adam Polyak, and Lior Wolf. Unsupervised cross-domain image generation. In *ICLR*, 2017. 2
- [64] Carl Vondrick, Hamed Pirsiavash, and Antonio Torralba. Generating videos with scene dynamics. In *NeurIPS*, 2016. 11
- [65] Kuan Wang, Zhijian Liu, Yujun Lin, Ji Lin, and Song Han. Haq: Hardware-aware automated quantization with mixed precision. In *CVPR*, 2019. 2
- [66] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Guilin Liu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. Video-to-video synthesis. In *NeurIPS*, 2018. 1, 11
- [67] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. High-resolution image synthesis and semantic manipulation with conditional gans. In *CVPR*, 2018. 2, 3
- [68] Shih-En Wei, Jason Saragih, Tomas Simon, Adam W Harley, Stephen Lombardi, Michal Perdoch, Alexander Hypes, Dawei Wang, Hernan Badino, and Yaser Sheikh. Vr facial animation via multiview image translation. *ACM Transactions on Graphics (TOG)*, 38(4):67, 2019. 1
- [69] Wei Wen, Chunpeng Wu, Yandan Wang, Yiran Chen, and Hai Li. Learning structured sparsity in deep neural networks. In *NeurIPS*, 2016. 2
- [70] Bichen Wu, Xiaoliang Dai, Peizhao Zhang, Yanghan Wang, Fei Sun, Yiming Wu, Yuandong Tian, Peter Vajda, Yangqing

- Jia, and Kurt Keutzer. Fbnet: Hardware-aware efficient convnet design via differentiable neural architecture search. In *CVPR*, 2019. 3
- [71] Zili Yi, Hao Zhang, Ping Tan, and Minglun Gong. Dualgan: Unsupervised dual learning for image-to-image translation. In *ICCV*, 2017. 2
- [72] Junho Yim, Donggyu Joo, Jihoon Bae, and Junmo Kim. A gift from knowledge distillation: Fast optimization, network minimization and transfer learning. In *CVPR*, pages 4133–4141, 2017. 3, 4, 15, 16
- [73] Aron Yu and Kristen Grauman. Fine-grained visual comparisons with local learning. In *CVPR*, 2014. 5
- [74] Fisher Yu, Vladlen Koltun, and Thomas Funkhouser. Dilated residual networks. In *CVPR*, 2017. 8
- [75] Sergey Zagoruyko and Nikos Komodakis. Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. *arXiv preprint arXiv:1612.03928*, 2016. 4
- [76] Han Zhang, Ian Goodfellow, Dimitris Metaxas, and Augustus Odena. Self-attention generative adversarial networks. In *ICML*, 2019. 15
- [77] Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiaogang Wang, Xiao lei Huang, and Dimitris Metaxas. Stackgan++: Realistic image synthesis with stacked generative adversarial networks. *PAMI*, 2018. 2
- [78] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018. 16, 17
- [79] Chenzhuo Zhu, Song Han, Huizi Mao, and William J Dally. Trained ternary quantization. In *ICLR*, 2017. 2
- [80] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *ICCV*, 2017. 1, 2, 3, 4, 5, 6, 7, 15
- [81] Zhuangwei Zhuang, Mingkui Tan, Bohan Zhuang, Jing Liu, Yong Guo, Qingyao Wu, Junzhou Huang, and Jinhui Zhu. Discrimination-aware channel pruning for deep neural networks. In *Advances in Neural Information Processing Systems*, pages 875–886, 2018. 5
- [82] Barret Zoph and Quoc V Le. Neural architecture search with reinforcement learning. In *ICLR*, 2017. 3

6. Appendix

6.1. Additional Implementation Details

Training Epochs. In all experiments, we adopt the Adam optimizer [33] and keep the same learning rate in the beginning and linearly decay the rate to zero over in the later stage of the training. We use different epochs for the from-scratch training, distillation, and fine-tuning from the “once-for-all” [7] network training. The specific epochs for each task are listed in Table 5.

Distillation Layers. We choose 4 intermediate activations for distillation in our experiments. We split the 9 residual blocks into groups of size 3, and use feature distillation every three layers. We empirically find that such a configuration can transfer enough knowledge while is easy for student network to learn as shown in Appendix 6.2.

