

# LOGAN: LATENT OPTIMISATION FOR GENERATIVE ADVERSARIAL NETWORKS

**Yan Wu, Jeff Donahue, David Balduzzi, Karen Simonyan, Timothy Lillicrap**

DeepMind

London, UK

{yanwu, jeffdonahue, dbalduzzi, simonyan, countzero}@google.com

## ABSTRACT

Training generative adversarial networks requires balancing of delicate adversarial dynamics. Even with careful tuning, training may diverge or end up in a bad equilibrium with dropped modes. In this work, we improve CS-GAN with natural gradient-based latent optimisation and show that it improves adversarial dynamics by enhancing interactions between the discriminator and the generator. Our experiments demonstrate that latent optimisation can significantly improve GAN training, obtaining state-of-the-art performance for the ImageNet ( $128 \times 128$ ) dataset. Our model achieves an Inception Score (IS) of 148 and an Fréchet Inception Distance (FID) of 3.4, an improvement of 17% and 32% in IS and FID respectively, compared with the baseline BigGAN-deep model with the same architecture and number of parameters.

## 1 INTRODUCTION

Generative Adversarial Nets (GANs) are implicit generative models that can be trained to match a given data distribution. GANs were originally developed by Goodfellow et al. (2014) for image data. As the field of generative modelling has advanced, GANs remain at the frontier, generating high-fidelity images at large scale (Brock et al., 2018; Karras et al., 2019). However, despite growing insights into the dynamics of GAN training, much of the progress in GAN-based image generation come from network architecture improvements (Radford et al., 2015; Zhang et al., 2019), or regularisation of particular parts of the model (Miyato et al., 2018; Miyato & Koyama, 2018).

Build on the compressed sensing view of GANs (CS-GAN; Wu et al., 2019), we improve the efficacy of latent optimisation in adversarial games, using natural gradient descent to optimise the latent variable (usually denoted  $z$ ) towards the direction favoured by the discriminator during training. This results in a scalable and easy to implement approach that improves the dynamic interaction between the discriminator and the generator. We generally call these approaches latent optimised GANs (LOGAN).

To summarise our contributions:

1. We propose an improved, efficient approach to latent optimisation using natural gradient descent.
2. Our algorithm improves the state-of-the-art BigGAN (Brock et al., 2018) by a significant margin, without introducing any architectural change, resulting in higher quality images and more diverse samples (see Table 1, Figure 1 and 2).
3. To provide theoretical insight, we analyse latent optimisation in GANs from the perspective of differentiable games (Balduzzi et al., 2018). We argue that latent optimisation can be viewed as improving the dynamics of adversarial training.



Figure 1: Samples from BigGAN-deep **(a)** and LOGAN **(b)** with similarly high IS. Samples from the two panels were drawn from truncation levels corresponding to points C and D in Figure 3 **b** respectively. (FID/IS: **(a)** 27.97/259.4, **(b)** 8.19/259.9)

