
Alias-Free Generative Adversarial Networks

Tero Karras
NVIDIA

Miika Aittala
NVIDIA

Samuli Laine
NVIDIA

Erik Härkönen*
Aalto University and NVIDIA

Janne Hellsten
NVIDIA

Jaakko Lehtinen
NVIDIA and Aalto University

Timo Aila
NVIDIA

Abstract

We observe that despite their hierarchical convolutional nature, the synthesis process of typical generative adversarial networks depends on absolute pixel coordinates in an unhealthy manner. This manifests itself as, e.g., detail appearing to be glued to image coordinates instead of the surfaces of depicted objects. We trace the root cause to careless signal processing that causes aliasing in the generator network. Interpreting all signals in the network as continuous, we derive generally applicable, small architectural changes that guarantee that unwanted information cannot leak into the hierarchical synthesis process. The resulting networks match the FID of StyleGAN2 but differ dramatically in their internal representations, and they are fully equivariant to translation and rotation even at subpixel scales. Our results pave the way for generative models better suited for video and animation.

1 Introduction

The resolution and quality of images produced by generative adversarial networks (GAN) [21] have seen rapid improvement recently [31, 11, 33, 34]. They have been used for a variety of applications, including image editing [49, 55, 43, 22, 39, 3], domain translation [70, 37, 61, 42], and video generation [57, 15, 24]. While several ways of controlling the generative process have been found [8, 29, 10, 42, 25, 2, 7, 48, 6], the foundations of the synthesis process remain only partially understood.

In the real world, details of different scale tend to transform hierarchically. For instance, moving a head causes the nose to move, which in turn moves the skin pores on it. The structure of a typical GAN generator is analogous: coarse, low-resolution features are hierarchically refined by upsampling layers, locally mixed by convolutions, and new detail is introduced through nonlinearities. We observe that despite this superficial similarity, current GAN architectures do not synthesize images in a natural hierarchical manner: the coarse features mainly control the *presence* of finer features, but not their precise positions. Instead, much of the fine detail appears to be fixed in pixel coordinates. This disturbing “texture sticking” is clearly visible in latent interpolations (see Figure 1 and our accompanying videos on the project page <https://nvlabs.github.io/stylegan3>), breaking the illusion of a solid and coherent object moving in space. Our goal is an architecture that exhibits a more natural transformation hierarchy, where the exact sub-pixel position of each feature is exclusively inherited from the underlying coarse features.

*This work was done during an internship at NVIDIA.

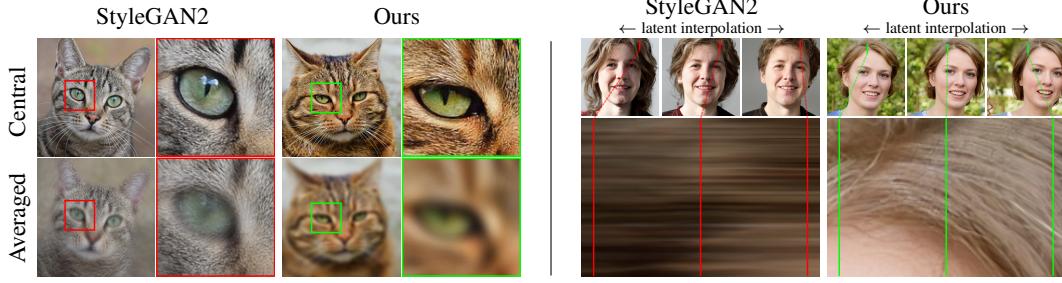


Figure 1: Examples of “texture sticking”. **Left:** The average of images generated from a small neighborhood around a central latent (top row). The intended result is uniformly blurry because all details should move together. However, with StyleGAN2 many details (e.g., fur) stick to the same pixel coordinates, showing unwanted sharpness. **Right:** From a latent space interpolation (top row), we extract a short vertical segment of pixels from each generated image and stack them horizontally (bottom). The desired result is hairs moving in animation, creating a time-varying field. With StyleGAN2 the hairs mostly stick to the same coordinates, creating horizontal streaks instead.

It turns out that current networks can partially bypass the ideal hierarchical construction by drawing on unintentional positional references available to the intermediate layers through image borders [28, 35, 66], per-pixel noise inputs [33] and positional encodings, and aliasing [5, 69]. Aliasing, despite being a subtle and critical issue [44], has received little attention in the GAN literature. We identify two sources for it: 1) faint after-images of the pixel grid resulting from non-ideal upsampling filters² such as nearest, bilinear, or strided convolutions, and 2) the pointwise application of nonlinearities such as ReLU [60] or swish [47]. We find that the network has the means and motivation to amplify even the slightest amount of aliasing and combining it over multiple scales allows it to build a basis for texture motifs that are fixed in screen coordinates. This holds for all filters commonly used in deep learning [69, 59], and even high-quality filters used in image processing.

How, then, do we eliminate the unwanted side information and thereby stop the network from using it? While borders can be solved by simply operating on slightly larger images, aliasing is much harder. We begin by noting that aliasing is most naturally treated in the classical Shannon-Nyquist signal processing framework, and switch focus to bandlimited functions on a continuous domain that are merely represented by discrete sample grids. Now, successful elimination of all sources of positional references means that details can be generated equally well regardless of pixel coordinates, which in turn is equivalent to enforcing continuous *equivariance* to sub-pixel translation (and optionally rotation) in all layers. To achieve this, we describe a comprehensive overhaul of all signal processing aspects of the StyleGAN2 generator [34]. Our contributions include the surprising finding that current upsampling filters are simply not aggressive enough in suppressing aliasing, and that extremely high-quality filters with over 100dB attenuation are required. Further, we present a principled solution to aliasing caused by pointwise nonlinearities [5] by considering their effect in the continuous domain and appropriately low-pass filtering the results. We also show that after the overhaul, a model based on 1×1 convolutions yields a strong, rotation equivariant generator.

Once aliasing is adequately suppressed to force the model to implement more natural hierarchical refinement, its mode of operation changes drastically: the emergent internal representations now include coordinate systems that allow details to be correctly attached to the underlying surfaces. This promises significant improvements to models that generate video and animation. The new StyleGAN3 generator matches StyleGAN2 in terms of FID [26], while being slightly heavier computationally. Our implementation and pre-trained models are available at <https://github.com/NVlabs/stylegan3>

Several recent works have studied the lack of translation equivariance in CNNs, mainly in the context of classification [28, 35, 66, 5, 38, 69, 12, 71, 59]. We significantly expand upon the antialiasing measures in this literature and show that doing so induces a fundamentally altered image generation behavior. Group-equivariant CNNs aim to generalize the efficiency benefits of translational weight

²Consider nearest neighbor upsampling. If we upsample a 4×4 image to 8×8 , the original pixels will be clearly visible, allowing one to reliably distinguish between even and odd pixels. Since the same is true on all scales, this (leaked) information makes it possible to reconstruct even the absolute pixel coordinates. With better filters such as bilinear or bicubic, the clues get less pronounced, but are nevertheless evident for the generator.

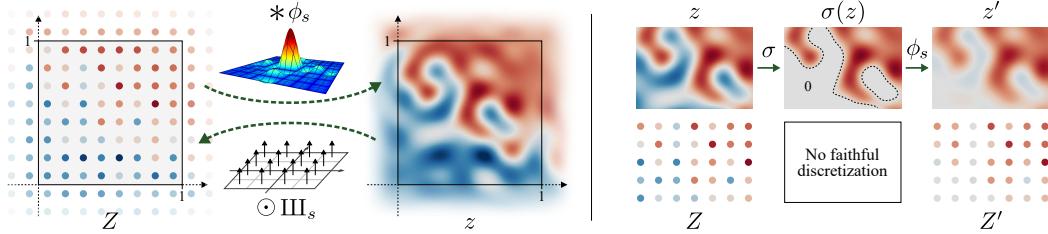


Figure 2: **Left:** Discrete representation Z and continuous representation z are related to each other via convolution with ideal interpolation filter ϕ_s and pointwise multiplication with Dirac comb III_s . **Right:** Nonlinearity σ , ReLU in this example, may produce arbitrarily high frequencies in the continuous-domain $\sigma(z)$. Low-pass filtering via ϕ_s is necessary to ensure that Z' captures the result.

sharing to, e.g., rotation [16, 65, 63, 62] and scale [64]. Our 1×1 convolutions can be seen as an instance of a continuously $E(2)$ -equivariant model [62] that remains compatible with, e.g., channel-wise ReLU nonlinearities and modulation. Dey et al. [17] apply 90° rotation-and-flip equivariant CNNs [16] to GANs and show improved data efficiency. Our work is complementary, and not motivated by efficiency. Recent implicit network [53, 56, 13] based GANs [4, 54] generate each pixel independently via similar 1×1 convolutions. While equivariant, these models do not help with texture sticking, as they do not use an upsampling hierarchy or implement a shallow non-antialiased one.

2 Equivariance via continuous signal interpretation

To begin our analysis of equivariance in CNNs, we shall first rethink our view of **what exactly is the signal that flows through a network**. Even though data may be stored as values in a pixel grid, we cannot naively hold these values to directly represent the signal. Doing so would prevent us from considering operations as trivial as translating the contents of a feature map by half a pixel.

According to the Nyquist–Shannon sampling theorem [51], a regularly sampled signal can represent any continuous signal containing frequencies between zero and half of the sampling rate. Let us consider a two-dimensional, discretely sampled feature map $Z[\mathbf{x}]$ that consists of a regular grid of Dirac impulses of varying magnitudes, spaced $1/s$ units apart where s is the sampling rate. This is analogous to an infinite two-dimensional grid of values.

Given $Z[\mathbf{x}]$ and s , the Whittaker–Shannon interpolation formula [51] states that the corresponding continuous representation $z(\mathbf{x})$ is obtained by convolving the discretely sampled Dirac grid $Z[\mathbf{x}]$ with an ideal interpolation filter ϕ_s , i.e., $z(\mathbf{x}) = (\phi_s * Z)(\mathbf{x})$, where $*$ denotes continuous convolution and $\phi_s(\mathbf{x}) = \text{sinc}(sx_0) \cdot \text{sinc}(sx_1)$ using the signal processing convention of defining $\text{sinc}(x) = \sin(\pi x)/(\pi x)$. ϕ_s has a bandlimit of $s/2$ along the horizontal and vertical dimensions, ensuring that the resulting continuous signal captures all frequencies that can be represented with sampling rate s .

Conversion from the continuous to the discrete domain corresponds to sampling the continuous signal $z(\mathbf{x})$ at the sampling points of $Z[\mathbf{x}]$ that we define to be offset by half the sample spacing to lie at the “pixel centers”, see Figure 2, left. This can be expressed as a pointwise multiplication with a two-dimensional Dirac comb $\text{III}_s(\mathbf{x}) = \sum_{X \in \mathbb{Z}^2} \delta(\mathbf{x} - (X + \frac{1}{2}))/s$.

We earmark the unit square $\mathbf{x} \in [0, 1]^2$ in $z(\mathbf{x})$ as our canvas for the signal of interest. In $Z[\mathbf{x}]$ there are s^2 discrete samples in this region, but the above convolution with ϕ_s means that values of $Z[\mathbf{x}]$ outside the unit square also influence $z(\mathbf{x})$ inside it. Thus storing an $s \times s$ -pixel feature map is not sufficient; in theory, we would need to store the entire infinite $Z[\mathbf{x}]$. As a practical solution, we store $Z[\mathbf{x}]$ as a two-dimensional array that covers a region slightly larger than the unit square (Section 3.2).

Having established correspondence between bandlimited, continuous feature maps $z(\mathbf{x})$ and discretely sampled feature maps $Z[\mathbf{x}]$, we can shift our focus away from the usual pixel-centric view of the signal. In the remainder of this paper, we shall interpret $z(\mathbf{x})$ as being the actual signal being operated on, and the discretely sampled feature map $Z[\mathbf{x}]$ as merely a convenient encoding for it.

Discrete and continuous representation of network layers Practical neural networks operate on the discretely sampled feature maps. Consider operation \mathbf{F} (convolution, nonlinearity, etc.) operating

on a discrete feature map: $Z' = \mathbf{F}(Z)$. The feature map has a corresponding continuous counterpart, so we also have a corresponding mapping in the continuous domain: $z' = \mathbf{f}(z)$. Now, an operation specified in one domain can be seen to perform a corresponding operation in the other domain:

$$\mathbf{f}(z) = \phi_{s'} * \mathbf{F}(\text{III}_s \odot z), \quad \mathbf{F}(Z) = \text{III}_{s'} \odot \mathbf{f}(\phi_s * Z), \quad (1)$$

where \odot denotes pointwise multiplication and s and s' are the input and output sampling rates. Note that in the latter case \mathbf{f} must not introduce frequency content beyond the output bandlimit $s'/2$.

2.1 Equivariant network layers

Operation \mathbf{f} is equivariant with respect to a spatial transformation \mathbf{t} of the 2D plane if it commutes with it in the continuous domain: $\mathbf{t} \circ \mathbf{f} = \mathbf{f} \circ \mathbf{t}$. We note that when inputs are bandlimited to $s/2$, an equivariant operation must not generate frequency content above the output bandlimit of $s'/2$, as otherwise no faithful discrete output representation exists.

We focus on two types of equivariance in this paper: translation and rotation. In the case of rotation the spectral constraint is somewhat stricter—rotating an image corresponds to rotating the spectrum, and in order to guarantee the bandlimit in both horizontal and vertical direction, the spectrum must be limited to a disc with radius $s/2$. This applies to both the initial network input as well as the bandlimiting filters used for downsampling, as will be described later.

We now consider the primitive operations in a typical generator network: convolution, upsampling, downsampling, and nonlinearity. Without loss of generality, we discuss the operations acting on a single feature map: pointwise linear combination of features has no effect on the analysis.

Convolution Consider a standard convolution with a discrete kernel K . We can interpret K as living in the same grid as the input feature map, with sampling rate s . The discrete-domain operation is simply $\mathbf{F}_{\text{conv}}(Z) = K * Z$, and we obtain the corresponding continuous operation from Eq. 1:

$$\mathbf{f}_{\text{conv}}(z) = \phi_s * (K * (\text{III}_s \odot z)) = K * (\phi_s * (\text{III}_s \odot z)) = K * z \quad (2)$$

due to commutativity of convolution and the fact that discretization followed by convolution with ideal low-pass filter, both with same sampling rate s , is an identity operation, i.e., $\phi_s * (\text{III}_s \odot z) = z$. In other words, the convolution operates by continuously sliding the discretized kernel over the continuous representation of the feature map. This convolution introduces no new frequencies, so the bandlimit requirements for both translation and rotation equivariance are trivially fulfilled.

Convolution also commutes with translation in the continuous domain, and thus the operation is equivariant to translation. For rotation equivariance, the discrete kernel K needs to be radially symmetric. We later show in Section 3.2 that trivially symmetric 1×1 convolution kernels are, despite their simplicity, a viable choice for rotation equivariant generative networks.

Upsampling and downsampling Ideal upsampling does not modify the continuous representation. Its only purpose is to increase the output sampling rate ($s' > s$) to add headroom in the spectrum where subsequent layers may introduce additional content. Translation and rotation equivariance follow directly from upsampling being an identity operation in the continuous domain. With $\mathbf{f}_{\text{up}}(z) = z$, the discrete operation according to Eq. 1 is $\mathbf{F}_{\text{up}}(Z) = \text{III}_{s'} \odot (\phi_s * Z)$. If we choose $s' = ns$ with integer n , this operation can be implemented by first interleaving Z with zeros to increase its sampling rate and then convolving it with a discretized filter $\text{III}_{s'} \odot \phi_s$.

