

Analyzing and Improving the Image Quality of StyleGAN

Tero Karras
NVIDIA

Samuli Laine
NVIDIA

Miika Aittala
NVIDIA

Janne Hellsten
NVIDIA

Jaakko Lehtinen
NVIDIA and Aalto University

Timo Aila
NVIDIA

Abstract

The style-based GAN architecture (StyleGAN) yields state-of-the-art results in data-driven unconditional generative image modeling. We expose and analyze several of its characteristic artifacts, and propose changes in both model architecture and training methods to address them. In particular, we redesign the generator normalization, revisit progressive growing, and regularize the generator to encourage good conditioning in the mapping from latent codes to images. In addition to improving image quality, this path length regularizer yields the additional benefit that the generator becomes significantly easier to invert. This makes it possible to reliably attribute a generated image to a particular network. We furthermore visualize how well the generator utilizes its output resolution, and identify a capacity problem, motivating us to train larger models for additional quality improvements. Overall, our improved model redefines the state of the art in unconditional image modeling, both in terms of existing distribution quality metrics as well as perceived image quality.

1. Introduction

The resolution and quality of images produced by generative methods, especially generative adversarial networks (GAN) [13], are improving rapidly [20, 26, 4]. The current state-of-the-art method for high-resolution image synthesis is StyleGAN [21], which has been shown to work reliably on a variety of datasets. Our work focuses on fixing its characteristic artifacts and improving the result quality further.

The distinguishing feature of StyleGAN [21] is its unconventional generator architecture. Instead of feeding the input latent code $\mathbf{z} \in \mathcal{Z}$ only to the beginning of a the network, the *mapping network* f first transforms it to an intermediate latent code $\mathbf{w} \in \mathcal{W}$. Affine transforms then produce *styles* that control the layers of the *synthesis network* g via adaptive instance normalization (AdaIN) [18, 8, 11, 7]. Additionally, stochastic variation is facilitated by providing

additional random noise maps to the synthesis network. It has been demonstrated [21, 33] that this design allows the intermediate latent space \mathcal{W} to be much less entangled than the input latent space \mathcal{Z} . In this paper, we focus all analysis solely on \mathcal{W} , as it is the relevant latent space from the synthesis network’s point of view.

Many observers have noticed characteristic artifacts in images generated by StyleGAN [3]. We identify two causes for these artifacts, and describe changes in architecture and training methods that eliminate them. First, we investigate the origin of common blob-like artifacts, and find that the generator creates them to circumvent a design flaw in its architecture. In Section 2, we redesign the normalization used in the generator, which removes the artifacts. Second, we analyze artifacts related to progressive growing [20] that has been highly successful in stabilizing high-resolution GAN training. We propose an alternative design that achieves the same goal—training starts by focusing on low-resolution images and then progressively shifts focus to higher and higher resolutions—without changing the network topology during training. This new design also allows us to reason about the effective resolution of the generated images, which turns out to be lower than expected, motivating a capacity increase (Section 4).

Quantitative analysis of the quality of images produced using generative methods continues to be a challenging topic. Fréchet inception distance (FID) [17] measures differences in the density of two distributions in the high-dimensional feature space of an InceptionV3 classifier [34]. Precision and Recall (P&R) [31, 22] provide additional visibility by explicitly quantifying the percentage of generated images that are similar to training data and the percentage of training data that can be generated, respectively. We use these metrics to quantify the improvements.

Both FID and P&R are based on classifier networks that have recently been shown to focus on textures rather than shapes [10], and consequently, the metrics do not accurately capture all aspects of image quality. We observe that the perceptual path length (PPL) metric [21], originally introduced as a method for estimating the quality of latent space



Figure 1. Instance normalization causes water droplet -like artifacts in StyleGAN images. These are not always obvious in the generated images, but if we look at the activations inside the generator network, the problem is always there, in all feature maps starting from the 64x64 resolution. It is a systemic problem that plagues all StyleGAN images.

interpolations, correlates with consistency and stability of shapes. Based on this, we regularize the synthesis network to favor smooth mappings (Section 3) and achieve a clear improvement in quality. To counter its computational expense, we also propose executing all regularizations less frequently, observing that this can be done without compromising effectiveness.

Finally, we find that projection of images to the latent space \mathcal{W} works significantly better with the new, path-length regularized StyleGAN2 generator than with the original StyleGAN. This makes it easier to attribute a generated image to its source (Section 5).

Our implementation and trained models are available at <https://github.com/NVlabs/stylegan2>

2. Removing normalization artifacts

We begin by observing that most images generated by StyleGAN exhibit characteristic blob-shaped artifacts that resemble water droplets. As shown in Figure 1, even when the droplet may not be obvious in the final image, it is present in the intermediate feature maps of the generator.¹ The anomaly starts to appear around 64×64 resolution, is present in all feature maps, and becomes progressively stronger at higher resolutions. The existence of such a consistent artifact is puzzling, as the discriminator should be able to detect it.

We pinpoint the problem to the AdaIN operation that normalizes the mean and variance of each feature map separately, thereby potentially destroying any information found in the magnitudes of the features relative to each other. We hypothesize that the droplet artifact is a result of the generator intentionally sneaking signal strength information past instance normalization: by creating a strong, localized spike that dominates the statistics, the generator can effectively scale the signal as it likes elsewhere. Our hypothesis is supported by the finding that when the normalization step is removed from the generator, as detailed below, the droplet artifacts disappear completely.

¹In rare cases (perhaps 0.1% of images) the droplet is missing, leading to severely corrupted images. See Appendix A for details.

2.1. Generator architecture revisited

We will first revise several details of the StyleGAN generator to better facilitate our redesigned normalization. These changes have either a neutral or small positive effect on their own in terms of quality metrics.