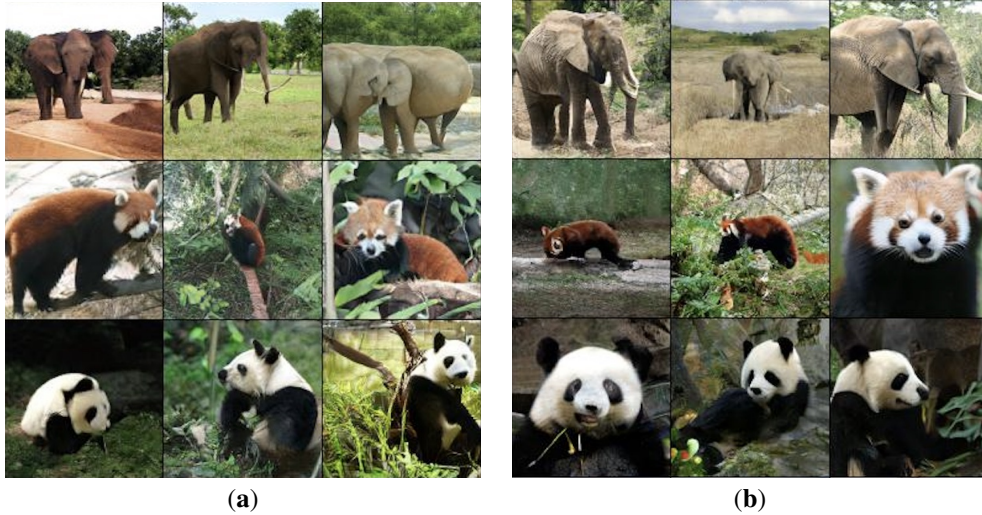


Figure 2: Samples from BigGAN-deep **(a)** and LOGAN **(b)** with similarly low FID. Samples from the two panels were drawn from truncation levels corresponding to points A and B in Figure 3 **b** respectively. (FID/IS: **(a)** 5.04/126.8, **(b)** 5.09/217.0)

## 2 BACKGROUND

### 2.1 NOTATION

We use  $\theta_D$  and  $\theta_G$  to denote the vectors representing parameters of the generator and discriminator. We use  $x$  for images, and  $z$  for the latent source generating an image. We use prime ' to denote a variable after one update step, e.g.,  $\theta'_D = \theta_D - \alpha \frac{\partial f(z; \theta_D, \theta_G)}{\partial \theta_D}$ .  $p(x)$  and  $p(z)$  denote the data distribution and source distribution respectively.  $\mathbb{E}_{p(x)} [f(x)]$  indicates taking the expectation of function  $f(x)$  over the distribution  $p(x)$ .

### 2.2 GENERATIVE ADVERSARIAL NETS

A GAN consists of a generator that generates image  $x = G(z; \theta_G)$  from a latent source  $z \sim p(z)$ , and a discriminator that scores the generated images as  $D(x; \theta_D)$  (Goodfellow et al., 2014). Training GANs involves an adversarial game: while the discriminator tries to distinguish generated samples

Table 1: Comparison of model scores. BigGAN-deep results are reproduced from Brock et al. (2018). “baseline” indicates our reproduced BigGAN-deep with small modifications. The 3rd and 4th columns are from the gradient descent (GD, ablated) and natural gradient descent (NGD) versions of LOGAN respectively. We report the Inception Score (IS, higher is better, Salimans et al. 2016) and Fréchet Inception Distance (FID, lower is better, Heusel et al. 2017).

	FID	IS
BigGAN-deep	$5.7 \pm 0.3$	$124.5 \pm 2.0$
baseline	$4.92 \pm 0.05$	$126.6 \pm 1.3$
LOGAN (GD)	$4.86 \pm 0.09$	$127.7 \pm 3.5$
LOGAN (NGD)	<b><math>3.36 \pm 0.14</math></b>	<b><math>148.2 \pm 3.1</math></b>

$x = G(z; \theta_G)$  from data  $x \sim p(x)$ , the generator tries to fool the discriminator. This procedure can be summarised as the following min-max game:

$$\min_{\theta_D} \max_{\theta_G} \mathbb{E}_{x \sim p(x)} [h_D(D(x; \theta_D))] + \mathbb{E}_{z \sim p(z)} [h_G(D(G(z; \theta_G); \theta_D))] \quad (1)$$

The exact form of  $h(\cdot)$  depends on the choice of loss function (Goodfellow et al., 2014; Arjovsky et al., 2017; Nowozin et al., 2016). To simplify our presentation and analysis, we use the Wasserstein loss (Arjovsky et al., 2017), so that  $h_D(t) = -t$  and  $h_G(t) = t$ . Our experiments with BigGAN-deep use the hinge loss (Lim & Ye, 2017; Tran et al., 2017), which is identical to this form in its linear regime. Our analysis can be generalised to other losses as in previous theoretical work (e.g., Arora et al. 2017). To simplify notation, we abbreviate  $f(z; \theta_D, \theta_G) = D(G(z; \theta_G); \theta_D)$ , which may be further simplified as  $f(z)$  when the explicit dependency on  $\theta_D$  and  $\theta_G$  can be omitted.

Training GANs requires carefully balancing updates to  $D$  and  $G$ , and is sensitive to both architecture and algorithm choices (Salimans et al., 2016; Radford et al., 2015). A recent milestone is BigGAN (and BigGAN-deep, Brock et al. 2018), which pushed the boundary of high fidelity image generation by scaling up GANs to an unprecedented level. BigGANs use an architecture based on residual blocks (He et al., 2016), in combination with regularisation mechanisms and self-attention (Saxe et al., 2014; Miyato et al., 2018; Zhang et al., 2019).