In downsampling, we must low-pass filter z to remove frequencies above the output bandlimit, so that the signal can be represented faithfully in the coarser discretization. The operation in continuous domain is $\mathbf{f}_{\text{down}}(z) = \psi_{s'} * z$, where an ideal low-pass filter $\psi_s := s^2 \cdot \phi_s$ is simply the corresponding interpolation filter normalized to unit mass. The discrete counterpart is $\mathbf{F}_{\text{down}}(Z) = \text{III}_{s'} \odot (\psi_{s'} * (\phi_s * Z)) = 1/s^2 \cdot \text{III}_{s'} \odot (\psi_{s'} * \psi_s * Z) = (s'/s)^2 \cdot \text{III}_{s'} \odot (\phi_{s'} * Z)$. The latter equality follows from $\psi_s * \psi_{s'} = \psi_{\min(s,s')}$. Similar to upsampling, downsampling by an integer fraction can be implemented with a discrete convolution followed by dropping sample points. Translation equivariance follows automatically from the commutativity of $\mathbf{f}_{\text{down}}(z)$ with translation, but for rotation equivariance we must replace $\phi_{s'}$ with a radially symmetric filter with disc-shaped frequency response. The ideal such filter [9] is given by $\phi_s^\circ(\mathbf{x}) = \text{jinc}(s\|\mathbf{x}\|) = 2J_1(\pi s\|\mathbf{x}\|)/(\pi s\|\mathbf{x}\|)$, where J_1 is the first order Bessel function of the first kind.

Configuration	FID↓	EQ-T↑	EQ-R↑	Parameter	FID↓	EQ-T↑	EQ-R↑	Time	Mem.
A StyleGAN2	5.14	—	—	Filter size $n = 4$	4.72	57.49	39.70	0.84 ×	0.99 ×
B + Fourier features	4.79	16.23	10.81	* Filter size $n = 6$	4.50	66.65	40.48	1.00×	1.00×
C + No noise inputs	4.54	15.81	10.84	Filter size $n = 8$	4.66	65.57	42.09	1.18×	1.01×
D + Simplified generator	5.21	19.47	10.41	Upsampling $m = 1$	4.38	39.96	36.42	0.65 ×	0.87 ×
E + Boundaries & upsampling	6.02	24.62	10.97	* Upsampling $m = 2$	4.50	66.65	40.48	1.00×	1.00×
F + Filtered nonlinearities	6.35	30.60	10.81	Upsampling $m = 4$	4.57	74.21	40.97	2.31×	1.62×
G + Non-critical sampling	4.78	43.90	10.84	Stopband $f_{t,0} = 2^{1.5}$	4.62	51.10	29.14	0.86 ×	0.90 ×
H + Transformed Fourier features	4.64	45.20	10.61	* Stopband $f_{t,0} = 2^{2.1}$	4.50	66.65	40.48	1.00×	1.00×
T + Flexible layers (StyleGAN3-T)	4.62	63.01	13.12	Stopband $f_{t,0} = 2^{3.1}$	4.68	73.13	41.63	1.36×	1.25×
R + Rotation equiv. (StyleGAN3-R)	4.50	66.65	40.48						

Figure 3: Results for FFHQ-U (unaligned FFHQ) at 256². **Left:** Training configurations. FID is computed between 50k generated images and all training images [26, 32]; lower is better. EQ-T and EQ-R are our equivariance metrics in decibels (dB); higher is better. **Right:** Parameter ablations using our final configuration (R) for the filter’s support, magnification around nonlinearities, and the minimum stopband frequency at the first layer. * indicates our default choices.

Nonlinearity Applying a pointwise nonlinearity σ in the discrete domain does not commute with fractional translation or rotation. However, in the continuous domain, any pointwise function commutes trivially with geometric transformations and is thus equivariant to translation and rotation. Fulfilling the bandlimit constraint is another question—applying, e.g., ReLU in the continuous domain may introduce arbitrarily high frequencies that cannot be represented in the output.

A natural solution is to eliminate the offending high-frequency content by convolving the continuous result with the ideal low-pass filter ψ_s . Then, the continuous representation of the nonlinearity becomes $\mathbf{f}_\sigma(z) = \psi_s * \sigma(z) = s^2 \cdot \phi_s * \sigma(z)$ and the discrete counterpart is $\mathbf{F}_\sigma(Z) = s^2 \cdot \text{III}_s \odot (\phi_s * \sigma(\phi_s * Z))$ (see Figure 2, right). This discrete operation cannot be realized without temporarily entering the continuous representation. We approximate this by upsampling the signal, applying the nonlinearity in the higher resolution, and downsampling it afterwards. Even though the nonlinearity is still performed in the discrete domain, we have found that only a 2× temporary resolution increase is sufficient for high-quality equivariance. For rotation equivariance, we must use the radially symmetric interpolation filter ϕ_s° in the downsampling step, as discussed above.

Note that nonlinearity is the only operation capable of generating novel frequencies in our formulation, and that we can limit the range of these novel frequencies by applying a reconstruction filter with a lower cutoff than $s/2$ before the final discretization operation. This gives us precise control over how much new information is introduced by each layer of a generator network (Section 3.2).

3 Practical application to generator network

We will now apply the theoretical ideas from the previous section in practice, by converting the well-established StyleGAN2 [34] generator to be fully equivariant to translation and rotation. We will introduce the necessary changes step-by-step, evaluating their impact in Figure 3. The discriminator remains unchanged in our experiments.

The StyleGAN2 generator consists of two parts. First, a *mapping network* transforms an initial, normally distributed latent to an intermediate latent code $\mathbf{w} \sim \mathcal{W}$. Then, a *synthesis network* \mathbf{G} starts from a learned 4×4×512 constant Z_0 and applies a sequence of N layers—consisting of convolutions, nonlinearities, upsampling, and per-pixel noise—to produce an output image $Z_N = \mathbf{G}(Z_0; \mathbf{w})$. The intermediate latent code \mathbf{w} controls the modulation of the convolution kernels in \mathbf{G} . The layers follow a rigid 2× upsampling schedule, where two layers are executed at each resolution and the number of feature maps is halved after each upsampling. Additionally, StyleGAN2 employs skip connections, mixing regularization [33], and path length regularization.

Our goal is to make every layer of \mathbf{G} equivariant w.r.t. the continuous signal, so that all finer details transform together with the coarser features of a local neighborhood. If this succeeds, the entire network becomes similarly equivariant. In other words, we aim to make the *continuous* operation \mathbf{g} of the synthesis network equivariant w.r.t. transformations \mathbf{t} (translations and rotations) applied on the continuous input z_0 : $\mathbf{g}(\mathbf{t}[z_0]; \mathbf{w}) = \mathbf{t}[\mathbf{g}(z_0; \mathbf{w})]$. To evaluate the impact of various architectural changes and practical approximations, we need a way to measure how well the network implements the equivariances. For translation equivariance, we report the peak signal-to-noise ratio (PSNR) in decibels (dB) between two sets of images, obtained by translating the input and output of the

synthesis network by a random amount, resembling the definition by Zhang [69]:

$$\text{EQ-T} = 10 \cdot \log_{10} \left(I_{\max}^2 / \mathbb{E}_{\mathbf{w} \sim \mathcal{W}, x \sim \mathcal{X}^2, p \sim \mathcal{V}, c \sim \mathcal{C}} \left[(\mathbf{g}(\mathbf{t}_x[z_0]; \mathbf{w})_c(p) - \mathbf{t}_x[\mathbf{g}(z_0; \mathbf{w})]_c(p))^2 \right] \right) \quad (3)$$

Each pair of images, corresponding to a different random choice of \mathbf{w} , is sampled at integer pixel locations p within their mutually valid region \mathcal{V} . Color channels c are processed independently, and the intended dynamic range of generated images $-1 \dots +1$ gives $I_{\max} = 2$. Operator \mathbf{t}_x implements spatial translation with 2D offset x , here drawn from distribution \mathcal{X}^2 of integer offsets. We define an analogous metric EQ-R for rotations, with the rotation angles drawn from $\mathcal{U}(0^\circ, 360^\circ)$. Appendix E gives implementation details and our accompanying videos highlight the practical relevance of different dB values.

3.1 Fourier features and baseline simplifications (configs B–D)

To facilitate exact continuous translation and rotation of the input z_0 , we replace the learned input constant in StyleGAN2 with Fourier features [56, 66], which also has the advantage of naturally defining a spatially infinite map. We sample the frequencies uniformly within the circular frequency band $f_c = 2$, matching the original 4×4 input resolution, and keep them fixed over the course of training. This change (configs A and B in Figure 3, left) slightly improves FID and, crucially, allows us to compute the equivariance metrics without having to approximate the operator \mathbf{t} . This baseline architecture is far from being equivariant; our accompanying videos show that the output images deteriorate drastically when the input features are translated or rotated from their original position.

Next, we remove the per-pixel noise inputs because they are strongly at odds with our goal of a natural transformation hierarchy, i.e., that the exact sub-pixel position of each feature is exclusively inherited from the underlying coarse features. While this change (config C) is approximately FID-neutral, it fails to improve the equivariance metrics when considered in isolation.

To further simplify the setup, we decrease the mapping network depth as recommended by Karras et al. [32] and disable mixing regularization and path length regularization [34]. Finally, we also eliminate the output skip connections. We hypothesize that their benefit is mostly related to gradient magnitude dynamics during training and address the underlying issue more directly using a simple normalization before each convolution. We track the exponential moving average $\sigma^2 = \mathbb{E}[x^2]$ over all pixels and feature maps during training, and divide the feature maps by $\sqrt{\sigma^2}$. In practice, we bake the division into the convolution weights to improve efficiency. These changes (config D) bring FID back to the level of original StyleGAN2, while leading to a slight improvement in translation equivariance.

3.2 Step-by-step redesign motivated by continuous interpretation

Boundaries and upsampling (config E) Our theory assumes an infinite spatial extent for the feature maps, which we approximate by maintaining a fixed-size margin around the target canvas, cropping to this extended canvas after each layer. This explicit extension is necessary as border padding is known to leak absolute image coordinates into the internal representations [28, 35, 66]. In practice, we have found a 10-pixel margin to be enough; further increase has no noticeable effect on the results.

Motivated by our theoretical model, we replace the bilinear $2 \times$ upsampling filter with a better approximation of the ideal low-pass filter. We use a windowed sinc filter with a relatively large Kaiser window [41] of size $n = 6$, meaning that each output pixel is affected by 6 input pixels in upsampling and each input pixel affects 6 output pixels in downsampling. Kaiser window is a particularly good choice for our purposes, because it offers explicit control over the transition band and attenuation (Figure 4a). In the remainder of this section, we specify the transition band explicitly and compute the remaining parameters using Kaiser’s original formulas (Appendix C). For now, we choose to employ *critical sampling* and set the filter cutoff $f_c = s/2$, i.e., exactly at the bandlimit, and transition band half-width $f_h = (\sqrt{2} - 1)(s/2)$. Recall that sampling rate s equals the width of the canvas in pixels, given our definitions in Section 2.

The improved handling of boundaries and upsampling (config E) leads to better translation equivariance. However, FID is compromised by 16%, probably because we started to constrain what the feature maps can contain. In a further ablation (Figure 3, right), smaller resampling filters ($n = 4$) hurt translation equivariance, while larger filters ($n = 8$) mainly increase training time.

Filtered nonlinearities (config F) Our theoretical treatment of nonlinearities calls for wrapping each leaky ReLU (or any other commonly used non-linearity) between $m \times$ upsampling and $m \times$

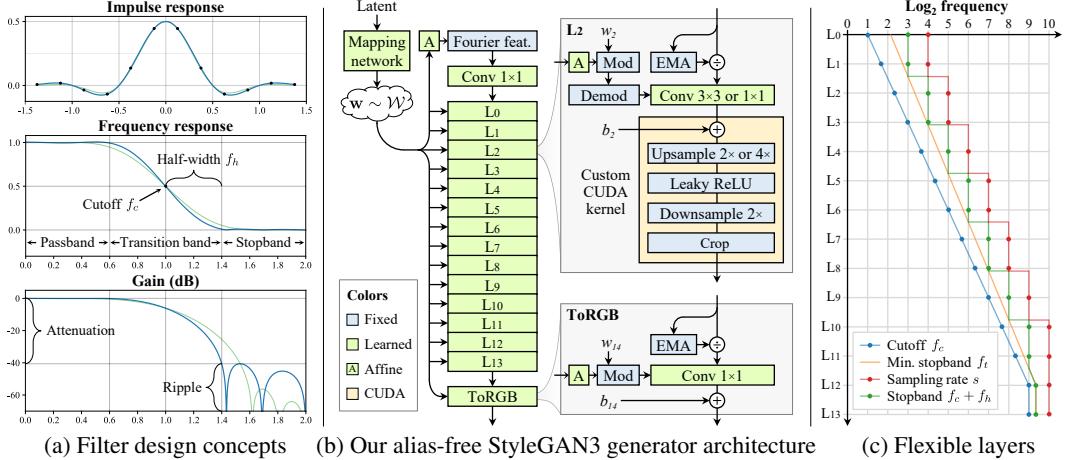


Figure 4: (a) 1D example of a $2\times$ upsampling filter with $n = 6$, $s = 2$, $f_c = 1$, and $f_h = 0.4$ (blue). Setting $f_h = 0.6$ makes the transition band wider (green), which reduces the unwanted stopband ripple and thus leads to stronger attenuation. (b) Our alias-free generator, corresponding to configs T and R in Figure 3. The main datapath consists of Fourier features and normalization (Section 3.1), modulated convolutions [34], and filtered nonlinearities (Section 3.2). (c) Flexible layer specifications (config T) with $N = 14$ and $s_N = 1024$. Cutoff f_c (blue) and minimum acceptable stopband frequency f_t (orange) obey geometric progression over the layers; sampling rate s (red) and actual stopband $f_c + f_h$ (green) are computed according to our design constraints.

downsampling, for some magnification factor m . We further note that the order of upsampling and convolution can be switched by virtue of the signal being bandlimited, allowing us to fuse the regular $2\times$ upsampling and a subsequent $m\times$ upsampling related to the nonlinearity into a single $2m\times$ upsampling. In practice, we find $m = 2$ to be sufficient (Figure 3, right), again improving EQ-T (config F). Implementing the upsample-LReLU-downsample sequence is not efficient using the primitives available in current deep learning frameworks [1, 45], and thus we implement a custom CUDA kernel (Appendix D) that combines these operations (Figure 4b), leading to $10\times$ faster training and considerable memory savings.

Non-critical sampling (config G) The critical sampling scheme — where filter cutoff is set exactly at the bandlimit — is ideal for many image processing applications as it strikes a good balance between antialiasing and the retention of high-frequency detail [58]. However, our goals are markedly different because aliasing is highly detrimental for the equivariance of the generator. While high-frequency detail is important in the output image and thus in the highest-resolution layers, it is less important in the earlier ones given that their exact resolutions are somewhat arbitrary to begin with.

To suppress aliasing, we can simply lower the cutoff frequency to $f_c = s/2 - f_h$, which ensures that all alias frequencies (above $s/2$) are in the stopband.³ For example, lowering the cutoff of the blue filter in Figure 4a would move its frequency response left so that the worst-case attenuation of alias frequencies improves from 6 dB to 40 dB. This *oversampling* can be seen as a computational cost of better antialiasing, as we now use the same number of samples to express a slower-varying signal than before. In practice, we choose to lower f_c on all layers except the highest-resolution ones, because in the end the generator must be able to produce crisp images to match the training data. As the signals now contain less spatial information, we modify the heuristic used for determining the number of feature maps to be inversely proportional to f_c instead of the sampling rate s . These changes (config G) further improve translation equivariance and push FID below the original StyleGAN2.