Loss function. For the pix2pix model [29], we replace the vanilla GAN loss [16] by a more stable Hinge GAN loss [37, 51, 76]. For the CycleGAN model [80] and GauGAN model [52], we follow the same setting of the original papers and use the LSGAN loss [49] and Hinge GAN loss term, respectively. We use the same GAN loss function for both teacher and student model as well as our baselines. The hyper-parameters λ_{recon} and λ_{distill} as mentioned in our paper are shown in Table 5.

Discriminator. A discriminator plays a critical role in the GAN training. We adopt the same discriminator architectures as the original work for each model. In our experiments, we did not compress the discriminator as it is not used at inference time. We also experimented with adjusting the capacity of discriminator but found it not helpful. We find that using the high-capacity discriminator with the compressed generator achieves better performance compared to using a compressed discriminator. Table 5 details the capacity of each discriminator.

Search Space. The search space design is critical for “once-for-all” [7] network training. Generally, a larger search space will produce more efficient models. To reduce the search cost, we remove certain channel number options in earlier layers, such as the input feature extractors. As we remove “pre-distillation” and use evolution search algorithm, Fast GAN Compression could support a much larger search space. The detailed search space sizes of GAN Compression and Fast GAN Compression are shown in Table 10. Please refer to our code for more details about the search space.

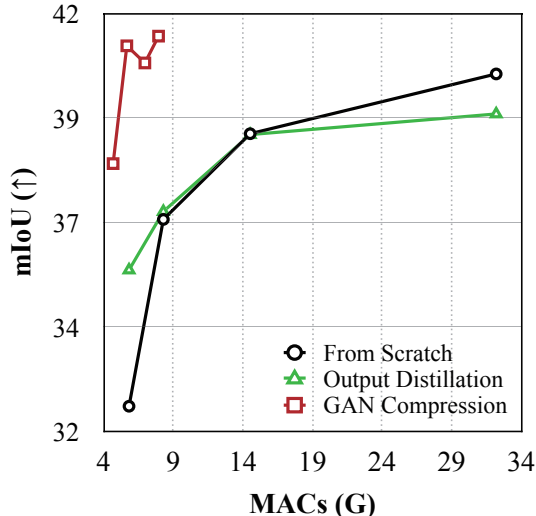


Figure 10: Performance vs. computation trade-off curve of pix2pix model on the Cityscapes dataset [14]. The output-only distillation method renders an even worse result than from-scratch training. Our GAN compression method significantly outperforms these two baselines.

6.2. Additional Ablation Study

Distillation. Recently, Aguinaldo *et al.* [2] adopts the knowledge distillation to accelerate the unconditional GANs inference. They enforce a student generator’s output to approximate a teacher generator’s output. However, in the paired conditional GAN training setting, the student generator can already learn enough information from its ground-truth target image. Therefore, the teacher’s output contains no extra information compared to the ground truth. Figure 10 empirically demonstrates this observation. We run the experiments for the pix2pix model on the cityscapes dataset [14]. The results from the distillation baseline [2] are even worse than models trained from scratch. Our GAN Compression method consistently outperforms these two baselines. We also compare our method with Yim *et al.* [72], a state-of-the-art distillation method used in recognition networks. Table 6 benchmarks different distillation methods on cityscapes dataset for pix2pix model. Our GAN Compression method outperforms other distillation methods by a large margin, paving the way for interactive image synthesis.

Network architecture for pix2pix. For pix2pix experiments, we replace the original U-net [56] by the ResNet-based generator [30]. Table 7 verifies our design choice. The ResNet generator achieves better performance on both edges→shoes and cityscapes datasets.