Figure 2a shows the original StyleGAN synthesis network g [21], and in Figure 2b we expand the diagram to full detail by showing the weights and biases and breaking the AdaIN operation to its two constituent parts: normalization and modulation. This allows us to re-draw the conceptual gray boxes so that each box indicates the part of the network where one style is active (i.e., “style block”). Interestingly, the original StyleGAN applies bias and noise within the style block, causing their relative impact to be inversely proportional to the current style’s magnitudes. We observe that more predictable results are obtained by moving these operations outside the style block, where they operate on normalized data. Furthermore, we notice that after this change it is sufficient for the normalization and modulation to operate on the standard deviation alone (i.e., the mean is not needed). The application of bias, noise, and normalization to the constant input can also be safely removed without observable drawbacks. This variant is shown in Figure 2c, and serves as a starting point for our redesigned normalization.

2.2. Instance normalization revisited

One of the main strengths of StyleGAN is the ability to control the generated images via *style mixing*, i.e., by feeding a different latent w to different layers at inference time. In practice, style modulation may amplify certain feature maps by an order of magnitude or more. For style mixing to work, we must explicitly counteract this amplification on a per-sample basis — otherwise the subsequent layers would not be able to operate on the data in a meaningful way.

If we were willing to sacrifice scale-specific controls (see video), we could simply remove the normalization, thus removing the artifacts and also improving FID slightly [22]. We will now propose a better alternative that removes the artifacts while retaining full controllability. The main idea is to base normalization on the *expected* statistics of the incoming feature maps, but without explicit forcing.

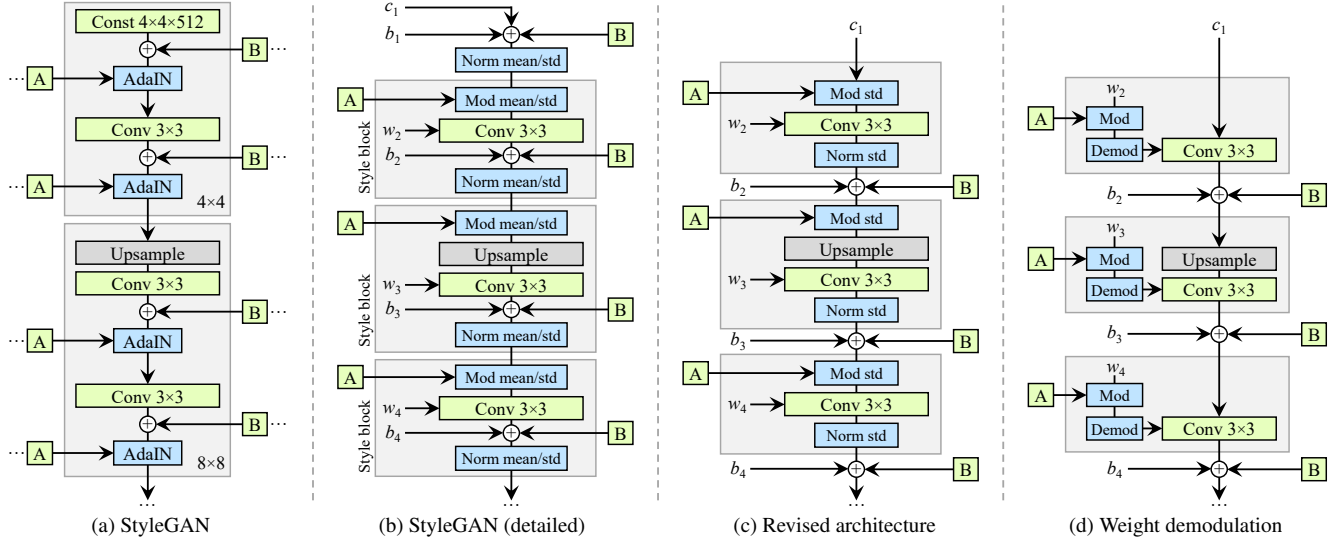


Figure 2. We redesign the architecture of the StyleGAN synthesis network. (a) The original StyleGAN, where \boxed{A} denotes a learned affine transform from \mathcal{W} that produces a style and \boxed{B} is a noise broadcast operation. (b) The same diagram with full detail. Here we have broken the AdaIN to explicit normalization followed by modulation, both operating on the mean and standard deviation per feature map. We have also annotated the learned weights (w), biases (b), and constant input (c), and redrawn the gray boxes so that one style is active per box. The activation function (leaky ReLU) is always applied right after adding the bias. (c) We make several changes to the original architecture that are justified in the main text. We remove some redundant operations at the beginning, move the addition of b and \boxed{B} to be outside active area of a style, and adjust only the standard deviation per feature map. (d) The revised architecture enables us to replace instance normalization with a “demodulation” operation, which we apply to the weights associated with each convolution layer.

Recall that a style block in Figure 2c consists of modulation, convolution, and normalization. Let us start by considering the effect of a modulation followed by a convolution. The modulation scales each input feature map of the convolution based on the incoming style, which can alternatively be implemented by scaling the convolution weights:

$$w'_{ijk} = s_i \cdot w_{ijk}, \quad (1)$$

where w and w' are the original and modulated weights, respectively, s_i is the scale corresponding to the i th input feature map, and j and k enumerate the output feature maps and spatial footprint of the convolution, respectively.