Transformed Fourier features (config H) Equivariant generator layers are well suited for modeling unaligned and arbitrarily oriented datasets, because any geometric transformation introduced to the intermediate features z_i will directly carry over to the final image z_N . Due to the limited capability of the layers themselves to introduce global transformations, however, the input features z_0 play a crucial role in defining the global orientation of z_N . To let the orientation vary on a per-image basis,

³Here, f_c and f_h correspond to the output (downsampling) filter of each layer. The input (upsampling) filters are based on the properties of the incoming signal, i.e., the output filter parameters of the previous layer.

the generator should have the ability to transform z_0 based on w . This motivates us to introduce a learned affine layer that outputs global translation and rotation parameters for the input Fourier features (Figure 4b and Appendix F). The layer is initialized to perform an identity transformation, but learns to use the mechanism over time when beneficial; in config H this improves the FID slightly.

Flexible layer specifications (config T) Our changes have improved the equivariance quality considerably, but some visible artifacts still remain as our accompanying videos demonstrate. On closer inspection, it turns out that the attenuation of our filters (as defined for config G) is still insufficient for the lowest-resolution layers. These layers tend to have rich frequency content near their bandlimit, which calls for extremely strong attenuation to completely eliminate aliasing.

So far, we have used the rigid sampling rate progression from StyleGAN2, coupled with simplistic choices for filter cutoff f_c and half-width f_h , but this need not be the case; we are free to specialize these parameters on a per-layer basis. In particular, we would like f_h to be high in the lowest-resolution layers to maximize attenuation in the stopband, but low in the highest-resolution layers to allow matching high-frequency details of the training data.

Figure 4c illustrates an example progression of filter parameters in a 14-layer generator with two critically sampled full-resolution layers at the end. The cutoff frequency grows geometrically from $f_c = 2$ in the first layer to $f_c = s_N/2$ in the first critically sampled layer. We choose the minimum acceptable stopband frequency to start at $f_{t,0} = 2^{2.1}$, and it grows geometrically but slower than the cutoff frequency. In our tests, the stopband target at the last layer is $f_t = f_c \cdot 2^{0.3}$, but the progression is halted at the first critically sampled layer. Next, we set the sampling rate s for each layer so that it accommodates frequencies up to f_t , rounding up to the next power of two without exceeding the output resolution. Finally, to maximize the attenuation of aliasing frequencies, we set the transition band half-width to $f_h = \max(s/2, f_t) - f_c$, i.e., making it as wide as possible within the limits of the sampling rate, but at least wide enough to reach f_t . The resulting improvement depends on how much slack is left between f_t and $s/2$; as an extreme example, the first layer stopband attenuation improves from 42 dB to 480 dB using this scheme.

The new layer specifications again improve translation equivariance (config T), eliminating the remaining artifacts. A further ablation (Figure 3, right) shows that $f_{t,0}$ provides an effective way to trade training speed for equivariance quality. Note that the number of layers is now a free parameter that does not directly depend on the output resolution. In fact, we have found that a fixed choice of N works consistently across multiple output resolutions and makes other hyperparameters such as learning rate behave more predictably. We use $N = 14$ in the remainder of this paper.

Rotation equivariance (config R) We obtain a rotation equivariant version of the network with two changes. First, we replace the 3×3 convolutions with 1×1 on all layers and compensate for the reduced capacity by doubling the number of feature maps. Only the upsampling and downsampling operations spread information between pixels in this config. Second, we replace the sinc-based downsampling filter with a radially symmetric jinc-based one that we construct using the same Kaiser scheme (Appendix C). We do this for all layers except the two critically sampled ones, where it is important to match the potentially non-radial spectrum of the training data. These changes (config R) improve EQ-R without harming FID, even though each layer has 56% fewer trainable parameters.

We also employ an additional stabilization trick in this configuration. Early on in the training, we blur all images the discriminator sees using a Gaussian filter. We start with $\sigma = 10$ pixels, which we ramp to zero over the first 200k images. This prevents the discriminator from focusing too heavily on high frequencies early on. Without this trick, config R is prone to early collapses because the generator sometimes learns to produce high frequencies with a small delay, trivializing the discriminator’s task.

4 Results

Figure 5 gives results for six datasets using StyleGAN2 [34] as well as our alias-free StyleGAN3-T and StyleGAN3-R generators. In addition to the standard FFHQ [33] and METFACES [32], we created unaligned versions of them. We also created a properly resampled version of AFHQ [14] and collected a new BEACHES dataset. Appendix B describes the datasets in detail. The results show that our FID remains competitive with StyleGAN2. StyleGAN3-T and StyleGAN3-R perform equally well in terms of FID, and both show a very high level of translation equivariance. As expected, only the latter provides rotation equivariance. In FFHQ (1024×1024) the three generators had 30.0M, 22.3M and 15.8M parameters, while the training times were 1106, 1576 (+42%) and 2248

Dataset	Config	FID ↓	EQ-T ↑	EQ-R ↑				
FFHQ-U 70000 img, 1024 ² Train from scratch	StyleGAN2	3.79	15.89	10.79				
	StyleGAN3-T (ours)	3.67	61.69	13.95				
	StyleGAN3-R (ours)	3.66	64.78	47.64				
FFHQ 70000 img, 1024 ² Train from scratch	StyleGAN2	2.70	13.58	10.22				
	StyleGAN3-T (ours)	2.79	61.21	13.82				
	StyleGAN3-R (ours)	3.07	64.76	46.62				
METFACES-U 1336 img, 1024 ² ADA, from FFHQ-U	StyleGAN2	18.98	18.77	13.19				
	StyleGAN3-T (ours)	18.75	64.11	16.63				
	StyleGAN3-R (ours)	18.75	66.34	48.57				
METFACES 1336 img, 1024 ² ADA, from FFHQ	StyleGAN2	15.22	16.39	12.89				
	StyleGAN3-T (ours)	15.11	65.23	16.82				
	StyleGAN3-R (ours)	15.33	64.86	46.81				
AFHQV2 15803 img, 512 ² ADA, from scratch	StyleGAN2	4.62	13.83	11.50				
	StyleGAN3-T (ours)	4.04	60.15	13.51				
	StyleGAN3-R (ours)	4.40	64.89	40.34				
BEACHES 20155 img, 512 ² ADA, from scratch	StyleGAN2	5.03	15.73	12.69				
	StyleGAN3-T (ours)	4.32	59.33	15.88				
	StyleGAN3-R (ours)	4.57	63.66	37.42				
Ablation		Translation eq.		+ Rotation eq.				
		FID ↓	EQ-T ↑	FID ↓	EQ-T ↑	EQ-R ↑		
* Main configuration		4.62	63.01	4.50	66.65	40.48		
With mixing reg.		4.60	63.48	4.67	63.59	40.90		
With noise inputs		4.96	24.46	5.79	26.71	26.80		
Without flexible layers		4.64	45.20	4.65	44.74	22.52		
Fixed Fourier features		5.93	64.57	6.48	66.20	41.77		
With path length reg.		5.00	68.36	5.98	71.64	42.18		
0.5× capacity		7.43	63.14	6.52	63.08	39.89		
* 1.0× capacity		4.62	63.01	4.50	66.65	40.48		
2.0× capacity		3.80	66.61	4.18	70.06	42.51		
* Kaiser filter, $n = 6$		4.62	63.01	4.50	66.65	40.48		
Lanczos filter, $a = 2$		4.69	51.93	4.44	57.70	25.25		
Gaussian filter, $\sigma = 0.4$		5.91	56.89	5.73	59.53	39.43		
G-CNN comparison		FID ↓	EQ-T ↑	EQ-R ↑	Params	Time		
* StyleGAN3-T (ours)		4.62	63.01	13.12	23.3M	1.00 ×		
+ p4 symmetry [16]		4.69	61.90	17.07	21.8M	2.48×		
* StyleGAN3-R (ours)		4.50	66.65	40.48	15.8M	1.37×		

Figure 5: **Left:** Results for six datasets. We use adaptive discriminator augmentation (ADA) [32] for the smaller datasets. “StyleGAN2” corresponds to our baseline config B with Fourier features. **Right:** Ablations and comparisons for FFHQ-U (unaligned FFHQ) at 256². * indicates our default choices.

(+103%) GPU hours. Our accompanying videos show side-by-side comparisons with StyleGAN2, demonstrating visually that the texture sticking problem has been solved. The resulting motion is much more natural, better sustaining an illusion that there is a coherent 3D scene being imaged.

Ablations and comparisons In Section 3.1 we disabled a number of StyleGAN2 features. We can now turn them on one by one to gauge their effect on our generators (Figure 5, right). While mixing regularization can be re-enabled without any ill effects, we also find that styles can be mixed quite reliably even without this explicit regularization (Appendix A). Re-enabling noise inputs or relying on StyleGAN2’s original layer specifications compromises equivariances significantly, and using fixed Fourier features or re-enabling path length regularization harms FID. Path length regularization is in principle at odds with translation equivariance, as it penalizes image changes upon latent space walk and thus encourages texture sticking. We suspect that the counterintuitive improvement in equivariance may come from slightly blurrier generated images, at a cost of poor FID.

In a scaling test we tried changing the number of feature maps, observing that equivariances remain at a high level, but FID suffers considerably when the capacity is halved. Doubling the capacity improves result quality in terms of FID, at the cost of almost 4× training time. Finally, we consider alternatives for our windowed Kaiser filter. Lanczos is competitive in terms of FID, but as a separable filter it compromises rotation equivariance in particular. Gaussian leads to clearly worse FIDs.

We compare StyleGAN3-R to an alternative where the rotation part is implemented using $p4$ symmetric G-CNN [16, 17] on top of our StyleGAN3-T. This approach provides only modest rotation equivariance while being slower to train. Steerable filters [63] could theoretically provide competitive EQ-R, but the memory and training time requirements proved infeasible with generator networks of this size.

Appendix A demonstrates that the spectral properties of generated images closely match training data, comparing favorably to several earlier architectures.

Internal representations Figure 6 visualizes typical internal representations from the networks. While in StyleGAN2 all feature maps seem to encode signal magnitudes, in our networks some of the maps take a different role and encode phase information instead. Clearly this is something that is needed when the network synthesizes detail *on the surfaces*; it needs to invent a coordinate system. In StyleGAN3-R, the emergent positional encoding patterns appear to be somewhat more well-defined. We believe that the existence of a coordinate system that allows precise localization on the surfaces of objects will prove useful in various applications, including advanced image and video editing.

5 Limitations, discussion, and future work

In this work we modified only the generator, but it seems likely that further benefits would be available by making the discriminator equivariant as well. For example, in our FFHQ results the teeth do not move correctly when the head turns, and we suspect that this is caused by the discriminator

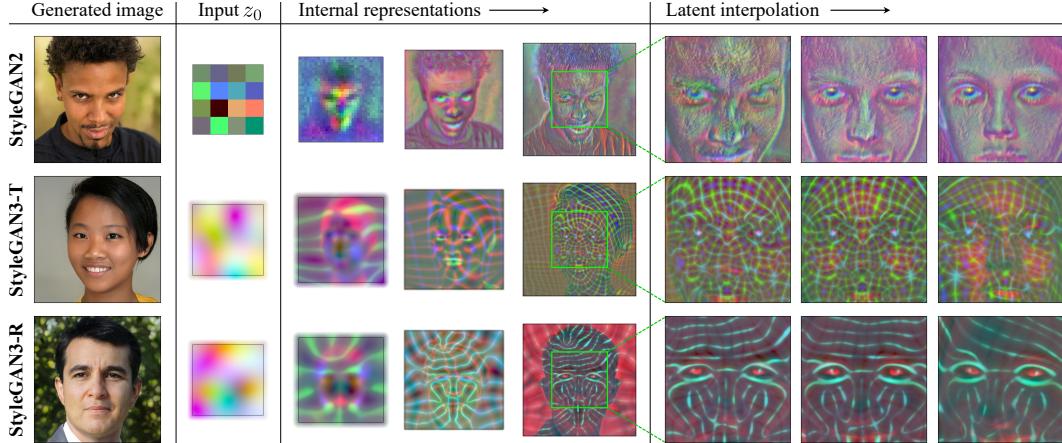


Figure 6: Example internal representations (3 feature maps as RGB) in StyleGAN2 and our generators.

accidentally preferring to see the front teeth at certain pixel locations. Concurrent work has identified that aliasing is detrimental for such generalization [59].

Our alias-free generator architecture contains implicit assumptions about the nature of the training data, and violating these may cause training difficulties. Let us consider an example. Suppose we have black-and-white cartoons as training data that we (incorrectly) pre-process using point sampling [44], leading to training images where almost all pixels are either black or white and the edges are jagged. This kind of badly aliased training data is difficult for GANs in general, but it is especially at odds with equivariance: on the one hand, we are asking the generator to be able to translate the output smoothly by subpixel amounts, but on the other hand, edges must still remain jagged and pixels only black/white, to remain faithful to the training data. The same issue can also arise with letterboxing of training images, low-quality JPEGs, or retro pixel graphics, where the jagged stair-step edges are a defining feature of the aesthetic. In such cases it may be beneficial for the generator to be aware of the pixel grid.

In future, it might be interesting to re-introduce noise inputs (stochastic variation) in a way that is consistent with hierarchical synthesis. A better path length regularization would encourage neighboring features to move together, not discourage them from moving at all. It might be beneficial to try to extend our approach to equivariance w.r.t. scaling, anisotropic scaling, or even arbitrary homeomorphisms. Finally, it is well known that antialiasing should be done before tone mapping. So far, all GANs — including ours — have operated in the sRGB color space (after tone mapping).

Attention layers in the middle of a generator [68] could likely be dealt with similarly to non-linearities by temporarily switching to higher resolution — although the time complexity of attention layers may make this somewhat challenging in practice. Recent attention-based GANs that start with a tokenizing transformer (e.g., VQGAN [18]) may be at odds with equivariance. Whether it is possible to make them equivariant is an important open question.

Potential negative societal impacts of (image-producing) GANs include many forms of disinformation, from fake portraits in social media [27] to propaganda videos of world leaders [50]. Our contribution eliminates certain characteristic artifacts from videos, potentially making them more convincing or deceiving, depending on the application. Viable solutions include model watermarking [67] along with large-scale authenticity assessment in major social media sites. This entire project consumed 92 GPU years and 225 MWh of electricity on an in-house cluster of NVIDIA V100s. The new StyleGAN3 generator is only marginally costlier to train or use than that of StyleGAN2.

6 Acknowledgments

We thank David Luebke, Ming-Yu Liu, Koki Nagano, Tuomas Kynkänniemi, and Timo Viitanen for reviewing early drafts and helpful suggestions. Frédéric Durand for early discussions. Tero Kuosmanen for maintaining our compute infrastructure. AFHQ authors for an updated version of their dataset. Getty Images for the training images in the BEACHES dataset. We did not receive external funding or additional revenues for this project.