Model	Dataset	Training Epochs		Once-for-all Epochs		λ_{recon}	λ_{distill}	λ_{feat}	GAN Loss	ngf		ndf
		Const	Decay	Const	Decay					Teacher	Student	
Pix2pix	edges→shoes	5	15	10	30	100	1	-	Hinge	64	48	128
	cityscapes	100	150	200	300	100	1	-	Hinge	96	48	128
	map→arial photo	100	200	200	400	10	0.01	-	Hinge	96	48	128
CycleGAN	horse→zebra	100	100	200	200	10	0.01	-	LSGAN	64	32	64
GauGAN	cityscapes	100	100	100	100	10	10	10	Hinge	64	48	64

Table 5: Hyper-parameters setting of GAN Compression. Step 1, “Training Epochs” means the epochs for the from-scratch training, distillation and fine-tuning. Step 2, “Once-for-all Epochs” means epochs for the “once-for-all” network training. “Const” means the epochs of keeping the same initial learning rate, while “Decay” means epochs of linearly decaying the learning rate to 0. λ_{recon} and λ_{distill} are the weights of the reconstruction loss term (in GauGAN, this means VGG loss term) and the distillation loss term. λ_{feat} is the weight of the extra GAN feature loss term for GauGAN. “GAN Loss” is the specific type of GAN loss we use for each model. ngf, ndf denotes the base number of filters in a generator and discriminator, respectively, which is an indicator of the model size. Model computation and model size are proportional to ngf^2 (or ndf^2).

Method	MACs		mIoU \uparrow	
Original Model	56.8G	–	42.06	–
From-scratch Training	5.82G	(9.5 \times)	32.57	(9.49 \odot)
Aguinaldo <i>et al.</i> [2]	5.82G	(9.5 \times)	35.67	(6.39 \odot)
Yim <i>et al.</i> [72]	5.82G	(9.5 \times)	36.69	(5.37 \odot)
Intermediate Distillation	5.82G	(9.5 \times)	38.26	(3.80 \odot)
GAN Compression	5.66G	(10.0\times)	40.77	(1.29 \odot)

Table 6: Comparison of GAN Compression and different distillation methods (without NAS) for pix2pix model on cityscapes dataset. Our intermediate distillation outperforms other methods.

Dataset	Arch.	MACs	#Params	Metric	
				FID (\downarrow)	mIoU (\uparrow)
Edges→shoes	U-net	18.1G	54.4M	59.3	–
	ResNet	14.5G	2.85M	30.1	–
Cityscapes	U-net	18.1G	54.4M	–	28.4
	ResNet	14.5G	2.85M	–	33.6

Table 7: The comparison of U-net generator and ResNet generator for pix2pix model. The ResNet generator outperforms the U-net generator on both edges→Shoes dataset and cityscapes dataset.

Retrained models vs. pre-trained models. We retrain the original models with minor modifications as mentioned in Section 4.1. Table 8 shows our retrained results. For CycleGAN model, our retrained model slightly outperforms the pre-trained models. For the GauGAN model, our retrained model with official codebase is slightly worse than the pre-trained model. But our compressed model can

Model	Dataset	Setting	Metric	
			FID (\downarrow)	mIoU (\uparrow)
CycleGAN	horse→zebra	Pre-trained	71.84	–
		Retrained	61.53	–
		Compressed	64.95	–
GauGAN	Cityscapes	Pre-trained	–	62.18
		Retrained	–	61.04
		Compressed	–	61.22
	COCO-Stuff	Pre-trained	21.38	38.78
		Retrained	21.95	38.39
		Compressed	25.06	35.34

Table 8: The comparison of the official pre-trained models, our retrained models, and our compressed models. Our retrained CycleGAN model outperforms the official pre-trained models, so we report our retrained model results in Table 1. For the GauGAN model, our retrained model with official codebase is slightly worse than the pre-trained model, so we report the pre-trained model in Table 1. However, our compressed model achieves 61.22 mIoU on Cityscapes compared to the pre-trained model.

also achieve 61.22 mIoU, which has only negligible 0.96 mIoU drop compared to the pre-trained model on Cityscapes.