Now, the purpose of instance normalization is to essentially remove the effect of s from the statistics of the convolution’s output feature maps. We observe that this goal can be achieved more directly. Let us assume that the input activations are i.i.d. random variables with unit standard deviation. After modulation and convolution, the output activations have standard deviation of

$$\sigma_j = \sqrt{\sum_{i,k} w'_{ijk}{}^2}, \quad (2)$$

i.e., the outputs are scaled by the L_2 norm of the corresponding weights. The subsequent normalization aims to restore the outputs back to unit standard deviation. Based on Equation 2, this is achieved if we scale (“demodulate”)

each output feature map j by $1/\sigma_j$. Alternatively, we can again bake this into the convolution weights:

$$w''_{ijk} = w'_{ijk} / \sqrt{\sum_{i,k} w'_{ijk}{}^2 + \epsilon}, \quad (3)$$

where ϵ is a small constant to avoid numerical issues.

We have now baked the entire style block to a single convolution layer whose weights are adjusted based on s using Equations 1 and 3 (Figure 2d). Compared to instance normalization, our demodulation technique is weaker because it is based on statistical assumptions about the signal instead of actual contents of the feature maps. Similar statistical analysis has been extensively used in modern network initializers [12, 16], but we are not aware of it being previously used as a replacement for data-dependent normalization. Our demodulation is also related to weight normalization [32] that performs the same calculation as a part of reparameterizing the weight tensor. Prior work has identified weight normalization as beneficial in the context of GAN training [38].

Our new design removes the characteristic artifacts (Figure 3) while retaining full controllability, as demonstrated in the accompanying video. FID remains largely unaffected (Table 1, rows A, B), but there is a notable shift from precision to recall. We argue that this is generally desirable, since recall can be traded into precision via truncation, whereas

Configuration	FFHQ, 1024×1024				LSUN Car, 512×384			
	FID ↓	Path length ↓	Precision ↑	Recall ↑	FID ↓	Path length ↓	Precision ↑	Recall ↑
A Baseline StyleGAN [21]	4.40	212.1	0.721	0.399	3.27	1484.5	0.701	0.435
B + Weight demodulation	4.39	175.4	0.702	0.425	3.04	862.4	0.685	0.488
C + Lazy regularization	4.38	158.0	0.719	0.427	2.83	981.6	0.688	0.493
D + Path length regularization	4.34	122.5	0.715	0.418	3.43	651.2	0.697	0.452
E + No growing, new G & D arch.	3.31	124.5	0.705	0.449	3.19	471.2	0.690	0.454
F + Large networks (StyleGAN2)	2.84	145.0	0.689	0.492	2.32	415.5	0.678	0.514
Config A with large networks	3.98	199.2	0.716	0.422	—	—	—	—

Table 1. Main results. For each training run, we selected the training snapshot with the lowest FID. We computed each metric 10 times with different random seeds and report their average. *Path length* corresponds to the PPL metric, computed based on path endpoints in \mathcal{W} [21], without the central crop used by Karras et al. [21]. The FFHQ dataset contains 70k images, and the discriminator saw 25M images during training. For LSUN CAR the numbers were 893k and 57M. \uparrow indicates that higher is better, and \downarrow that lower is better.

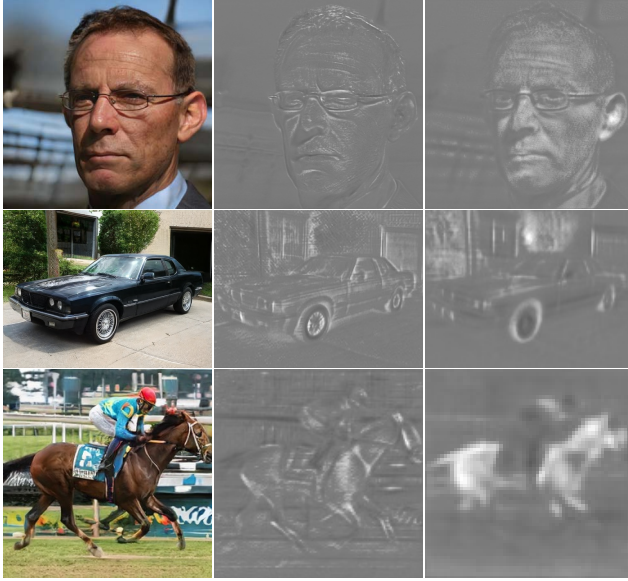


Figure 3. Replacing normalization with demodulation removes the characteristic artifacts from images and activations.

the opposite is not true [22]. In practice our design can be implemented efficiently using grouped convolutions, as detailed in Appendix B. To avoid having to account for the activation function in Equation 3, we scale our activation functions so that they retain the expected signal variance.

3. Image quality and generator smoothness

While GAN metrics such as FID or Precision and Recall (P&R) successfully capture many aspects of the generator, they continue to have somewhat of a blind spot for image quality. For an example, refer to Figures 3 and 4 in the Supplement that contrast generators with identical FID and P&R scores but markedly different overall quality.²

²We believe that the key to the apparent inconsistency lies in the particular choice of feature space rather than the foundations of FID or P&R. It was recently discovered that classifiers trained using ImageNet [30] tend to base their decisions much more on texture than shape [10], while humans strongly focus on shape [23]. This is relevant in our context because

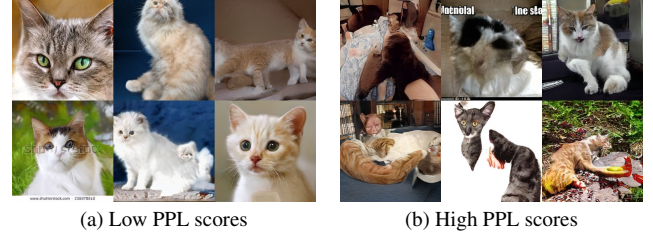


Figure 4. Connection between perceptual path length and image quality using baseline StyleGAN (config A) with LSUN CAT. (a) Random examples with low PPL ($\leq 10^{\text{th}}$ percentile). (b) Examples with high PPL ($\geq 90^{\text{th}}$ percentile). There is a clear correlation between PPL scores and semantic consistency of the images.