References

- [1] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, M. Kudlur, J. Levenberg, R. Monga, S. Moore, D. G. Murray, B. Steiner, P. Tucker, V. Vasudevan, P. Warden, M. Wicke, Y. Yu, and X. Zheng. TensorFlow: A system for large-scale machine learning. In *Proc. 12th USENIX Conference on Operating Systems Design and Implementation*, OSDI'16, pages 265–283, 2016.
- [2] R. Abdal, P. Zhu, N. J. Mitra, and P. Wonka. StyleFlow: Attribute-conditioned exploration of StyleGAN-generated images using conditional continuous normalizing flows. *ACM Trans. Graph.*, 40(3), 2021.
- [3] Y. Alaluf, O. Patashnik, and D. Cohen-Or. Only a matter of style: Age transformation using a style-based regression model. *CoRR*, abs/2102.02754, 2021.
- [4] I. Anokhin, K. Demochkin, T. Khakhulin, G. Sterkin, V. Lempitsky, and D. Korzhenkov. Image generators with conditionally-independent pixel synthesis. In *Proc. CVPR*, 2021.
- [5] A. Azulay and Y. Weiss. Why do deep convolutional networks generalize so poorly to small image transformations? *Journal of Machine Learning Research*, 20(184):1–25, 2019.
- [6] D. Bau, A. Andonian, A. Cui, Y. Park, A. Jahanian, A. Oliva, and A. Torralba. Paint by word. *CoRR*, abs/2103.10951, 2021.
- [7] D. Bau, S. Liu, T. Wang, J.-Y. Zhu, and A. Torralba. Rewriting a deep generative model. In *Proc. ECCV*, 2020.
- [8] D. Bau, J. Zhu, H. Strobelt, B. Zhou, J. B. Tenenbaum, W. T. Freeman, and A. Torralba. GAN dissection: Visualizing and understanding generative adversarial networks. In *Proc. ICLR*, 2019.
- [9] R. E. Blahut. *Theory of remote image formation*. Cambridge University Press, 2004.
- [10] T. Broad, F. F. Leymarie, and M. Grierson. Network bending: Expressive manipulation of deep generative models. In *Proc. EvoMUSART*, pages 20–36, 2021.
- [11] A. Brock, J. Donahue, and K. Simonyan. Large scale GAN training for high fidelity natural image synthesis. In *Proc. ICLR*, 2019.
- [12] A. Chaman and I. Dokmanić. Truly shift-invariant convolutional neural networks. In *Proc. CVPR*, 2021.
- [13] Y. Chen, S. Liu, and X. Wang. Learning continuous image representation with local implicit image function. In *Proc. CVPR*, 2021.
- [14] Y. Choi, Y. Uh, J. Yoo, and J.-W. Ha. StarGAN v2: Diverse image synthesis for multiple domains. In *Proc. CVPR*, 2020.
- [15] M. Chu, Y. Xie, J. Mayer, L. Leal-Taixé, and N. Thuerey. Learning temporal coherence via self-supervision for GAN-based video generation. *ACM Trans. Graph.*, 39(4), 2020.
- [16] T. S. Cohen and M. Welling. Group equivariant convolutional networks. In *Proc. ICML*, 2016.
- [17] N. Dey, A. Chen, and S. Ghafurian. Group equivariant generative adversarial networks. In *Proc. ICLR*, 2021.
- [18] P. Esser, R. Rombach, and B. Ommer. Taming transformers for high-resolution image synthesis. In *Proc. CVPR*, 2021.
- [19] R. Gal, D. Cohen, A. Bermano, and D. Cohen-Or. SWAGAN: A style-based wavelet-driven generative model. *CoRR*, abs/2102.06108, 2021.
- [20] R. Ge, X. Feng, H. Pyla, K. Cameron, and W. Feng. Power measurement tutorial for the Green500 list. <https://www.top500.org/green500/resources/tutorials/>, Accessed March 1, 2020.
- [21] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial networks. In *Proc. NIPS*, 2014.
- [22] J. Gu, Y. Shen, and B. Zhou. Image processing using multi-code GAN prior. In *Proc. CVPR*, 2020.
- [23] Gwern. Making anime faces with stylegan. <https://www.gwern.net/Faces#stylegan2-ext-modifications>, Accessed June 4, 2021.
- [24] Z. Hao, A. Mallya, S. J. Belongie, and M. Liu. GANcraft: Unsupervised 3D neural rendering of minecraft worlds. *CoRR*, abs/2104.07659, 2021.
- [25] E. Härkönen, A. Hertzmann, J. Lehtinen, and S. Paris. GANSpace: Discovering interpretable GAN controls. In *Proc. NeurIPS*, 2020.
- [26] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter. GANs trained by a two time-scale update rule converge to a local Nash equilibrium. In *Proc. NIPS*, 2017.
- [27] K. Hill and J. White. Designed to deceive: Do these people look real to you? *The New York Times*, 11 2020.
- [28] M. A. Islam, S. Jia, and N. D. B. Bruce. How much position information do convolutional neural networks encode? In *Proc. ICLR*, 2020.
- [29] A. Jahanian, L. Chai, and P. Isola. On the "steerability" of generative adversarial networks. In *Proc. ICLR*, 2020.
- [30] J. F. Kaiser. Nonrecursive digital filter design using the I_0 -sinh window function. In *Proc. 1974 IEEE International Symposium on Circuits & Systems*, pages 20–23, 1974.
- [31] T. Karras, T. Aila, S. Laine, and J. Lehtinen. Progressive growing of GANs for improved quality, stability, and variation. In *Proc. ICLR*, 2018.

- [32] T. Karras, M. Aittala, J. Hellsten, S. Laine, J. Lehtinen, and T. Aila. Training generative adversarial networks with limited data. In *Proc. NeurIPS*, 2020.
- [33] T. Karras, S. Laine, and T. Aila. A style-based generator architecture for generative adversarial networks. In *Proc. CVPR*, 2018.
- [34] T. Karras, S. Laine, M. Aittala, J. Hellsten, J. Lehtinen, and T. Aila. Analyzing and improving the image quality of StyleGAN. In *Proc. CVPR*, 2020.
- [35] O. S. Kayhan and J. C. van Gemert. On translation invariance in CNNs: Convolutional layers can exploit absolute spatial location. In *Proc. CVPR*, 2020.
- [36] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *Proc. ICLR*, 2015.
- [37] M. Liu, T. Breuel, and J. Kautz. Unsupervised image-to-image translation networks. In *Proc. NIPS*, 2017.
- [38] M. Manfredi and Y. Wang. Shift equivariance in object detection. In *Proc. ECCV 2020 Workshops*, 2020.
- [39] S. Menon, A. Damian, S. Hu, N. Ravi, and C. Rudin. PULSE: Self-supervised photo upsampling via latent space exploration of generative models. In *Proc. CVPR*, 2020.
- [40] L. Mescheder, A. Geiger, and S. Nowozin. Which training methods for GANs do actually converge? In *Proc. ICML*, 2018.
- [41] A. V. Oppenheim and R. W. Schafer. *Discrete-Time Signal Processing*. Prentice Hall Press, USA, 3rd edition, 2009.
- [42] T. Park, M. Liu, T. Wang, and J. Zhu. Semantic image synthesis with spatially-adaptive normalization. In *Proc. CVPR*, 2019.
- [43] T. Park, J.-Y. Zhu, O. Wang, J. Lu, E. Shechtman, A. A. Efros, and R. Zhang. Swapping autoencoder for deep image manipulation. In *Proc. NeurIPS*, 2020.
- [44] G. Parmar, R. Zhang, and J. Zhu. On buggy resizing libraries and surprising subtleties in FID calculation. *CoRR*, abs/2104.11222, 2021.
- [45] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala. PyTorch: An imperative style, high-performance deep learning library. In *Proc. NeurIPS*, 2019.
- [46] O. Patashnik, Z. Wu, E. Shechtman, D. Cohen-Or, and D. Lischinski. StyleCLIP: Text-driven manipulation of StyleGAN imagery. *CoRR*, abs/2103.17249, 2021.
- [47] P. Ramachandran, B. Zoph, and Q. V. Le. Swish: a self-gated activation function. *CoRR*, abs/1710.05941, 2017.
- [48] A. Ramesh, M. Pavlov, G. Goh, S. Gray, C. Voss, A. Radford, M. Chen, and I. Sutskever. Zero-shot text-to-image generation. *CoRR*, abs/2102.12092, 2021.
- [49] E. Richardson, Y. Alaluf, O. Patashnik, Y. Nitzan, Y. Azar, S. Shapiro, and D. Cohen-Or. Encoding in style: A StyleGAN encoder for image-to-image translation. In *Proc. CVPR*, 2021.
- [50] M. Seymour. Canny AI: Imagine world leaders singing. *fxguide*, 4 2019.
- [51] C. E. Shannon. Communication in the presence of noise. *Proc. Institute of Radio Engineers*, 37(1):10–21, 1949.
- [52] Y. Shen and B. Zhou. Closed-form factorization of latent semantics in GANs. In *CVPR*, 2021.
- [53] V. Sitzmann, J. N. Martel, A. W. Bergman, D. B. Lindell, and G. Wetzstein. Implicit neural representations with periodic activation functions. In *Proc. NeurIPS*, 2020.
- [54] I. Skorokhodov, S. Ignatyev, and M. Elhoseiny. Adversarial generation of continuous images. In *Proc. CVPR*, 2021.
- [55] R. Suzuki, M. Koyama, T. Miyato, T. Yonetsuji, and H. Zhu. Spatially controllable image synthesis with internal representation collaging. *CoRR*, abs/1811.10153, 2019.
- [56] M. Tancik, P. P. Srinivasan, B. Mildenhall, S. Fridovich-Keil, N. Raghavan, U. Singhal, R. Ramamoorthi, J. T. Barron, and R. Ng. Fourier features let networks learn high frequency functions in low dimensional domains. In *Proc. NeurIPS*, 2020.
- [57] S. Tulyakov, M. Liu, X. Yang, and J. Kautz. MoCoGAN: Decomposing motion and content for video generation. In *Proc. CVPR*, 2018.
- [58] K. Turkowski. *Filters for Common Resampling Tasks*, pages 147–165. Academic Press Professional, Inc., USA, 1990.
- [59] C. Vasconcelos, H. Larochelle, V. Dumoulin, R. Romijnders, N. L. Roux, and R. Goroshin. Impact of aliasing on generalization in deep convolutional networks. In *ICCV*, 2021.
- [60] C. von der Malsburg. Self-organization of orientation sensitive cells in striate cortex. *Biological Cybernetics*, 14(2):85–100, 1973.
- [61] T. Wang, M. Liu, J. Zhu, A. Tao, J. Kautz, and B. Catanzaro. High-resolution image synthesis and semantic manipulation with conditional GANs. In *Proc. CVPR*, 2018.
- [62] M. Weiler and G. Cesa. General $E(2)$ -equivariant steerable CNNs. In *Proc. NeurIPS*, 2019.
- [63] M. Weiler, F. A. Hamprecht, and M. Storath. Learning steerable filters for rotation equivariant CNNs. In *Proc. CVPR*, 2018.
- [64] D. Worrall and M. Welling. Deep scale-spaces: Equivariance over scale. In *Proc. NeurIPS*, 2019.

- [65] D. E. Worrall, S. J. Garbin, D. Turmukhambetov, and G. J. Brostow. Harmonic networks: Deep translation and rotation equivariance. In *Proc. CVPR*, 2017.
- [66] R. Xu, X. Wang, K. Chen, B. Zhou, and C. C. Loy. Positional encoding as spatial inductive bias in GANs. In *Proc. CVPR*, 2021.
- [67] N. Yu, V. Skripniuk, S. Abdelnabi, and M. Fritz. Artificial fingerprinting for generative models: Rooting deepfake attribution in training data. *CoRR*, abs/2007.08457, 2021.
- [68] H. Zhang, I. Goodfellow, D. Metaxas, and A. Odena. Self-attention generative adversarial networks. In *Proc. ICML*, 2019.
- [69] R. Zhang. Making convolutional networks shift-invariant again. In *Proc. ICML*, 2019.
- [70] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proc. ICCV*, 2017.
- [71] X. Zou, F. Xiao, Z. Yu, and Y. J. Lee. Delving deeper into anti-aliasing in ConvNets. In *Proc. BMVC*, 2020.

Appendices

A Additional results

Uncurated sets of samples for StyleGAN2 (baseline config B with Fourier features) and our alias-free generators StyleGAN3-T and StyleGAN3-R are shown in Figures 7 (FFHQ-U), 8 (METFACES-U), 9 (AFHQv2), and 10 (BEACHES). Truncation trick was not used when generating the images.

StyleGAN2 and our generators yield comparable FIDs in all of these datasets. Visual inspection did not reveal anything surprising in the first three datasets, but in BEACHES our new generators seem to generate a somewhat reduced set of possible scene layouts properly. We suspect that this is related to the lack of noise inputs, which forces the generators to waste capacity for what is essentially random number generation [34]. Finding a way to reintroduce noise inputs without breaking equivariances is therefore an important avenue of future work.

The accompanying interpolation videos reveal major differences between StyleGAN2 and StyleGAN3-R. For example, in METFACES much of details such as brushstrokes or cracked paint seems to be glued to the pixel coordinates in StyleGAN2, whereas with StyleGAN3 all details move together with the depicted model. The same is evident in AFHQv2 with the fur moving credibly in StyleGAN3 interpolations, while mostly sticking to the image coordinates in StyleGAN2. In BEACHES we furthermore observe that StyleGAN2 tends to “fade in” details while retaining a mostly fixed viewing position, while StyleGAN3 creates plenty of apparent rotations and movement. The videos use hand-picked seeds to better showcase the relevant effects.

In a further test we created two example cinemagraphs that mimic small-scale head movement and facial animation in FFHQ. The geometric head motion was generated as a random latent space walk along hand-picked directions from GANSpace [25] and SeFa [52]. The changes in expression were realized by applying the “global directions” method of StyleCLIP [46], using the prompts “angry face”, “laughing face”, “kissing face”, “sad face”, “singing face”, and “surprised face”. The differences between StyleGAN2 and StyleGAN3 are again very prominent, with the former displaying jarring sticking of facial hair and skin texture, even under subtle movements.

The equivariance quality videos illustrate the practical relevance of the PSNR numbers in Figures 3 and 5 of the main paper. We observe that for EQ-T numbers over ~ 50 dB indicate high-quality results, and for EQ-R ~ 40 dB look good.

We also provide an animated version of the nonlinearity visualization in Figure 2.

In style mixing [34] two or more independently chosen latent codes are fed into different layers of the generator. Ideally all combinations would produce images that are not obviously broken, and furthermore, it would be desirable that specific layers end up controlling well-defined semantic aspects in the images. StyleGAN uses mixing regularization [34] during training to achieve these goals. We observe that mixing regularization continues to work similarly in StyleGAN3, but we also wanted to know whether it is truly necessary because the regularization is known to be detrimental for many complex and multi-modal datasets [23]. When we disable the regularization, obviously broken images remain rare, based on a visual inspection of a large number of images. The semantically meaningful controls are somewhat compromised, however, as Figure 11 shows.

Figure 12 compares the convergence of our main configurations (config T and R) against the results of Karras et al. [34, 32]. The overall shape of the curves is similar; introducing translation and rotation equivariance in the generator does not appear to significantly alter the training dynamics.

Following recent works that address signal processing issues in GANs [4, 19], we show average power spectra of the generated and real images in Figure 13. The plots are computed from images that are whitened with the overall training dataset mean and standard deviation. Because FFT interprets the signal as periodic, we eliminate the sharp step edge across the image borders by windowing the pixel values prior to the transform. This eliminates the axis-aligned cross artifact which may obscure meaningful detail in the spectrum. We display the average 2D spectrum as a contour plot, which makes the orientation-dependent falloff apparent, and highlights detail like regularly spaced residuals of upsampling grids, and fixed noise patterns. We also plot 1D slices of the spectrum along the horizontal and diagonal angle without azimuthal integration, so as to not average out the detail. The code for reproducing these steps is included in the public release.