Perceptual similarity and user study. For a paired dataset, we evaluate the perceptual photorealism of our results. We use the LPIPS [78] metric to measure the perceptual similarity of generated images and the corresponding real images. Lower LPIPS indicates a better quality of the generated images. For the CycleGAN model, we conduct a human preference test on the horse→zebra dataset on Amazon Mechanical Turk (AMT) as there are no paired images. We basically follow the protocol of [29], except we ask the

Model	Dataset	Method	Preference	LPIPS (\downarrow)
Pix2pix	edges \rightarrow shoes	Original	–	0.185
		Ours	–	0.193(- 0.008)
		0.28 Pix2pix	–	0.201(-0.016)
	cityscapes	Original	–	0.0435
		Ours	–	0.436(- 0.001)
		0.31 Pix2pix	–	0.442(-0.007)
CycleGANhorse \rightarrow zebra	Ours	72.4%	–	
	0.25 CycleGAN	27.6%	–	

Table 9: Perceptual study: The LPIPS [78] is a perceptual metric for evaluating the similarity between a generated image and its corresponding ground truth real image. The lower LPIPS indicates a better perceptual photorealism of the results. This reference-based metric requires paired data. For unpaired setting, such as the horse \rightarrow zebra dataset, we conduct a human study for our GAN compression method and the 0.25 CycleGAN. We ask human participants which generated image looks more like a zebra. 72.4% workers favor results from our model.

Model	Dataset	Number of Sub-networks	
		GAN Comp.	Fast GAN Comp.
CycleGAN	horse \rightarrow zebra	6.6×10^3	3.9×10^5
	edges \rightarrow shoes	5.2×10^3	5.8×10^6
Pix2pix	cityscapes	5.2×10^3	1.58×10^6
	map \rightarrow arial photo	5.2×10^3	5.8×10^6
GauGAN	Cityscapes	1.6×10^4	1.3×10^5
	COCO-Stuff	–	3.9×10^5

Table 10: The detailed search space sizes of **GAN Compression (GAN Comp.)** and **Fast GAN Compression (Fast GAN Comp.)**. We use the number of sub-networks within a search space to measure the size of the search space. For the newest experiment of GauGAN on COCO-Stuff, we directly apply Fast GAN Compression, so we do not have the search space for the GAN Compression. Benefiting from removing “pre-distillation” and the evolution search algorithm, Fast GAN Compression could support a much larger search space than GAN Compression.

workers to decide which image is more like a real zebra image between our GAN Compression model and the 0.25 CycleGAN. Table 9 shows our perceptual study results on both the pix2pix model and the CycleGAN model. Our GAN Compression method significantly outperforms the straightforward from-scratch training baseline.

6.3. Additional Results

In Figure 11, we show additional visual results of our proposed GAN Compression and Fast GAN Compression methods for the CycleGAN model in horse \rightarrow zebra dataset.

In Figure 12, 13 and 14, we show additional visual results of our proposed methods for the pix2pix model on edges \rightarrow shoes, map \rightarrow arial photo and cityscapes datasets.

In Figure 15, we show additional visual results of our proposed methods for the GauGAN model on cityscapes.

6.4. Changelog

v1 Initial preprint release (CVPR 2020)

v2 (a) Correct the metric naming (mAP to mIoU). Update Cityscapes mIoU evaluation protocol ($\text{DRN}(\text{upsample}(G(x))) \rightarrow \text{upsample}(\text{DRN}(G(x)))$). See Section 4.2 and Table 1 and 3. (b) Add more details regarding Horse2zebra FID evaluation (Section 4.2). (c) Compare the official pre-trained models and our retrained models (Table 8).

v3 (a) Introduce Fast GAN Compression, a more efficient training method with a simplified training pipeline and a faster search strategy (Section 4.1). (b) Add the results of GauGAN on COCO-Stuff dataset (Table 1).