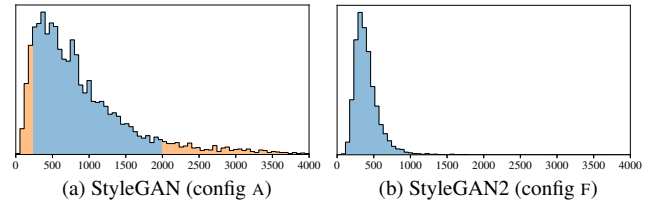


Figure 5. (a) Distribution of PPL scores of individual images generated using baseline StyleGAN (config A) with LSUN CAT (FID=8.53, PPL=924). The percentile ranges corresponding to Figure 4 are highlighted in orange. (b) StyleGAN2 (config F) improves the PPL distribution considerably (showing a snapshot with the same FID=8.53, PPL=387).

We observe a correlation between perceived image quality and perceptual path length (PPL) [21], a metric that was originally introduced for quantifying the smoothness of the mapping from a latent space to the output image by measuring average LPIPS distances [44] between generated images under small perturbations in latent space. Again consulting Figures 3 and 4 in the Supplement, a smaller PPL (smoother generator mapping) appears to correlate with higher over-

FID and P&R use high-level features from InceptionV3 [34] and VGG-16 [34], respectively, which were trained in this way and are thus expected to be biased towards texture detection. As such, images with, e.g., strong cat textures may appear more similar to each other than a human observer would agree, thus partially compromising density-based metrics (FID) and manifold coverage metrics (P&R).

all image quality, whereas other metrics are blind to the change. Figure 4 examines this correlation more closely through per-image PPL scores on LSUN CAT, computed by sampling the latent space around $\mathbf{w} \sim f(\mathbf{z})$. Low scores are indeed indicative of high-quality images, and vice versa. Figure 5a shows the corresponding histogram and reveals the long tail of the distribution. The overall PPL for the model is simply the expected value of these per-image PPL scores. We always compute PPL for the entire image, as opposed to Karras et al. [21] who use a smaller central crop.

It is not immediately obvious why a low PPL should correlate with image quality. We hypothesize that during training, as the discriminator penalizes broken images, the most direct way for the generator to improve is to effectively stretch the region of latent space that yields good images. This would lead to the low-quality images being squeezed into small latent space regions of rapid change. While this improves the average output quality in the short term, the accumulating distortions impair the training dynamics and consequently the final image quality.

Clearly, we cannot simply encourage minimal PPL since that would guide the generator toward a degenerate solution with zero recall. Instead, we will describe a new regularizer that aims for a smoother generator mapping without this drawback. As the resulting regularization term is somewhat expensive to compute, we first describe a general optimization that applies to any regularization technique.

3.1. Lazy regularization

Typically the main loss function (e.g., logistic loss [13]) and regularization terms (e.g., R_1 [25]) are written as a single expression and are thus optimized simultaneously. We observe that the regularization terms can be computed less frequently than the main loss function, thus greatly diminishing their computational cost and the overall memory usage. Table 1, row C shows that no harm is caused when R_1 regularization is performed only once every 16 minibatches, and we adopt the same strategy for our new regularizer as well. Appendix B gives implementation details.

3.2. Path length regularization

We would like to encourage that a fixed-size step in \mathcal{W} results in a non-zero, fixed-magnitude change in the image. We can measure the deviation from this ideal empirically by stepping into random directions in the image space and observing the corresponding \mathbf{w} gradients. These gradients should have close to an equal length regardless of \mathbf{w} or the image-space direction, indicating that the mapping from the latent space to image space is well-conditioned [28].

At a single $\mathbf{w} \in \mathcal{W}$, the local metric scaling properties of the generator mapping $g(\mathbf{w}) : \mathcal{W} \mapsto \mathcal{Y}$ are captured by the Jacobian matrix $\mathbf{J}_{\mathbf{w}} = \partial g(\mathbf{w}) / \partial \mathbf{w}$. Motivated by the desire to preserve the expected lengths of vectors regardless

of the direction, we formulate our regularizer as

$$\mathbb{E}_{\mathbf{w}, \mathbf{y} \sim \mathcal{N}(0, \mathbf{I})} (\|\mathbf{J}_{\mathbf{w}}^T \mathbf{y}\|_2 - a)^2, \quad (4)$$

where \mathbf{y} are random images with normally distributed pixel intensities, and $\mathbf{w} \sim f(\mathbf{z})$, where \mathbf{z} are normally distributed. We show in Appendix C that, in high dimensions, this prior is minimized when $\mathbf{J}_{\mathbf{w}}$ is orthogonal (up to a global scale) at any \mathbf{w} . An orthogonal matrix preserves lengths and introduces no squeezing along any dimension.

To avoid explicit computation of the Jacobian matrix, we use the identity $\mathbf{J}_{\mathbf{w}}^T \mathbf{y} = \nabla_{\mathbf{w}}(g(\mathbf{w}) \cdot \mathbf{y})$, which is efficiently computable using standard backpropagation [5]. The constant a is set dynamically during optimization as the long-running exponential moving average of the lengths $\|\mathbf{J}_{\mathbf{w}}^T \mathbf{y}\|_2$, allowing the optimization to find a suitable global scale by itself.

Our regularizer is closely related to the Jacobian clamping regularizer presented by Odena et al. [28]. Practical differences include that we compute the products $\mathbf{J}_{\mathbf{w}}^T \mathbf{y}$ analytically whereas they use finite differences for estimating $\mathbf{J}_{\mathbf{w}} \delta$ with $\mathcal{Z} \ni \delta \sim \mathcal{N}(0, \mathbf{I})$. It should be noted that spectral normalization [26] of the generator [40] only constrains the largest singular value, posing no constraints on the others and hence not necessarily leading to better conditioning. We find that enabling spectral normalization in addition to our contributions—or instead of them—invariably compromises FID, as detailed in Appendix E.