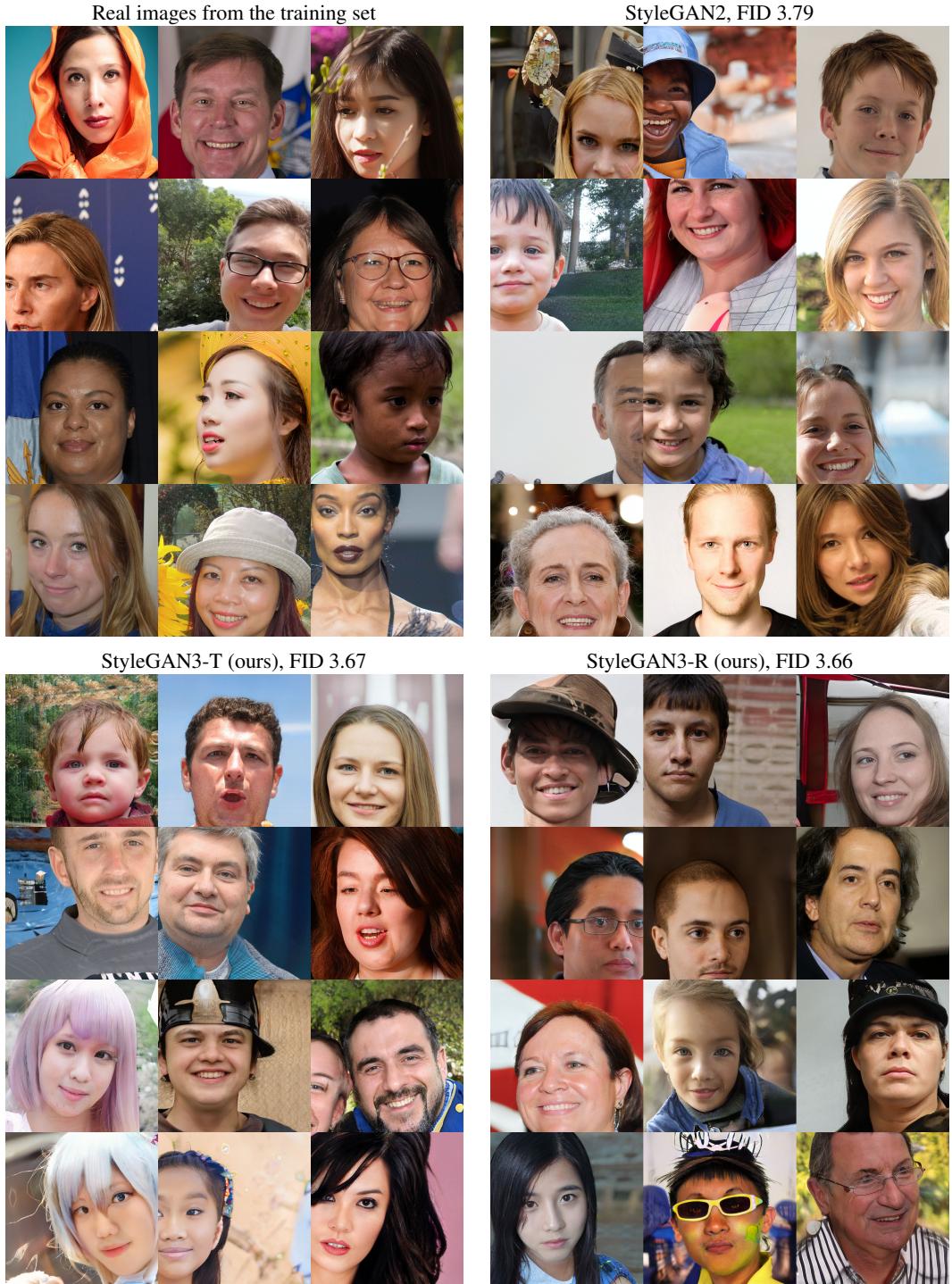


Figure 7: Uncurated samples for unaligned FFHQ (FFHQ-U). Truncation was not used.

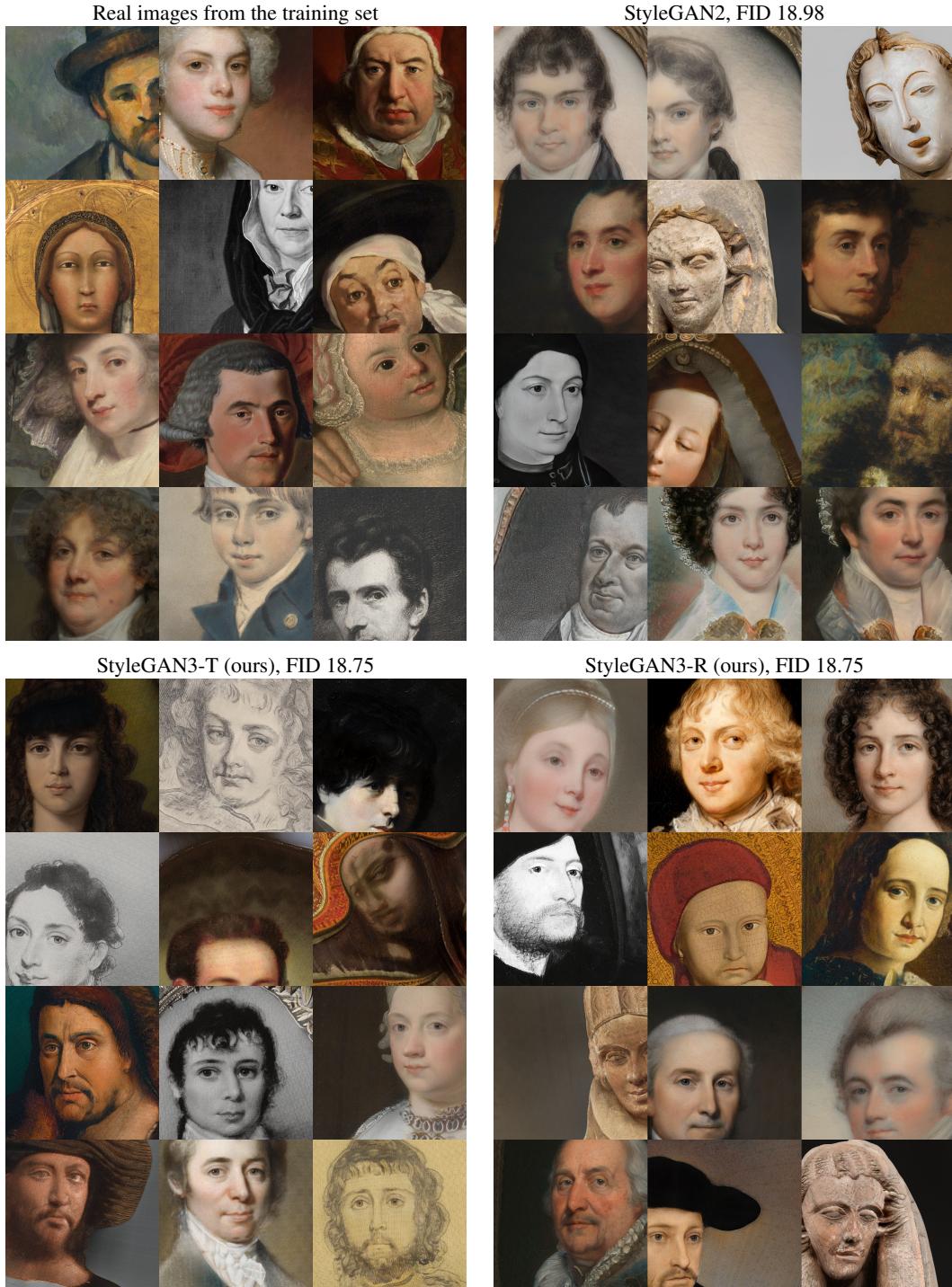


Figure 8: Uncurated samples for unaligned METFACES (METFACES-U). Truncation was not used.



Figure 9: Uncurated samples for AFHQv2. Truncation was not used.

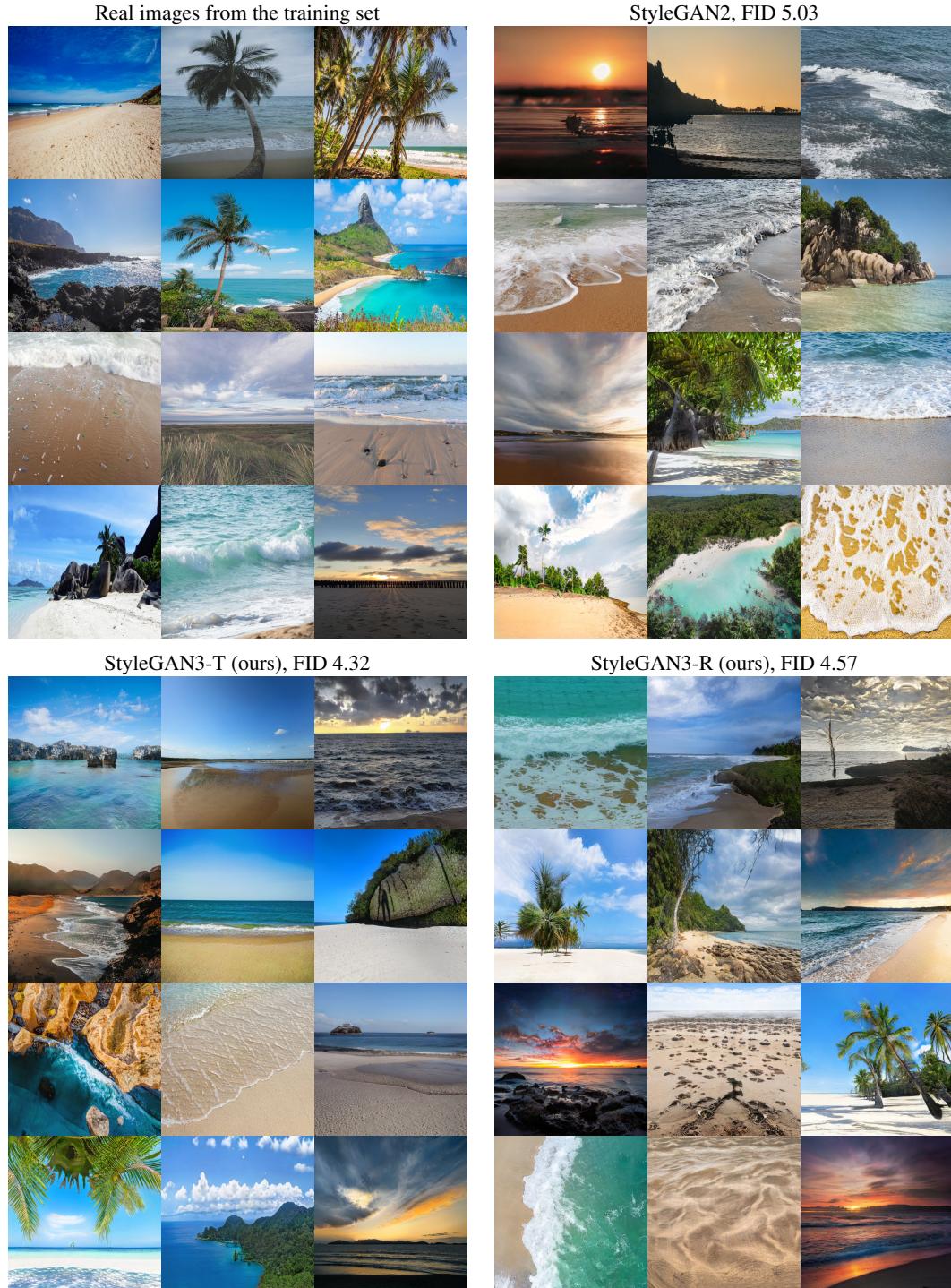




Figure 11: Hand-picked style mixing examples where the coarse (0–6) and fine (7–14) layers use a different latent code. Mixing regularization was not used during training. Head pose, coarse facial shape, hair length and glasses seem to get inherited from the coarse layers (top row), while coloring and finer facial features are mostly inherited from the fine layers (leftmost column). The control is not quite perfect: e.g., feminine/masculine features are not reliably copied from exactly one of the sources. Moving the fine/coarse boundary fixes this particular issue, but other similar problems persist.

B Datasets

In this section, we describe the new datasets and list the licenses of all datasets.

B.1 FFHQ-U and MetFaces-U

We built unaligned variants of the existing FFHQ [33] and METFACES [32] datasets. The originals are available at <https://github.com/NVlabs/ffhq-dataset> and <https://github.com/NVlabs/metfaces-dataset>, respectively. The datasets were rebuilt with a modification of the original procedure based on the original code, raw uncropped images, and facial landmark metadata. The code required to reproduce the modified datasets is included in the public release.

We use axis-aligned crop rectangles, and do not rotate them to match the orientation of the face. This retains the natural variation of camera and head tilt angles. Note that the images are still generally

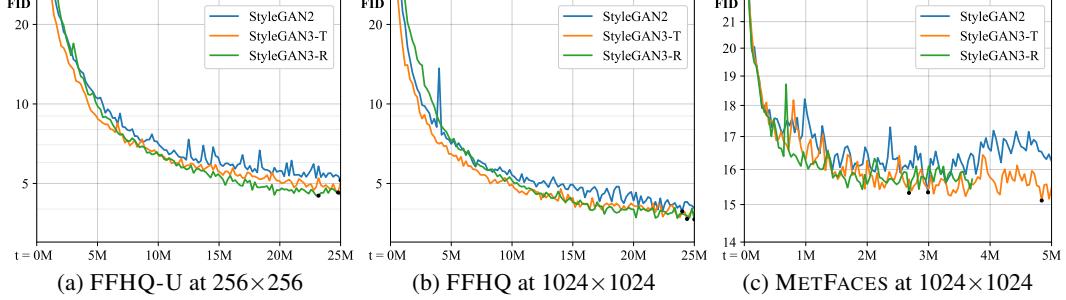


Figure 12: Training convergence with three datasets using StyleGAN2 and our main configurations (config T and R). x -axis corresponds to the total number of real images shown to the discriminator and y -axis is the Fréchet inception distance (FID), computed between 50k generated images and all training images [26, 32]; lower is better. The black dots indicate the best FID for each training run, matching the corresponding cases in Figures 3 and 5. METFACES was trained using adaptive discriminator augmentation (ADA) [32], starting from the corresponding FFHQ snapshot with the lowest FID.

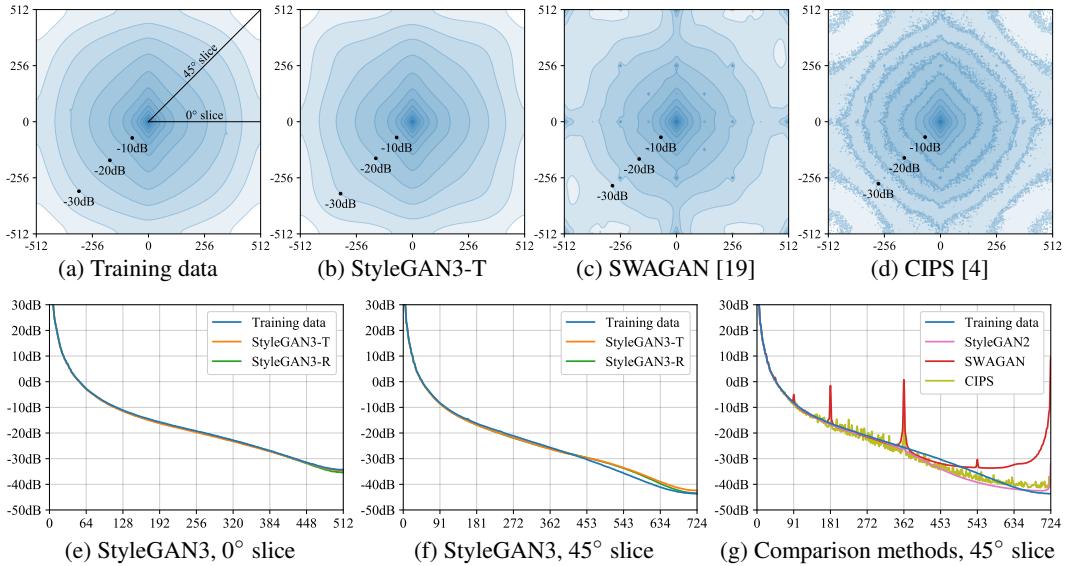


Figure 13: **Top:** Average 2D power spectrum of the training images in FFHQ at 1024×1024 resolution, along with the corresponding spectra of random images generated using StyleGAN3-T, SWAGAN [19], and CIPS [4]. Each plot represents the average power over 70k images, computed as follows. From each image, we subtract the training dataset mean, after which we divide it by the training dataset standard deviation. Note that these normalizing quantities represent the entire dataset reduced to two scalars, and do not vary by color channel or pixel coordinate. The image is then multiplied with a separable Kaiser window with $\beta = 8$, and its power spectrum is computed as the absolute values of the FFT raised to the second power. This processing is applied to each color channel separately and the result is averaged over them. These spectra are then averaged over all the images in the dataset. The result is plotted on the decibel scale. **Bottom:** One-dimensional slices of the power spectra at 0° and 45° angles.

upright, i.e., never upside down or at 90° angle. The scale of the rectangle is determined as before. For each image, the crop rectangle is randomly shifted from its original face-centered position, with the horizontal and vertical offset independently drawn from a normal distribution. The standard deviation is chosen as 20% of the crop rectangle dimension. If the crop rectangle falls partially outside the original image boundaries, we keep drawing new random offsets until we find one that does not. This removes the need to pad the images with fictional mirrored content, and we explicitly disabled this feature of the original build script.