Here we aim to improve the adversarial dynamics during training. We focus on the second term in eq. 1 which is at the heart of the min-max game. For clarity, we explicitly write the losses for  $D$  as  $L_D(z) = h_D(f(z))$  and  $G$  as  $L_G(z) = h_G(f(z))$ , so the total loss vector can be written as

$$L(z) = [L_D(z), L_G(z)]^T = [f(z), -f(z)]^T \quad (2)$$

Computing the gradients with respect to  $\theta_D$  and  $\theta_G$  gives the following vector field, which *cannot* be expressed as the gradient of any single function (Balduzzi et al., 2018):

$$g = \left[ \frac{\partial L_D(z)}{\partial \theta_D}, \frac{\partial L_G(z)}{\partial \theta_G} \right]^T = \left[ \frac{\partial f(z)}{\partial \theta_D}, -\frac{\partial f(z)}{\partial \theta_G} \right]^T \quad (3)$$

The fact that  $g$  is not the gradient of a function implies that gradient updates in GANs can exhibit cycling behaviour which can slow down or prevent convergence. Balduzzi et al. (2018) refer to vector fields of this form as the *simultaneous gradient*. Although many GAN models use alternating update rules (e.g., Goodfellow et al. 2014; Brock et al. 2018), following the gradient with respect to  $\theta_D$  and  $\theta_G$  alternatively in each step, they share the same problem from gradients of this form. Therefore, we use the simpler simultaneous gradient (eq. 3) for our analysis (see also Mescheder et al. 2017; 2018).

### 2.3 LATENT OPTIMISED GANS

Inspired by compressed sensing (Candes et al., 2006; Donoho, 2006), Wu et al. (2019) introduced latent optimisation for GANs. Latent optimisation exploits knowledge from  $D$  to refine the latent source  $z$ . Intuitively, the gradient  $\nabla_z f(z) = \frac{\partial f(z)}{\partial z}$  points in the direction that better satisfies the discriminator  $D$ , which implies better samples. Therefore, instead of using the randomly sampled

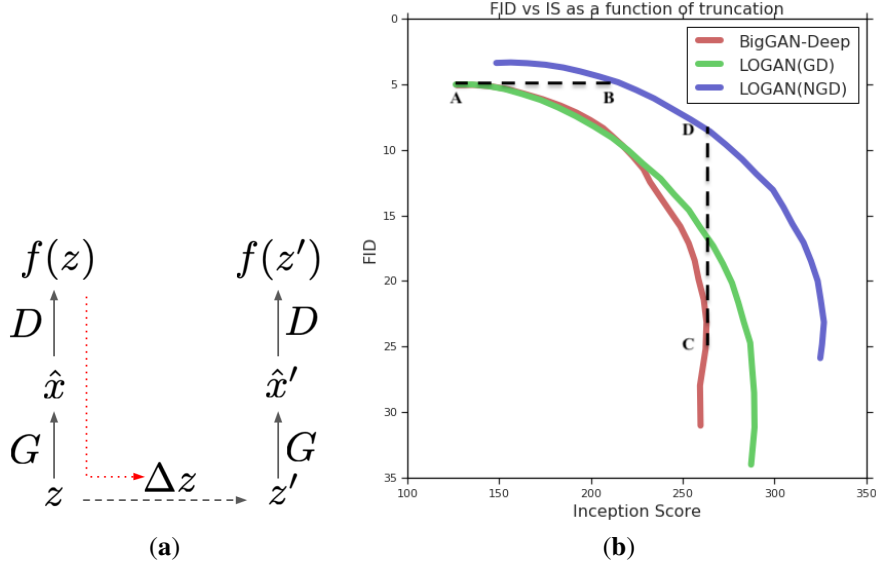


Figure 3: **(a)** Schematic of LOGAN. We first compute a forward pass through  $G$  and  $D$  with a sampled latent  $z$ . Then, we use gradients from the generator loss (dashed red arrow) to compute an improved latent,  $z'$ . After we use this optimised latent code in a second forward pass, we compute gradients of the discriminator back through the latent optimisation into the model parameters  $\theta_D$ ,  $\theta_G$ . We use these gradients to update the model. **(b)** Truncation curves illustrate the FID/IS trade-off for each model by altering the range of the noise source  $p(z)$ . GD: gradient descent. NGD: natural gradient descent. Points A, B, C, D correspond to samples shown in Figure 1 and 2.