In practice, we notice that path length regularization leads to more reliable and consistently behaving models, making architecture exploration easier. We also observe that the smoother generator is significantly easier to invert (Section 5). Figure 5b shows that path length regularization clearly tightens the distribution of per-image PPL scores, without pushing the mode to zero. However, Table 1, row D points toward a tradeoff between FID and PPL in datasets that are less structured than FFHQ.

4. Progressive growing revisited

Progressive growing [20] has been very successful in stabilizing high-resolution image synthesis, but it causes its own characteristic artifacts. The key issue is that the progressively grown generator appears to have a strong location preference for details; the accompanying video shows that when features like teeth or eyes should move smoothly over the image, they may instead remain stuck in place before jumping to the next preferred location. Figure 6 shows a related artifact. We believe the problem is that in progressive growing each resolution serves momentarily as the output resolution, forcing it to generate maximal frequency details, which then leads to the trained network to have excessively high frequencies in the intermediate layers, compromising shift invariance [43]. Appendix A shows an example. These

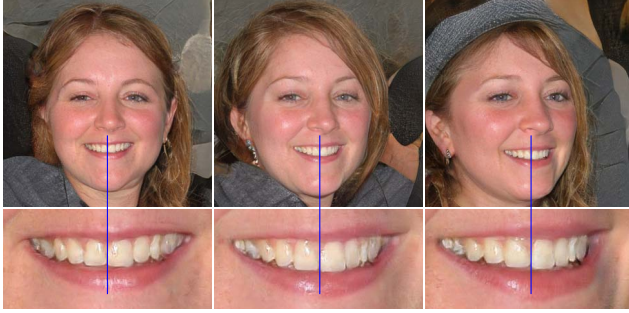


Figure 6. Progressive growing leads to “phase” artifacts. In this example the teeth do not follow the pose but stay aligned to the camera, as indicated by the blue line.

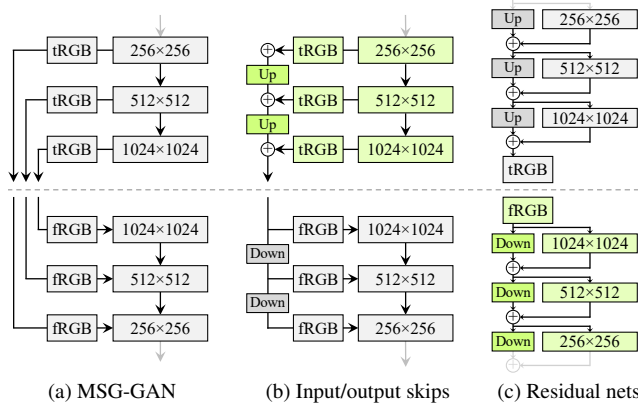


Figure 7. Three generator (above the dashed line) and discriminator architectures. **Up** and **Down** denote bilinear up and down-sampling, respectively. In residual networks these also include 1×1 convolutions to adjust the number of feature maps. **tRGB** and **fRGB** convert between RGB and high-dimensional per-pixel data. Architectures used in configs E and F are shown in green.

issues prompt us to search for an alternative formulation that would retain the benefits of progressive growing without the drawbacks.

4.1. Alternative network architectures

While StyleGAN uses simple feedforward designs in the generator (synthesis network) and discriminator, there is a vast body of work dedicated to the study of better network architectures. Skip connections [29, 19], residual networks [15, 14, 26], and hierarchical methods [6, 41, 42] have proven highly successful also in the context of generative methods. As such, we decided to re-evaluate the network design of StyleGAN and search for an architecture that produces high-quality images without progressive growing.

Figure 7a shows MSG-GAN [19], which connects the matching resolutions of the generator and discriminator using multiple skip connections. The MSG-GAN generator is modified to output a mipmap [37] instead of an image, and a similar representation is computed for each real im-

FFHQ	D original		D input skips		D residual	
	FID	PPL	FID	PPL	FID	PPL
G original	4.32	265	4.18	235	3.58	269
G output skips	4.33	169	3.77	127	3.31	125
G residual	4.35	203	3.96	229	3.79	243

LSUN Car	D original		D input skips		D residual	
	FID	PPL	FID	PPL	FID	PPL
G original	3.75	905	3.23	758	3.25	802
G output skips	3.77	544	3.86	316	3.19	471
G residual	3.93	981	3.40	667	2.66	645

Table 2. Comparison of generator and discriminator architectures without progressive growing. The combination of generator with output skips and residual discriminator corresponds to configuration E in the main result table.

age as well. In Figure 7b we simplify this design by up-sampling and summing the contributions of RGB outputs corresponding to different resolutions. In the discriminator, we similarly provide the downsampled image to each resolution block of the discriminator. We use bilinear filtering in all up and downsampling operations. In Figure 7c we further modify the design to use residual connections.³ This design is similar to LAPGAN [6] without the per-resolution discriminators employed by Denton et al.

Table 2 compares three generator and three discriminator architectures: original feedforward networks as used in StyleGAN, skip connections, and residual networks, all trained without progressive growing. FID and PPL are provided for each of the 9 combinations. We can see two broad trends: skip connections in the generator drastically improve PPL in all configurations, and a residual discriminator network is clearly beneficial for FID. The latter is perhaps not surprising since the structure of discriminator resembles classifiers where residual architectures are known to be helpful. However, a residual architecture was harmful in the generator—the lone exception was FID in LSUN CAR when both networks were residual.

For the rest of the paper we use a skip generator and a residual discriminator, without progressive growing. This corresponds to configuration E in Table 1, and it significantly improves FID and PPL.