Aside from the exact image content, the number of images and other specifications match the original dataset exactly. While FFHQ-U contains identifiable images of persons, it does not introduce new images beyond those already in the original FFHQ.

B.2 AFHQv2

We used an updated version of the AFHQ dataset [14] where the resampling filtering has been improved. The original dataset suffers from pixel-level artifacts caused by inadequate downsampling filters [44]. This caused convergence problems with our models, as the sharp “stair-step” aliasing artifacts are difficult to reproduce without direct access to the pixel grid.

The dataset was rebuilt using the original uncropped images and crop rectangle metadata, using the PIL library implementation of Lanczos resampling as recommended by Parmar et al. [44]. In a minority of cases, the crop rectangles were modified to remove non-isotropic scaling and other unnecessary transformations. A small amount ($\sim 2\%$) of images were dropped for technical reasons, leaving a total of 15803 images. Aside from this, the specifications of the dataset match the original. We use all images of all the three classes (cats, dogs, and wild animals) as one training dataset.

B.3 Beaches

BEACHES is a new dataset of 20155 photographs of beaches at resolution 512×512 . The training images were provided by Getty Images. BEACHES is a proprietary dataset that we are licensed to use, but not to redistribute. We are therefore unable to release the full training data or pre-trained models for this dataset.

B.4 Licenses

The FFHQ dataset is available under Creative Commons BY-NC-SA 4.0 license by NVIDIA Corporation, and consist of images published by respective authors under Creative Commons BY 2.0, Creative Commons BY-NC 2.0, Public Domain Mark 1.0, Public Domain CC0 1.0, and U.S. Government Works license.

The METFACES dataset is available under Creative Commons BY-NC 2.0 license by NVIDIA Corporation, and consists of images available under the Creative Commons Zero (CC0) license by the Metropolitan Museum of Art.

The original AFHQ dataset is available at <https://github.com/clovaai/stargan-v2> under Creative Commons BY-NC 4.0 license by NAVER Corporation.

C Filter details

In this section, we review basic FIR filter design methodology and detail the recipe used to construct the upsampling and downsampling filters in our generator. We start with simple Kaiser filters in one dimension, discussing parameter selection and the necessary modifications needed for upsampling and downsampling. We then proceed to extend the filters to two dimensions and conclude by detailing the alternative filters evaluated in Figure 5, right. Our definitions are consistent with standard signal processing literature (e.g., Oppenheim [41]) as well as widely used software packages (e.g., `scipy.signal.firwin`).

C.1 Kaiser low-pass filters

In one dimension, the ideal continuous-time low-pass filter with cutoff f_c is given by $\psi(x) = 2f_c \cdot \text{sinc}(2f_c x)$, where $\text{sinc}(x) = \sin(\pi x)/(\pi x)$. The ideal filter has infinite attenuation in the stopband, i.e., it completely eliminates all frequencies above f_c . However, its impulse response is also infinite, which makes it impractical for three reasons: implementation efficiency, border artifacts, and *ringing* caused by long-distance interactions. The most common way to overcome these issues is to limit the spatial extent of the filter using the window method [41]:

$$h_K(x) = 2f_c \cdot \text{sinc}(2f_c x) \cdot w_K(x), \quad (4)$$

where $w_K(x)$ is a *window function* and $h_K(x)$ is the resulting practical approximation of $\psi(x)$. Different window functions represent different tradeoffs between the frequency response and spatial extent; the smaller the spatial extent, the weaker the attenuation. In this paper we use the Kaiser window [30], also known as the Kaiser–Bessel window, that provides explicit control over this tradeoff. The Kaiser window is defined as

$$w_K(x) = \begin{cases} I_0\left(\beta\sqrt{1 - (2x/L)^2}\right)/I_0(\beta), & \text{if } |x| \leq L/2, \\ 0, & \text{if } |x| > L/2, \end{cases} \quad (5)$$

where L is the desired spatial extent, β is a free parameter that controls the shape of the window, and I_0 is the zeroth-order modified Bessel function of the first kind. Note that the window has discontinuities at $\pm L/2$; the value is strictly positive at $x = L/2$ but zero at $x = L/2 + \epsilon$.

When operating on discretely sampled signals, it is necessary to discretize the filter as well:

$$h_K[i] = h_K\left((i - (n-1)/2)/s\right)/s, \quad \text{for } i \in \{0, 1, \dots, n-1\}, \quad (6)$$

where $h_K[i]$ is the discretized version of $h_K(x)$ and s is the sampling rate. The filter is defined at n discrete spatial locations, i.e., *taps*, located $1/s$ units apart and placed symmetrically around zero. Given the values of n and s , the spatial extent can be expressed as $L = (n-1)/s$. An odd value of n results in a *zero-phase* filter that preserves the original sample locations, whereas an even value shifts the sample locations by $1/(2s)$ units.

The filters considered in this paper are approximately normalized by construction, i.e., $\int_x h_K(x) \approx \sum_i h_K[i] \approx 1$. Nevertheless, we have found it beneficial to explicitly normalize them after discretization. In other words, we strictly enforce $\sum_i h_K[i] = 1$ by scaling the filter taps to reduce the risk of introducing cumulative scaling errors when the signal is passed through several consecutive layers.

C.2 Selecting window parameters

Kaiser [30] provides convenient empirical formulas to connect the parameters of w_K to the properties of h_K . Given the number of taps and the desired transition band width, the maximum attenuation achievable with $h_K[i]$ is approximated by

$$A = 2.285 \cdot (n-1) \cdot \pi \cdot \Delta f + 7.95, \quad (7)$$

where A is the attenuation measured in decibels and Δf is the width of the transition band expressed as a fraction of $s/2$. We choose to define the transition band using half-width f_h , which gives $\Delta f = (2f_h)/(s/2)$. Given the value of A , the optimal choice for the shape parameter β is then approximated [30] by

$$\beta = \begin{cases} 0.1102 \cdot (A - 8.7), & \text{if } A > 50, \\ 0.5842 \cdot (A - 21)^{0.4} + 0.07886 \cdot (A - 21), & \text{if } 21 \leq A \leq 50, \\ 0, & \text{if } A < 21, \end{cases} \quad (8)$$

This leaves us with two free parameters: n controls the spatial extent while f_h controls the transition band. The choice of these parameters directly influences the resulting attenuation; increasing either parameter yields a higher value for A .

C.3 Upsampling and downsampling

When upsampling a signal, i.e., $\mathbf{F}_{\text{up}}(Z) = \text{III}_{s'} \odot (\phi_s * Z) = 1/s^2 \cdot \text{III}_{s'} \odot (\psi_s * Z)$, we are concerned not only with the input sampling rate s , but also with the output sampling rate s' . With an integer upsampling factor m , we can think of the upsampling operation as consisting of two steps: we first increase the sampling rate to $s' = s \cdot m$ by interleaving $m-1$ zeros between each input sample and then low-pass filter the resulting signal to eliminate the alias frequencies above $s/2$. In order to keep the signal magnitude unchanged, we must also scale the result by m with one-dimensional signals, or by m^2 with two-dimensional signals. Since the filter now operates under s' instead of s , we must adjust its parameters accordingly:

$$n' = n \cdot m, \quad L' = (n'-1)/s', \quad \Delta f' = (2f_h)/(s'/2), \quad (9)$$

which gives us the final upsampling filter

$$h'_K[i] = h'_K\left((i - (n' - 1)/2)/s'\right)/s', \quad \text{for } i \in \{0, 1, \dots, n' - 1\}. \quad (10)$$

Multiplying the number of taps by m keeps the spatial extent of the filter unchanged with respect to the input samples, and it also compensates for the reduced attenuation from $\Delta f' < \Delta f$. Note that if the upsampling factor is even, n' will be even as well, meaning that h'_K shifts the sample locations by $1/(2s')$. This is the desired behavior—if we consider sample i to represent the continuous interval $[i \cdot s, (i + 1) \cdot s]$ in the input signal, the same interval will be represented by m consecutive samples $m \cdot i, \dots, m \cdot i + m - 1$ in the output signal. Using a zero-phase upsampling filter, i.e., an odd value for n' , would break this symmetry, leading to inconsistent behavior with respect to the boundaries. Note that our symmetric interpretation is common in many computer graphics APIs, such as OpenGL, and it is also reflected in our definition of the Dirac comb III in Section 2.

Upsampling and downsampling are adjoint operations with respect to each other, disregarding the scaling of the signal magnitude. This means that the above definitions are readily applicable to downsampling as well; to downsample a signal by factor m , we first filter it by h'_K and then discard the last $m - 1$ samples within each group of m consecutive samples. The interpretation of all filter parameters, as well as the sample locations, is analogous to the upsampling case.

C.4 Two-dimensional filters

Any one-dimensional filter, including h_K , can be trivially extended to two dimensions by defining the corresponding separable filter

$$h_K^+(\mathbf{x}) = h_K(x_0) \cdot h_K(x_1) = (2f_c)^2 \cdot \text{sinc}(2f_c x_0) \cdot \text{sinc}(2f_c x_1) \cdot w_K(x_0) \cdot w_K(x_1), \quad (11)$$

where $\mathbf{x} = (x_0, x_1)$. h_K^+ has the same cutoff as h_K along the coordinate axes, i.e., $\mathbf{f}_{c,x} = (f_c, 0)$ and $\mathbf{f}_{c,y} = (0, f_c)$, and its frequency response forms a square shape over the 2D plane, implying that the cutoff frequency along the diagonal is $\mathbf{f}_{c,d} = (f_c, f_c)$. In practice, a separable filter can be implemented efficiently by first filtering each row of the two-dimensional signal independently with h_K and then doing the same for each column. This makes h_K^+ an ideal choice for all upsampling filters in our generator, as well as the downsampling filters in configs A–T (Figure 3, left).

The fact that the spectrum of h_K^+ is not radially symmetric, i.e., $\|\mathbf{f}_{c,d}\| \neq \|\mathbf{f}_{c,x}\|$, is problematic considering config R. If we rotate the input feature maps of a given layer, their frequency content will rotate as well. To enforce rotation equivariant behavior, we must ensure that the effective cutoff frequencies remain unchanged by this. The ideal radially symmetric low-pass filter [9] is given by $\psi_s^\circ(\mathbf{x}) = (2f_c)^2 \cdot \text{jinc}(2f_c \|\mathbf{x}\|)$. The jinc function, also known as besinc, sombrero function, or Airy disk, is defined as $\text{jinc}(x) = 2J_1(\pi x)/(\pi x)$, where J_1 is the first order Bessel function of the first kind. Using the same windowing scheme as before, we define the corresponding practical filter as

$$h_K^\circ(\mathbf{x}) = (2f_c)^2 \cdot \text{jinc}(2f_c \|\mathbf{x}\|) \cdot w_K(x_0) \cdot w_K(x_1). \quad (12)$$

Note that even though jinc is radially symmetric, we still treat the window function as separable in order to retain its spectral properties. In config R, we perform all downsampling operations using h_K° , except for the last two critically sampled layers where we revert to h_K^+ .

C.5 Alternative filters

In Figure 5, right, we compare the effectiveness of Kaiser filters against two alternatives: Lanczos and Gaussian. These filters are typically defined using prototypical filter kernels k_L and k_G , respectively:

$$k_L(x) = \begin{cases} \text{sinc}(x) \cdot \text{sinc}(x/a), & \text{if } |x| < a, \\ 0, & \text{if } |x| \geq a, \end{cases} \quad (13)$$

$$k_G(x) = \exp\left(-\frac{1}{2}(x/\sigma)^2\right) / (\sigma\sqrt{2\pi}), \quad (14)$$

where a is the spatial extent of the Lanczos kernel, typically set to 2 or 3, and σ is the standard deviation of the Gaussian kernel. In Figure 5 of the main paper we set $a = 2$ and $\sigma = 0.4$; we tested several different values and found these choices to work reasonably well.

The main shortcoming of the prototypical kernels is that they do not provide an explicit way to control the cutoff frequency. In order to enable apples-to-apples comparison, we assume that the kernels have an implicit cutoff frequency at 0.5 and scale their impulse responses to account for the varying f_c :

$$h_L(x) = 2f_c \cdot k_L(2f_c x), \quad h_G(x) = 2f_c \cdot k_G(2f_c x). \quad (15)$$

We limit the computational complexity of the Gaussian filter by enforcing $h_G(x) = 0$ when $|x| > 8/s$, with respect to the input sampling rate in the upsampling case. In practice, $h_G(x)$ is already very close to zero in this range, so the effect of this approximation is negligible. Finally, we extend the filters to two dimensions by defining the corresponding separable filters:

$$h_L^+(x) = (2f_c)^2 \cdot k_L(2f_c x_0) \cdot k_L(2f_c x_1), \quad h_G^+(x) = (2f_c)^2 \cdot k_G(2f_c x_0) \cdot k_G(2f_c x_1). \quad (16)$$

Note that h_G^+ is radially symmetric by construction, which makes it ideal for rotation equivariance. h_L^+ , however, has no widely accepted radially symmetric counterpart, so we simply use the same separable filter in config R as well.

D Custom CUDA kernel for filtered nonlinearity

Implementing the upsample-nonlinearity-downsample sequence is inefficient using the standard primitives available in modern deep learning frameworks. The intermediate feature maps have to be transferred between on-chip and off-chip GPU memory multiple times and retained for the backward pass. This is especially costly because the intermediate steps operate on upsampled, high-resolution data. To overcome this, we implement the entire sequence as a single operation using a custom CUDA kernel. This improves training performance by approximately an order of magnitude thanks to reduced memory traffic, and also decreases GPU memory usage significantly.

The combined kernel consists of four phases: input, upsampling, nonlinearity, and downsampling. The computation is parallelized by subdividing the output feature maps into non-overlapping tiles, and computing one output tile per CUDA thread block. First, in input phase, the corresponding input region is read into on-chip shared memory of the thread block. Note that the input regions for neighboring output tiles will overlap spatially due to the spatial extent of filters.

The execution of up-/downsampling phases depends on whether the corresponding filters are separable or not. For a separable filter, we perform vertical and horizontal 1D convolutions sequentially, whereas a non-separable filter requires a single 2D convolution. All these convolutions and the nonlinearity operate in on-chip shared memory, and only the final output of the downsampling phase is written to off-chip GPU memory.