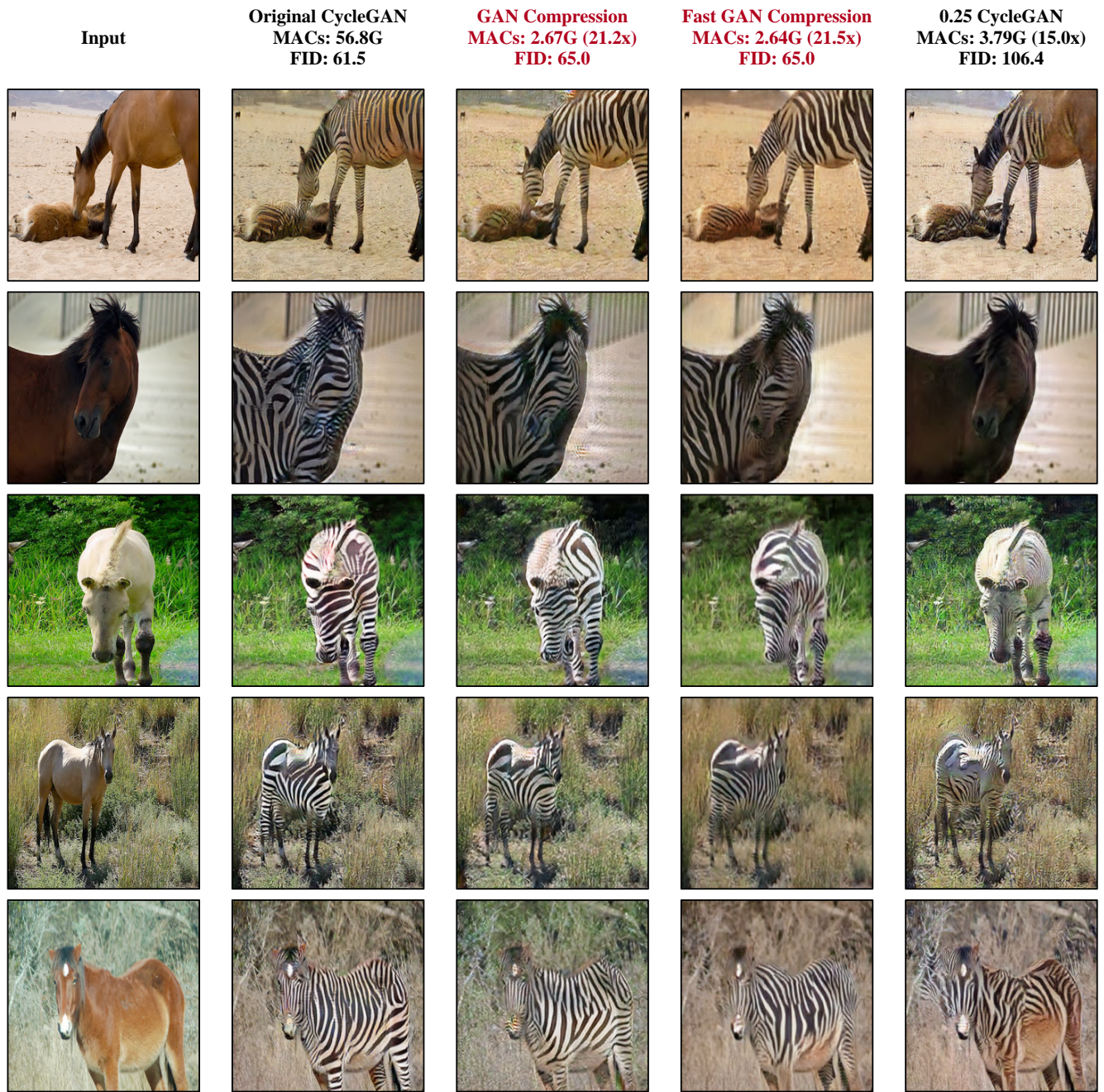


Figure 11: Additional results of GAN Compression and Fast GAN Compression with comparison to the 0.25 CycleGAN model on the horse→zebra dataset.

Input	Ground-truth	Original Pix2pix MACs: 56.8G FID: 24.2	GAN Compression MACs: 4.81G (11.8x) FID: 26.6	Fast GAN Compression MACs: 4.86G (11.7x) FID: 25.8	0.28 Pix2pix MACs: 4.75G (12.0x) FID: 37.1
					
					
					
					
					
					

Figure 12: Additional results of GAN Compression and Fast GAN Compression with comparison to the 0.28 pix2pix model on the edges→shoes dataset.