4.2. Resolution usage

The key aspect of progressive growing, which we would like to preserve, is that the generator will initially focus on low-resolution features and then slowly shift its attention to finer details. The architectures in Figure 7 make it possible for the generator to first output low resolution images that are not affected by the higher-resolution layers in a significant way, and later shift the focus to the higher-resolution

³In residual network architectures, the addition of two paths leads to a doubling of signal variance, which we cancel by multiplying with $1/\sqrt{2}$. This is crucial for our networks, whereas in classification resnets [15] the issue is typically hidden by batch normalization.

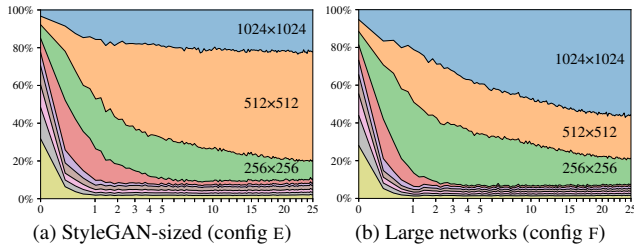


Figure 8. Contribution of each resolution to the output of the generator as a function of training time. The vertical axis shows a breakdown of the relative standard deviations of different resolutions, and the horizontal axis corresponds to training progress, measured in millions of training images shown to the discriminator. We can see that in the beginning the network focuses on low-resolution images and progressively shifts its focus on larger resolutions as training progresses. In (a) the generator basically outputs a 512^2 image with some minor sharpening for 1024^2 , while in (b) the larger network focuses more on the high-resolution details.

layers as the training proceeds. Since this is not enforced in any way, the generator will do it only if it is beneficial. To analyze the behavior in practice, we need to quantify how strongly the generator relies on particular resolutions over the course of training.

Since the skip generator (Figure 7b) forms the image by explicitly summing RGB values from multiple resolutions, we can estimate the relative importance of the corresponding layers by measuring how much they contribute to the final image. In Figure 8a, we plot the standard deviation of the pixel values produced by each tRGB layer as a function of training time. We calculate the standard deviations over 1024 random samples of \mathbf{w} and normalize the values so that they sum to 100%.

At the start of training, we can see that the new skip generator behaves similar to progressive growing—now achieved without changing the network topology. It would thus be reasonable to expect the highest resolution to dominate towards the end of the training. The plot, however, shows that this fails to happen in practice, which indicates that the generator may not be able to “fully utilize” the target resolution. To verify this, we inspected the generated images manually and noticed that they generally lack some of the pixel-level detail that is present in the training data—the images could be described as being sharpened versions of 512^2 images instead of true 1024^2 images.

This leads us to hypothesize that there is a capacity problem in our networks, which we test by doubling the number of feature maps in the highest-resolution layers of both networks.⁴ This brings the behavior more in line with expecta-

⁴We double the number of feature maps in resolutions 64^2 – 1024^2 while keeping other parts of the networks unchanged. This increases the total number of trainable parameters in the generator by 22% ($25\text{M} \rightarrow 30\text{M}$) and in the discriminator by 21% ($24\text{M} \rightarrow 29\text{M}$).

Dataset	Resolution	StyleGAN (A)		StyleGAN2 (F)	
		FID	PPL	FID	PPL
LSUN CAR	512×384	3.27	1485	2.32	416
LSUN CAT	256×256	8.53	924	6.93	439
LSUN CHURCH	256×256	4.21	742	3.86	342
LSUN HORSE	256×256	3.83	1405	3.43	338

Table 3. Improvement in LSUN datasets measured using FID and PPL. We trained CAR for 57M images, CAT for 88M, CHURCH for 48M, and HORSE for 100M images.

tions: Figure 8b shows a significant increase in the contribution of the highest-resolution layers, and Table 1, row F shows that FID and Recall improve markedly. The last row shows that baseline StyleGAN also benefits from additional capacity, but its quality remains far below StyleGAN2.

Table 3 compares StyleGAN and StyleGAN2 in four LSUN categories, again showing clear improvements in FID and significant advances in PPL. It is possible that further increases in the size could provide additional benefits.

5. Projection of images to latent space

Inverting the synthesis network g is an interesting problem that has many applications. Manipulating a given image in the latent feature space requires finding a matching latent code \mathbf{w} for it first. Previous research [1, 9] suggests that instead of finding a common latent code \mathbf{w} , the results improve if a separate \mathbf{w} is chosen for each layer of the generator. The same approach was used in an early encoder implementation [27]. While extending the latent space in this fashion finds a closer match to a given image, it also enables projecting arbitrary images that should have no latent representation. Instead, we concentrate on finding latent codes in the original, unextended latent space, as these correspond to images that the generator could have produced.

Our projection method differs from previous methods in two ways. First, we add ramped-down noise to the latent code during optimization in order to explore the latent space more comprehensively. Second, we also optimize the stochastic noise inputs of the StyleGAN generator, regularizing them to ensure they do not end up carrying coherent signal. The regularization is based on enforcing the autocorrelation coefficients of the noise maps to match those of unit Gaussian noise over multiple scales. Details of our projection method can be found in Appendix D.

5.1. Attribution of generated images

Detection of manipulated or generated images is a very important task. At present, classifier-based methods can quite reliably detect generated images, regardless of their exact origin [24, 39, 35, 45, 36]. However, given the rapid pace of progress in generative methods, this may not be a lasting situation. Besides general detection of fake images, we may also consider a more limited form of the problem:



Figure 9. Example images and their projected and re-synthesized counterparts. For each configuration, top row shows the target images and bottom row shows the synthesis of the corresponding projected latent vector and noise inputs. With the baseline StyleGAN, projection often finds a reasonably close match for generated images, but especially the backgrounds differ from the originals. The images generated using StyleGAN2 can be projected almost perfectly back into generator inputs, while projected real images (from the training set) show clear differences to the originals, as expected. All tests were done using the same projection method and hyperparameters.