D.1 Gradient computation

To compute gradients of the combined operation, they need to propagate through each of the phases in reverse order. Fortunately, the combined upsample-nonlinearity-downsample operation is mostly self-adjoint with proper changes in parameters, e.g., swapping the up-/downsampling factors and the associated filters. The only problematic part is the nonlinearity that is performed in the upsampled resolution. A naive but general solution would be to store the intermediate high-resolution input to the nonlinearity, but the memory consumption would be infeasible for training large models.

Our kernel is specialized to use leaky ReLU as the nonlinearity, which offers a straightforward way to conserve memory: to propagate gradients, it is sufficient to know whether the corresponding input value to nonlinearity was positive or negative. When using 16-bit floating-point datatypes, there is an additional complication because the outputs of the nonlinearity need to be clamped [32], and when this occurs, the corresponding gradients must be zero. Therefore, in the forward pass we store two bits of auxiliary information per value to cover the three possible cases: positive, negative, or clamped. In the backward pass, reading these bits is sufficient for correct gradient computation — no other information from the forward pass is needed.

D.2 Optimizations for common upsampling factors

Let us consider one-dimensional $2 \times$ upsampling where the input is (virtually) interleaved with zeros and convolved with an n' -tap filter where $n' = 2n$ (cf. Equation 9). There are n nonzero input values

	upsample 2×				upsample 4×				upsample 2×			
Sep. up	yes	yes	no	no	yes	yes	no	no	yes	yes	no	no
Sep. down	yes	no	yes	no	yes	no	yes	no	yes	no	yes	no
PyTorch (ms)	7.88	12.40	12.68	17.12	10.07	31.51	14.96	36.33	39.35	56.73	125.83	143.15
Ours (ms)	0.42	0.59	0.66	0.92	0.49	0.84	0.80	1.01	1.20	1.89	3.04	3.66
Speedup ×	19	21	19	19	21	38	19	36	33	30	41	39

Figure 14: Upsample-nonlinearity-downsample timings in milliseconds using native PyTorch operations vs our optimized CUDA kernel. The benchmarks were run on NVIDIA Titan V GPU, using input size $512 \times 512 \times 32$ and filter size $n = 6$, i.e., $n' = 12$ and $n' = 24$ for up-/downsampling rates of 2 and 4, respectively. Sep. up and Sep. down indicate the use of separable up-/downsampling filters.

under the n' -tap kernel, so if each output pixel is computed separately, the convolution requires n multiply-add operations per pixel and equally many shared memory load instructions, for a total of $2n$ instructions per output pixel.⁴ However, note that the computation of two neighboring output pixels accesses only $n + 1$ input pixels in total. By computing two output pixels at a time and avoiding redundant shared memory load instructions, we obtain an average cost of $\frac{3}{2}n + \frac{1}{2}$ instructions per pixel—close to 25% savings. For $4 \times$ upsampling, we can similarly reduce the instruction count by up to 37.5% by computing four output pixels at a time. We apply these optimizations in $2 \times$ and $4 \times$ upsampling for both separable and non-separable filters.

Figure 14 benchmarks the performance of our kernel with various up-/downsampling factors and with separable and non-separable filters. In network layers that keep the sampling rate fixed, both factors are $2 \times$, whereas layers that increase the sampling rate by a factor of two, $4 \times$ upsampling is combined with $2 \times$ downsampling. The remaining combination of $2 \times$ upsampling and $4 \times$ downsampling is needed when computing gradients of the latter case. The speedup over native PyTorch operations varies between $\sim 20\text{--}40\times$, which yields an overall training speedup of approximately $10\times$.

E Equivariance metrics

In this section, we describe our equivariance metrics, EQ-T and EQ-R, in detail. We also present additional results using an alternative translation metric, EQ-T_{frac}, based on fractional sub-pixel translation.

We express each of our metrics as the *peak signal-to-noise ratio* (PSNR) between two sets of images, measured in decibels (dB). PSNR is a commonly used metric in image restoration literature. In the typical setting we have two signals, reference I and its noisy approximation K , defined over discrete domain \mathcal{D} —usually a two-dimensional pixel grid. The PSNR between I and K is then defined via the mean squared error (MSE):

$$\text{MSE}_{\mathcal{D}}(I, K) = \frac{1}{\|\mathcal{D}\|} \sum_{i \in \mathcal{D}} (I[i] - K[i])^2, \quad (17)$$

$$\text{PSNR}_{\mathcal{D}}(I, K) = 10 \cdot \log_{10} \left(\frac{I_{max}^2}{\text{MSE}_{\mathcal{D}}(I, K)} \right), \quad (18)$$

where $\text{MSE}_{\mathcal{D}}(I, K)$ is the average squared difference between matching elements of I and K . I_{max} is the expected dynamic range of the reference signal, i.e., $I_{max} \approx \max_{i \in \mathcal{D}}(I[i]) - \min_{i \in \mathcal{D}}(I[i])$. The dynamic range is usually considered to be a global constant, e.g., the range of valid RGB values, as opposed to being dependent on the content of I . In our case, I and K represent desired and actual outputs of the synthesis network, respectively, with a dynamic range of $[-1, 1]$. This implies that $I_{max} = 2$. High PSNR values indicate that K is close to I ; in the extreme case, where $K = I$, we have $\text{PSNR}_{\mathcal{D}}(I, K) = \infty$ dB.

Since we are interested in *sets* of images, we use a slightly extended definition for MSE that allows I and K to be defined over an arbitrary, potentially uncountable domain:

$$\text{MSE}_{\mathcal{D}}(I, K) = \mathbb{E}_{i \sim \mathcal{D}} [(I(i) - K(i))^2]. \quad (19)$$

⁴Input of the upsampling is stored in shared memory, but the filter weights can be stored in CUDA constant memory where they can be accessed without a separate load instruction.

Configuration	FID	EQ-T	EQ-T _{frac}	Parameter	FID	EQ-T	EQ-T _{frac}
A StyleGAN2	5.14	—	—	Filter size $n = 4$	4.72	57.49	44.65
B + Fourier features	4.79	16.23	16.28	* Filter size $n = 6$	4.50	66.65	45.92
C + No noise inputs	4.54	15.81	15.84	Filter size $n = 8$	4.66	65.57	46.57
D + Simplified generator	5.21	19.47	19.57	Upsampling $m = 1$	4.38	39.96	37.55
E + Boundaries & upsampling	6.02	24.62	24.70	* Upsampling $m = 2$	4.50	66.65	45.92
F + Filtered nonlinearities	6.35	30.60	30.68	Upsampling $m = 4$	4.57	74.21	46.81
G + Non-critical sampling	4.78	43.90	42.24	Stopband $f_{t,0} = 2^{1.5}$	4.62	51.10	44.46
H + Transformed Fourier features	4.64	45.20	42.78	* Stopband $f_{t,0} = 2^{2.1}$	4.50	66.65	45.92
T + Flexible layers (StyleGAN3-T)	4.62	63.01	46.40	Stopband $f_{t,0} = 2^{3.1}$	4.68	73.13	46.27
R + Rotation equiv. (StyleGAN3-R)	4.50	66.65	45.92				

Figure 15: Results with our alternative translation equivariance metric EQ-T_{frac}; higher is better.

E.1 Integer translation

The goal of our integer translation metric, EQ-T, is to measure how closely, on average, the output synthesis network \mathbf{G} matches a translated reference image when we translate the input of \mathbf{G} . In other words,

$$\begin{aligned} \text{EQ-T} &= \text{PSNR}_{\mathcal{W} \times \mathcal{X}^2 \times \mathcal{V} \times \mathcal{C}}(I_t, K_t), \\ I_t(\mathbf{w}, \mathbf{x}, \mathbf{p}, c) &= \mathbf{T}_x[\mathbf{G}(z_0; \mathbf{w})][\mathbf{p}, c], \\ K_t(\mathbf{w}, \mathbf{x}, \mathbf{p}, c) &= \mathbf{G}(\mathbf{t}_x[z_0]; \mathbf{w})[\mathbf{p}, c], \end{aligned} \quad (20)$$

where $\mathbf{w} \sim \mathcal{W}$ is a random intermediate latent code produced by the mapping network, $\mathbf{x} = (x_0, x_1) \sim \mathcal{X}^2$ is a random translation offset, \mathbf{p} enumerates pixel locations in the mutually valid region \mathcal{V} , $c \sim \mathcal{C}$ is the color channel, and z_0 represents the input Fourier features. For integer translations, we sample the translation offsets x_0 and x_1 from $\mathcal{X} = \mathcal{U}[-s_N/8, s_N/8]$, where s_N is the width of the image in pixels.

In practice, we estimate the expectation in Equation 20 as an average over 50,000 random samples of $(\mathbf{w}, \mathbf{x}) \sim \mathcal{W} \times \mathcal{X}^2$. For given \mathbf{w} and \mathbf{x} , we generate the reference image I_t by running the synthesis network and translating the resulting image by \mathbf{x} pixels (operator \mathbf{T}_x). We then obtain the approximate result image K_t by translating the input Fourier features by the corresponding amount (operator \mathbf{t}_x), as discussed in Appendix F.1, and running the synthesis network again. The mutually valid region of I_t (translated by (x_0, x_1)) and K_t (translated by $(0, 0)$) is given by

$$\begin{aligned} \mathcal{V} = \{ &\max(x_0, 0), \dots, s_N + \min(x_0, 0) - 1 \} \times \\ &\{ \max(x_1, 0), \dots, s_N + \min(x_1, 0) - 1 \}. \end{aligned} \quad (21)$$

E.2 Fractional translation

Our translation equivariance metric has the nice property that, for a perfectly equivariant generator, the value of EQ-T converges to ∞ dB when the number of samples tends to infinity. However, this comes at the cost of completely ignoring subpixel effects. In fact, it is easy to imagine a generator that is perfectly equivariant to integer translation but fails with subpixel translation; in principle, this is true for any generator whose output is not properly bandlimited, including, e.g., implicit coordinate-based MLPs [4].

To verify that our generators *are* able to handle subpixel translation, we define an alternative translation equivariance metric, EQ-T_{frac}, where the translation offsets x_0 and x_1 are sampled from a continuous distribution $\mathcal{X} = \mathcal{U}(-s_N/8, s_N/8)$. While the continuous operator \mathbf{t}_x readily supports this new definition with fractional offsets, extending the discrete \mathbf{T}_x is slightly more tricky.

In practice, we define \mathbf{T}_x via standard Lanczos resampling, by filtering the image produced by \mathbf{G} using the prototypical Lanczos filter (Equation 15) with $a = 3$, evaluated at integer tap locations offset by \mathbf{x} . We explicitly normalize the resulting discretized filter to enforce the partition of unity property. We also shrink the mutually valid region to account for the spatial extent a by redefining

$$\begin{aligned} \mathcal{V} = \{ &\max(x_0 + a, 0), \dots, s_N + \min(x_0 - a, -1) \} \times \\ &\{ \max(x_1 + a, 0), \dots, s_N + \min(x_1 - a, -1) \}. \end{aligned} \quad (22)$$

Figure 15 compares the results of the two metrics, EQ-T and EQ-T_{frac}, using the same training configurations as Figure 3 in the main paper. The metrics agree reasonably well up until ~ 40 dB,

after which the fractional metric starts to saturate; it consistently fails to rise above 50 dB in our tests. This is due to the fact that the definition of subpixel translation is inherently ambiguous. The choice of the resampling filter represents a tradeoff between aliasing, ringing, and retention of high frequencies; there is no reason to assume that the generator would necessarily have to make the same tradeoff as the metric. Based on the results, we conclude that our configs G–R are essentially perfectly equivariant to subpixel translation within the limits of Lanczos resampling’s accuracy. However, due to its inherent limitations, we refrain from choosing EQ-T_{frac} as our primary metric.

E.3 Rotation

Measuring equivariance with respect to arbitrary rotations has the same fundamental limitation as our EQ-T_{frac} metric: the resampling operation is inherently ambiguous, so we cannot expect the results to be perfectly accurate beyond ~ 40 dB. Arbitrary rotations also have the additional complication that the bandlimit of a discretely sampled image is not radially symmetric.

Consider rotating the continuous representation of a discretely sampled image by 45° . The original frequency content of the image is constrained within the rectangular bandlimit $f \in [-s_N/2, +s_N/2]^2$. The frequency content of the rotated image, however, forms a diamond shape that extends all the way to $\|f\| = \sqrt{2}s_N/2$ along the main axes but only to $\|f\| = s_N/2$ along the diagonals. In other words, it simultaneously has too much frequency content, but also too little. This has two implications. First, in order to obtain a valid discretized result image, we have to low-pass filter the image *both* before *and* after the rotation to completely eliminate aliasing. Second, even if we are successful in eliminating the aliasing, the rotated image will still lack the highest representable diagonal frequencies. The second point further implies that when computing PSNR, our reference image I will inevitably lack some frequencies that are present in the output of \mathbf{G} . To obtain the correct result, we must eliminate these extraneous frequencies — without modifying the output image in any other way.

Based on the above reasoning, we define our EQ-R metric as follows:

$$\begin{aligned} \text{EQ-R} &= \text{PSNR}_{W \times A \times V \times C}(I_r, K_r), \\ I_r(\mathbf{w}, \alpha, \mathbf{p}, c) &= \mathbf{R}_\alpha[\mathbf{G}(z_0; \mathbf{w})][\mathbf{p}, c], \\ K_r(\mathbf{w}, \alpha, \mathbf{p}, c) &= \mathbf{R}_\alpha^*[\mathbf{G}(\mathbf{r}_\alpha[z_0]; \mathbf{w})][\mathbf{p}, c], \end{aligned} \quad (23)$$

where the random rotation angle α is drawn from $\mathcal{A} = \mathcal{U}(0^\circ, 360^\circ)$ and operator \mathbf{r}_α corresponds to continuous rotation of the input Fourier features by α with respect to the center of the canvas $[0, 1]^2$. \mathbf{R}_α corresponds to high-quality rotation of the reference image, and \mathbf{R}_α^* represents a *pseudo-rotation* operator that modifies the frequency content of the image *as if* it had undergone \mathbf{R}_α — but without *actually* rotating it.

The *ideal* rotation operator $\hat{\mathbf{R}}$ is easily defined under our theoretical framework presented in Section 2.1:

$$\hat{\mathbf{R}}_\alpha[Z] = \text{III} \odot (\psi * \mathbf{r}_\alpha[\phi * Z]) = 1/s^2 \cdot \text{III} \odot (\psi * \mathbf{r}_\alpha[\psi * Z]). \quad (24)$$

In other words, we first convolve the discretely sampled input image Z with ϕ to obtain the corresponding continuous representation. We then rotate this continuous representation using \mathbf{r}_α , bandlimit the result by convolving with ψ , and finally extract the corresponding discrete representation by multiplying with III . To reduce notational clutter, we omit the subscripts denoting the sampling rate s . We can swap the order of the rotation and a convolution in the above formula by rotating the kernel in the opposite direction to compensate:

$$\hat{\mathbf{R}}_\alpha[Z] = 1/s^2 \cdot \text{III} \odot \mathbf{r}_\alpha[\hat{h}_R * Z], \quad \hat{h}_R = \mathbf{r}_{-\alpha}[\psi] * \psi, \quad (25)$$

where \hat{h}_R represents an ideal “rotation filter” that bandlimits the signal with respect to both the input and the output. Its spectrum is the eight-sided polygonal intersection of the original and the rotated rectangle.

In order to obtain a practical approximation \mathbf{R}_α , we must replace \hat{h}_R with an approximate filter h_R that has finite support. Given such a filter, we get $\mathbf{R}_\alpha[Z] = 1/s^2 \cdot \text{III} \odot \mathbf{r}_\alpha[h_R * Z]$. In practice, we implement this operation using two additional approximations. First, we approximate $1/s^2 \cdot h_R * Z$ by an upsampling operation to a higher temporary resolution, using h_R as the upsampling filter and $m = 4$. Second, we approximate $\text{III} \odot \mathbf{r}_\alpha$ by performing a set of bilinear lookups from the temporary high-resolution image.