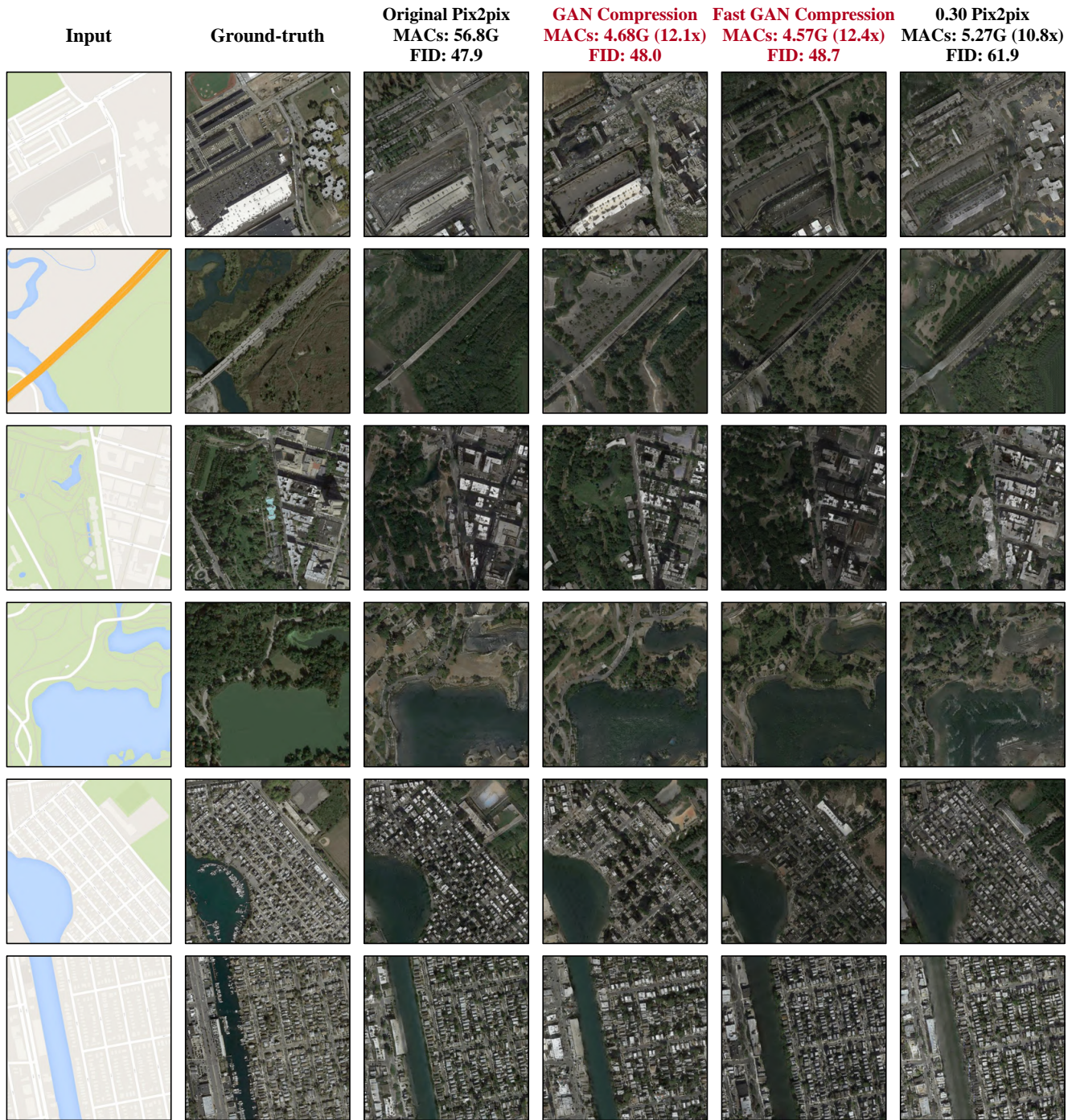


Figure 13: Additional results of GAN Compression and Fast GAN Compression with comparison to the 0.30 pix2pix model on the map→arial photo dataset.



Figure 14: Additional results of GAN Compression and Fast GAN Compression with comparison to the 0.31 pix2pix model on the cityscapes dataset.



Figure 15: Additional results of compressing GauGAN model on the cityscapes dataset.