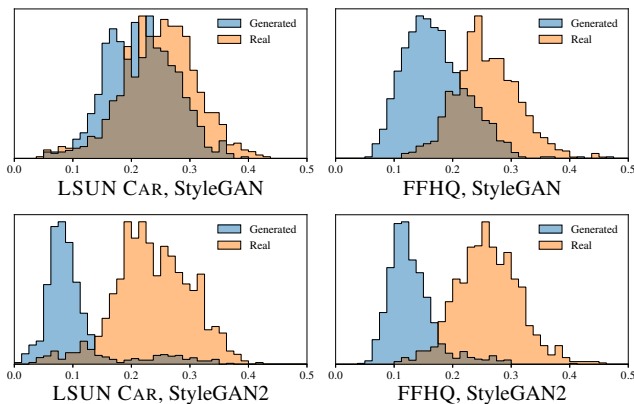


Figure 10. LPIPS distance histograms between original and projected images for generated (blue) and real images (orange). Despite the higher image quality of our improved generator, it is much easier to project the generated images into its latent space \mathcal{W} . The same projection method was used in all cases.

being able to attribute a fake image to its specific source [2]. With StyleGAN, this amounts to checking if there exists a $\mathbf{w} \in \mathcal{W}$ that re-synthesis the image in question.

We measure how well the projection succeeds by computing the LPIPS [44] distance between original and re-synthesized image as $D_{\text{LPIPS}}[\mathbf{x}, g(\tilde{g}^{-1}(\mathbf{x}))]$, where \mathbf{x} is the image being analyzed and \tilde{g}^{-1} denotes the approximate projection operation. Figure 10 shows histograms of these distances for LSUN CAR and FFHQ datasets using the original StyleGAN and StyleGAN2, and Figure 9 shows example projections. The images generated using StyleGAN2 can be projected into \mathcal{W} so well that they can be almost unambiguously attributed to the generating network. However, with the original StyleGAN, even though it should technically be possible to find a matching latent code, it appears that the mapping from \mathcal{W} to images is too complex for this to succeed reliably in practice. We find it encouraging that StyleGAN2 makes source attribution easier even though the image quality has improved significantly.

6. Conclusions and future work

We have identified and fixed several image quality issues in StyleGAN, improving the quality further and considerably advancing the state of the art in several datasets. In some cases the improvements are more clearly seen in motion, as demonstrated in the accompanying video. Appendix A includes further examples of results obtainable using our method. Despite the improved quality, StyleGAN2 makes it easier to attribute a generated image to its source.

Training performance has also improved. At 1024² resolution, the original StyleGAN (config A in Table 1) trains at 37 images per second on NVIDIA DGX-1 with 8 Tesla V100 GPUs, while our config E trains 40% faster at 61 img/s. Most of the speedup comes from simplified dataflow due to weight demodulation, lazy regularization, and code optimizations. StyleGAN2 (config F, larger networks) trains at 31 img/s, and is thus only slightly more expensive to train than original StyleGAN. Its total training time was 9 days for FFHQ and 13 days for LSUN CAR.

The entire project, including all exploration, consumed 132 MWh of electricity, of which 0.68 MWh went into training the final FFHQ model. In total, we used about 51 single-GPU years of computation (Volta class GPU). A more detailed discussion is available in Appendix F.

In the future, it could be fruitful to study further improvements to the path length regularization, e.g., by replacing the pixel-space L_2 distance with a data-driven feature-space metric. Considering the practical deployment of GANs, we feel that it will be important to find new ways to reduce the training data requirements. This is especially crucial in applications where it is infeasible to acquire tens of thousands of training samples, and with datasets that include a lot of intrinsic variation.

Acknowledgements We thank Ming-Yu Liu for an early review, Timo Viitanen for help with the public release, David Luebke for in-depth discussions and helpful comments, and Tero Kuosmanen for technical support with the compute infrastructure.