To obtain h_R , we again utilize the standard Lanczos window with $a = 3$:

$$h_R = (\mathbf{r}_{-\alpha}[\psi] * \psi) \odot (\mathbf{r}_{-\alpha}[w_L^+] * w_L^+), \quad (26)$$

where we apply the same rotation-convolution to both the filter and the window function. w_L^+ corresponds to the canonical separable Lanczos window, similar to the one used in Equation 15:

$$w_L^+(\mathbf{x}) = \begin{cases} \text{sinc}(x_0/a) \cdot \text{sinc}(x_1/a), & \text{if } \max(|x_0|, |x_1|) < a, \\ 0, & \text{if } \max(|x_0|, |x_1|) \geq a, \end{cases} \quad (27)$$

We can now define the pseudo-rotation operator $\mathbf{R}_\alpha^*[Z]$ as a simple convolution with another filter that resembles h_R :

$$\begin{aligned} \mathbf{R}_\alpha^*[Z] &= 1/s^2 \cdot \text{III} \odot (h_R^* * Z) = H_R^* * Z, \\ h_R^* &= (\psi * \mathbf{r}_\alpha[\psi]) \odot (w_L^+ * \mathbf{r}_\alpha[w_L^+]), \end{aligned} \quad (28)$$

where the discrete version H_R^* is obtained from h_R^* using Equation 6.

Finally, we define the valid region \mathcal{V} the same way as in Appendix E.2: the set of pixels for which both filter footprints fall within the bounds of the corresponding original images.

F Implementation details

We implemented our alias-free generator on top of the official PyTorch implementation of StyleGAN2-ADA, available at <https://github.com/NVlabs/stylegan2-ada-pytorch>. We kept most of the details unchanged, including discriminator architecture [34], weight demodulation [34], equalized learning rate for all trainable parameters [31], minibatch standard deviation layer at the end of the discriminator [31], exponential moving average of generator weights [31], mixed-precision FP16/FP32 training [32], non-saturating logistic loss [21], R_1 regularization [40], lazy regularization [34], and Adam optimizer [36] with $\beta_1 = 0$, $\beta_2 = 0.99$, and $\epsilon = 10^{-8}$.

We ran all experiments on NVIDIA DGX-1 with 8 Tesla V100 GPUs using PyTorch 1.7.1, CUDA 11.0, and cuDNN 8.0.5. We computed FID between 50k generated images and all training images using the official pre-trained Inception network, available at <http://download.tensorflow.org/models/image/imagenet/inception-2015-12-05.tgz>

Our implementation and pre-trained models are available at <https://github.com/NVlabs/stylegan3>

F.1 Generator architecture

Normalization (configs D–R) We have observed that eliminating the output skip connections in StyleGAN2 [34] results in uncontrolled drift of signal magnitudes over the generator layers. This does not necessarily lead to lower-quality results, but it generally increases the amount of random variation between training runs and may occasionally lead to numerical issues with mixed-precision training. We eliminate the drift by tracking a long-term exponential moving average of the input signal magnitude on each layer and normalizing the feature maps accordingly. We update the moving average once per training iteration, based on the mean of squares over the entire input tensor, and freeze its value after training. We initialize the moving average to 1 and decay it at a constant rate, resulting in 50% decay per 20k real images shown to the discriminator. With this explicit normalization in place, we have found it beneficial to slightly adjust the dynamic range of the output RGB colors. StyleGAN2 uses -1 and $+1$ to represent black and white, respectively; we change these values to -4 and $+4$ starting from config D and, for consistency with the original generator, divide the color channels by 4 afterwards.

Transformed Fourier features (configs H–R) We enable the orientation of the input features z_0 to vary on a per-image basis by introducing an additional affine layer (Figure 4b) and applying a geometric transformation based on its output. The affine layer produces a four-dimensional vector $\mathbf{t} = (r_c, r_s, t_x, t_y)$ based on \mathbf{w} . We initialize its weights so that $\mathbf{t} = (1, 0, 0, 0)$ at the beginning, but allow them to change freely over the course of training. To interpret \mathbf{t} as a geometric transformation, we first normalize its value based on the first two components, i.e., $\mathbf{t}' = (r'_c, r'_s, t'_x, t'_y) = \mathbf{t} / \sqrt{r_c^2 + r_s^2}$. This makes the transformation independent of the magnitude of \mathbf{w} , similar to the weight modulation and

Parameter	Config	Datasets (Figure 5, left)			Ablations at 256×256			R_1 regularization γ	B	T	R
		B	T	R	A–C	D–T	R				
Batch size		32	32	32	64	64	64	FFHQ-U	256 ²	1.0	1.0
Moving average		10k	10k	10k	20k	20k	20k	FFHQ-U	1024 ²	10.0	32.8
Mapping net depth		8	2	2	8	2	2	FFHQ	1024 ²	10.0	32.8
Minibatch stddev		4	4	4	8	4	4	METFACES-U	1024 ²	10.0	16.4
G layers		15/17	14	14	13	14	14	METFACES	1024 ²	5.0	6.6
G capacity: C_{base}		2^{15}	2^{15}	2^{16}	2^{14}	2^{14}	2^{15}	AFHQv2	512 ²	5.0	8.2
G capacity: C_{max}		512	512	1024	512	512	1024	BEACHES	512 ²	2.0	4.1
G learning rate		0.0020	0.0025	0.0025	0.0025	0.0025	0.0025				
D learning rate		0.0020	0.0020	0.0020	0.0025	0.0025	0.0025				

Figure 16: **Left:** Hyperparameters used in each experiment. **Right:** R_1 regularization weights.

demodulation [34] on the other layers. We then interpret the first two components as rotation around the center of the canvas $[0, 1]^2$, with the rotation angle α defined by $r'_c = \cos \alpha$ and $r'_s = \sin \alpha$. Finally, we interpret the remaining two components as translation by (t'_x, t'_y) units, so that the translation is performed after the rotation. In practice, we implement the resulting geometric transformation by modifying the phases and two-dimensional frequencies of the Fourier features, which is equivalent to applying the same transformation to the continuous representation of z_0 analytically.

Flexible layer specifications In configs T and R, we define the per-layer filter parameters (Figure 4c) as follows. The cutoff frequency f_c and the minimum acceptable stopband frequency f_t obey geometric progression until the first critically sampled layer:

$$f_c[i] = f_{c,0} \cdot (f_{c,N}/f_{c,0})^{\min(i/(N-N_{\text{crit}}), 1)}, \quad f_t[i] = f_{t,0} \cdot (f_{t,N}/f_{t,0})^{\min(i/(N-N_{\text{crit}}), 1)}, \quad (29)$$

where $N = 14$ is the total number of layers, $N_{\text{crit}} = 2$ is the number of critically sampled layers at the end, $f_{c,0} = 2$ corresponds to the frequency content of the input Fourier features, and $f_{c,N} = s_N/2$ is defined by the output resolution. $f_{t,0}$ and $f_{t,N}$ are free parameters; we use $f_{t,0} = 2^{2.1}$ and $f_{t,N} = f_{c,N} \cdot 2^{0.3}$ in most of our tests. Given the values of $f_c[i]$ and $f_t[i]$, the sampling rate $s[i]$ and transition band half-width $f_h[i]$ are then determined by

$$s[i] = \exp_2 \lceil \log_2 (\min(2 \cdot f_t[i], s_N)) \rceil, \quad f_h[i] = \max(f_t[i], s[i]/2) - f_c[i]. \quad (30)$$

The sampling rate is rounded up to the nearest power of two that satisfies $s[i] \geq 2f_t[i]$, but it is not allowed to exceed the output resolution. The transition band half-width is selected to satisfy either $f_c[i] + f_h[i] = f_t[i]$ or $f_c[i] + f_h[i] = s[i]/2$, whichever yields a higher value.

We consider $f_c[i]$ to represent the output frequency content of layer i , for $i \in \{0, 1, \dots, N-1\}$, whereas the input is represented by $f_c[\max(i-1, 0)]$. Thus, we construct the corresponding upsampling filter according to $f_c[\max(i-1, 0)]$ and $f_h[\max(i-1, 0)]$ and the downsampling filter according to $f_c[i]$ and $f_h[i]$. The nonlinearity is evaluated at a temporary sampling rate $s' = \max(s[i], s[\max(i-1, 0)]) \cdot m$, where m is the upsampling parameter discussed in Section 3.2 that we set to 2 in most of our tests.

F.2 Hyperparameters and training configurations

We used 8 GPUs for all our training runs and continued the training until the discriminator had seen a total of 25M real images when training from scratch, or 5M images when using transfer learning. Figure 16 shows the hyperparameters used in each experiment. We performed the baseline runs (configs A–C) using the corresponding standard configurations: StyleGAN2 config F [34] for the high-resolution datasets in Figure 5, left, and ADA 256×256 baseline config [32] for the ablations in Figure 3 and Figure 5, right.

Many of our hyperparameters, including discriminator capacity and learning rate, batch size, and generator moving average decay, are inherited directly from the baseline configurations, and kept unchanged in all experiments. In configs C and D, we disable noise inputs [33], path length regularization [34], and mixing regularization [33]. In config D, we also decrease the mapping network depth to 2 and set the minibatch standard deviation group size to 4 as recommended in the StyleGAN2-ADA documentation. The introduction of explicit normalization in config D allows us to use the same generator learning rate, 0.0025, for all output resolutions. In Figure 5, right, we show results for path length regularization with weight 0.5 and mixing regularization with probability 0.5.

Augmentation Since our datasets are horizontally symmetric in nature, we enable dataset x -flip augmentation in all our experiments. To prevent the discriminator from overfitting, we enable adaptive discriminator augmentation (ADA) [32] with default settings for METFACES, METFACES-U, AFHQV2, and BEACHES, but disable it for FFHQ and FFHQ-U. Furthermore, we train METFACES and METFACES-U using transfer learning from the corresponding FFHQ or FFHQ-U snapshot with the lowest FID, similar to Karras et al. [32], but start the training from scratch in all other experiments.

Generator capacity StyleGAN2 defines the number of feature maps on a given layer to be inversely proportional to its resolution, i.e., $C[i] = C(s[i]) = \min(\text{round}(C_{\text{base}}/s[i]), C_{\text{max}})$, where $s[i]$ is the output resolution of layer i . Parameters C_{base} and C_{max} control the overall capacity of the generator; our baseline configurations use $C_{\text{max}} = 512$ and $C_{\text{base}} = 2^{14}$ or 2^{15} depending on the output resolution. Since StyleGAN2 can be considered to employ critical sampling on all layers, i.e., $f_c[i] = s[i]/2$, we can equally well define the number of feature maps as $C[i] = C(2f_c[i])$. These two definitions are equivalent for configs A–F, but in configs G–R we explicitly set $f_c[i] \leq s[i]/2$, which necessitates using the latter definition. In config R, we double the value of both C_{base} and C_{max} to compensate for the reduced capacity of the 1×1 convolutions. In Figure 5, right, we sweep the capacity by multiplying both parameters by 0.5, 1.0, and 2.0.

R_1 regularization The optimal choice for the R_1 regularization weight γ is highly dependent on the dataset, necessitating a grid search [34, 32]. For the baseline config B, we tested $\gamma \in \{1, 2, 5, 10, 20\}$ and selected the value that gave the best FID for each dataset. For our configs T and R, we followed the recommendation of Karras et al. [32] to define $\gamma = \gamma_0 \cdot N/M$, where $N = s_N^2$ is the number of output pixels and M is the batch size, and performed a grid search over $\gamma_0 \in \{0.0002, 0.0005, 0.0010, 0.0020, 0.0050\}$. For the low-resolution ablations, we chose to use a fixed value $\gamma = 1$ for simplicity. The resulting values of γ are shown in Figure 16, right.

Training of config R In this configuration, we blur all images the discriminator sees in the beginning of the training. This Gaussian blur is executed just before the ADA augmentation. We start with $\sigma = 10$ pixels, which we ramp to zero over the first 200k images. This prevents the discriminator from focusing too heavily on high frequencies early on. It seems that in this configuration the generator sometimes learns to produce high frequencies with a small delay, allowing the discriminator to trivially tell training data from the generated images without providing useful feedback to the generator. As such, config R is prone to random training failures in the beginning of the training without this trick. The other configurations do not have this issue.

E.3 G-CNN comparison

In Figure 5, bottom, we compare our config R with config T extended with $p4$ -symmetric group convolutions [16, 17]. $p4$ symmetry makes the generator equivariant to 0° , 90° , 180° , and 270° rotations, but not to arbitrary rotation angles. In practice, we implement the group convolutions by extending all intermediate activation tensors in the synthesis network with an additional group dimension of size 4 and introducing appropriate redundancy in the convolution weights. We keep the input layer unchanged and introduce the group dimension by replicating each element of z_0 four times. Similarly, we eliminate the group dimension after the last layer by computing an average of the four elements. $p4$ -symmetric group convolutions have $4 \times$ as many trainable parameters as the corresponding regular convolutions. To enable an apples-to-apples comparison, we compensate for this increase by halving the values of C_{base} and C_{max} , which brings the number of parameters back to the original level.

G Energy consumption

Computation is an essential resource in machine learning projects: its availability and cost, as well as the associated energy consumption, are key factors in both choosing research directions and practical adoption. We provide a detailed breakdown for our entire project in Table 17 in terms of both GPU time and electricity consumption. We report expended computational effort as single-GPU years (Volta class GPU). We used a varying number of NVIDIA DGX-1s for different stages of the project, and converted each run to single-GPU equivalents by simply scaling by the number of GPUs used.

Item	Number of training runs	GPU years (Volta)	Electricity (MWh)
Early exploration	233	18.02	42.45
Project exploration	1207	48.93	118.13
Setting up ablations	297	13.30	32.48
Per-dataset tuning	63	4.54	13.28
Producing results in the paper	53	5.26	14.35
StyleGAN3-R at 1024×1024	1	0.30	0.87
Other runs in the dataset table	17	2.35	6.88
Ablation tables	35	2.61	6.60
Results intentionally left out	23	1.72	3.93
Total	1876	91.77	224.62

Figure 17: Computational effort expenditure and electricity consumption data for this project. The unit for computation is GPU-years on a single NVIDIA V100 GPU—it would have taken approximately 92 years to execute this project using a single GPU. See the text for additional details about the computation and energy consumption estimates. **Early exploration** includes early training runs that affected our decision to start this project. **Project exploration** includes training runs that were done specifically for this project, leading to the final StyleGAN3-T and StyleGAN3-R configurations. These runs were not intended to be used in the paper as-is. **Setting up ablations** includes hyperparameter tuning for the intermediate configurations and ablation experiments in Figures 3 and 5. **Per-dataset tuning** includes hyperparameter tuning for individual datasets, mainly the grid search for R_1 regularization weight. **Config R at 1024×1024** corresponds to one training run in Figure 5, left, and **Other runs in the dataset table** includes the remaining runs. **Ablation tables** includes the low-resolution ablations in Figures 3 and Figure 5. **Results intentionally left out** includes additional results that were initially planned, but then left out to improve focus and clarity.

We followed the Green500 power measurements guidelines [20]. The entire project consumed approximately 225 megawatt hours (MWh) of electricity. Approximately 70% of it was used for exploratory runs, where we gradually built the new configurations; first in an unstructured manner and then specifically ironing out the new StyleGAN3-T and StyleGAN3-R configurations. Setting up the intermediate configurations between StyleGAN2 and our generators, as well as, the key parameter ablations was also quite expensive at $\sim 15\%$. Training a single instance of StyleGAN3-R at 1024×1024 is only slightly more expensive (0.9MWh) than training StyleGAN2 (0.7MWh) [34].