References

- [1] Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2StyleGAN: How to embed images into the StyleGAN latent space? In *ICCV*, 2019. 7
- [2] Michael Albright and Scott McCloskey. Source generator attribution via inversion. In *CVPR Workshops*, 2019. 8
- [3] Carl Bergstrom and Jevin West. Which face is real? <http://www.whichfaceisreal.com/learn.html>, Accessed November 15, 2019. 1
- [4] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale GAN training for high fidelity natural image synthesis. *CoRR*, abs/1809.11096, 2018. 1
- [5] Yann N. Dauphin, Harm de Vries, and Yoshua Bengio. Equilibrated adaptive learning rates for non-convex optimization. *CoRR*, abs/1502.04390, 2015. 5
- [6] Emily L. Denton, Soumith Chintala, Arthur Szlam, and Robert Fergus. Deep generative image models using a Laplacian pyramid of adversarial networks. *CoRR*, abs/1506.05751, 2015. 6
- [7] Vincent Dumoulin, Ethan Perez, Nathan Schucher, Florian Strub, Harm de Vries, Aaron Courville, and Yoshua Bengio. Feature-wise transformations. *Distill*, 2018. <https://distill.pub/2018/feature-wise-transformations>. 1
- [8] Vincent Dumoulin, Jonathon Shlens, and Manjunath Kudlur. A learned representation for artistic style. *CoRR*, abs/1610.07629, 2016. 1
- [9] Aviv Gabbay and Yedid Hoshen. Style generator inversion for image enhancement and animation. *CoRR*, abs/1906.11880, 2019. 7
- [10] Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A. Wichmann, and Wieland Brendel. ImageNet-trained CNNs are biased towards texture; increasing shape bias improves accuracy and robustness. *CoRR*, abs/1811.12231, 2018. 1, 4
- [11] Golnaz Ghiasi, Honglak Lee, Manjunath Kudlur, Vincent Dumoulin, and Jonathon Shlens. Exploring the structure of a real-time, arbitrary neural artistic stylization network. *CoRR*, abs/1705.06830, 2017. 1
- [12] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, pages 249–256, 2010. 3
- [13] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. In *NIPS*, 2014. 1, 5
- [14] Ishaan Gulrajani, Faruk Ahmed, Martín Arjovsky, Vincent Dumoulin, and Aaron C. Courville. Improved training of Wasserstein GANs. *CoRR*, abs/1704.00028, 2017. 6
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *CoRR*, abs/1512.03385, 2015. 6
- [16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on ImageNet classification. *CoRR*, abs/1502.01852, 2015. 3
- [17] 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 *Proc. NIPS*, pages 6626–6637, 2017. 1
- [18] Xun Huang and Serge J. Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. *CoRR*, abs/1703.06868, 2017. 1
- [19] Animesh Karnewar and Oliver Wang. MSG-GAN: multi-scale gradients for generative adversarial networks. In *Proc. CVPR*, 2020. 6
- [20] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of GANs for improved quality, stability, and variation. *CoRR*, abs/1710.10196, 2017. 1, 5
- [21] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proc. CVPR*, 2018. 1, 2, 4, 5
- [22] Tuomas Kynkäänniemi, Tero Karras, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Improved precision and recall metric for assessing generative models. In *Proc. NeurIPS*, 2019. 1, 2, 4
- [23] Barbara Landau, Linda B. Smith, and Susan S. Jones. The importance of shape in early lexical learning. *Cognitive Development*, 3(3), 1988. 4
- [24] Haodong Li, Han Chen, Bin Li, and Shunquan Tan. Can forensic detectors identify GAN generated images? In *Proc. Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*, 2018. 7
- [25] Lars Mescheder, Andreas Geiger, and Sebastian Nowozin. Which training methods for GANs do actually converge? *CoRR*, abs/1801.04406, 2018. 5
- [26] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative adversarial networks. *CoRR*, abs/1802.05957, 2018. 1, 5, 6
- [27] Dmitry Nikitko. StyleGAN – Encoder for official TensorFlow implementation. <https://github.com/Puzer/stylegan-encoder/>, 2019. 7
- [28] Augustus Odena, Jacob Buckman, Catherine Olsson, Tom B. Brown, Christopher Olah, Colin Raffel, and Ian Goodfellow. Is generator conditioning causally related to GAN performance? *CoRR*, abs/1802.08768, 2018. 5
- [29] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-Net: Convolutional networks for biomedical image segmentation. In *Proc. Medical Image Computing and Computer-Assisted Intervention (MICCAI)*, pages 234–241, 2015. 6
- [30] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael S. Bernstein, Alexander C. Berg, and Fei-Fei Li. ImageNet large scale visual recognition challenge. In *Proc. CVPR*, 2015. 4
- [31] Mehdi S. M. Sajjadi, Olivier Bachem, Mario Lucic, Olivier Bousquet, and Sylvain Gelly. Assessing generative models via precision and recall. *CoRR*, abs/1806.00035, 2018. 1
- [32] Tim Salimans and Diederik P. Kingma. Weight normalization: A simple reparameterization to accelerate training of deep neural networks. *CoRR*, abs/1602.07868, 2016. 3
- [33] Yujun Shen, Jinjin Gu, Xiaou Tang, and Bolei Zhou. Interpreting the latent space of GANs for semantic face editing. *CoRR*, abs/1907.10786, 2019. 1

- [34] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556, 2014. 1, 4
- [35] Run Wang, Lei Ma, Felix Juefei-Xu, Xiaofei Xie, Jian Wang, and Yang Liu. FakeSpotter: A simple baseline for spotting AI-synthesized fake faces. *CoRR*, abs/1909.06122, 2019. 7
- [36] Sheng-Yu Wang, Oliver Wang, Richard Zhang, Andrew Owens, and Alexei A. Efros. CNN-generated images are surprisingly easy to spot... for now. *CoRR*, abs/1912.11035, 2019. 7
- [37] Lance Williams. Pyramidal parametrics. *SIGGRAPH Comput. Graph.*, 17(3):1–11, 1983. 6
- [38] Sitao Xiang and Hao Li. On the effects of batch and weight normalization in generative adversarial networks. *CoRR*, abs/1704.03971, 2017. 3
- [39] Ning Yu, Larry Davis, and Mario Fritz. Attributing fake images to GANs: Analyzing fingerprints in generated images. *CoRR*, abs/1811.08180, 2018. 7
- [40] Han Zhang, Ian Goodfellow, Dimitris Metaxas, and Augustus Odena. Self-attention generative adversarial networks. *CoRR*, abs/1805.08318, 2018. 5
- [41] Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiao lei Huang, Xiaogang Wang, and Dimitris N. Metaxas. StackGAN: text to photo-realistic image synthesis with stacked generative adversarial networks. In *ICCV*, 2017. 6
- [42] Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiaogang Wang, Xiao lei Huang, and Dimitris N. Metaxas. StackGAN++: realistic image synthesis with stacked generative adversarial networks. *CoRR*, abs/1710.10916, 2017. 6
- [43] Richard Zhang. Making convolutional networks shift-invariant again. In *Proc. ICML*, 2019. 5
- [44] Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proc. CVPR*, 2018. 4, 8
- [45] Xu Zhang, Svebor Karaman, and Shih-Fu Chang. Detecting and simulating artifacts in GAN fake images. *CoRR*, abs/1907.06515, 2019. 7