---

**Algorithm 1** Latent Optimised GANs with Automatic Differentiation

---

**Input:** data distribution  $p(x)$ , latent distribution  $p(z)$ ,  $D(\cdot; \theta_D)$ ,  $G(\cdot; \theta_G)$ , learning rate  $\alpha$ , batch size  $N$

**repeat**

    Initialise discriminator and generator parameters  $\theta_D, \theta_G$

**for**  $i = 1$  **to**  $N$  **do**

        Sample  $z \sim p(z)$ ,  $x \sim p(x)$

        Compute the gradient  $\frac{\partial D(G(z))}{\partial z}$  and use it to obtain  $\Delta z$  from eq. 4 (GD) or eq. 16 (NGD)

        Optimise the latent  $z' \leftarrow [z + \Delta z]$ ,  $[\cdot]$  indicates clipping the value between  $-1$  and  $1$

        Compute generator loss  $L_G^{(i)} = -D(G(z'))$

        Compute discriminator loss  $L_D^{(i)} = D(G(z')) - D(x)$

**end for**

    Compute batch losses  $L_G = \frac{1}{N} \sum_{i=1}^N L_G^{(i)}$  and  $L_D = \frac{1}{N} \sum_{i=1}^N L_D^{(i)}$

    Update  $\theta_D$  and  $\theta_G$  with the gradients  $\frac{\partial L_D}{\partial \theta_D}$ ,  $\frac{\partial L_G}{\partial \theta_G}$

**until** reaches the maximum training steps

---

$z \sim p(z)$ , Wu et al. (2019) uses the optimised latent

$$\Delta z = \alpha \frac{\partial f(z)}{\partial z} \quad z' = z + \Delta z \quad (4)$$

in eq. 1 for training <sup>1</sup>.

Historically, compressed sensing has been developed as a signal processing technique mostly without any concern on training. However, here we emphasise the influence of this procedure on training, which we will show dominates the effects on large scale models — in contrast, the run-time optimisation that is central in compressed sensing may be unnecessary after training. Therefore, we call

---

<sup>1</sup>Although multiple gradient descent steps can be employed for optimising  $z$ , we found one step works well in training and justify this choice in section 3.



this type of models latent-optimised GANs (LOGAN) to avoid any confusion, except when explicitly referring to the results from Wu et al. (2019). Latent optimisation has been shown to improve the stability of training as well as the final performance for medium-sized models such as DCGANs and Spectral Normalised GANs (Radford et al., 2015; Miyato et al., 2018). The general algorithm is summarised in Algorithm 1 and illustrated in Figure 3 a. However, we found that the potential of latent optimisation remained largely untapped in this setting, and develop the natural gradient descent form of latent update in Section 4.

### 3 ANALYSIS OF THE ALGORITHM

To understand how latent optimisation interacts with GAN training, we analyse LOGAN as a differentiable game following Balduzzi et al. (2018); Gemp & Mahadevan (2018); Letcher et al. (2019). The Appendix A provides a complementary analysis from the perspective of stochastic approximation (Heusel et al., 2017; Borkar, 1997). We can explicitly compute the gradients for the discriminator and generator at  $z'$  after one step of latent optimisation by differentiating  $[L_D(z'), L_G(z')]^T = [f(z'), -f(z')]^T$  (where  $z' = z + \Delta z$  from eq. 4):

$$\left[ \frac{dL_D}{d\theta_D}, \frac{dL_G}{d\theta_G} \right]^T = \left[ \frac{\partial f(z')}{\partial \theta_D} + \left( \frac{\partial \Delta z}{\partial \theta_D} \right)^T \frac{\partial f(z')}{\partial \Delta z}, -\frac{\partial f(z')}{\partial \theta_G} - \left( \frac{\partial \Delta z}{\partial \theta_G} \right)^T \frac{\partial f(z')}{\partial \Delta z} \right]^T \quad (5)$$

$$= \left[ \frac{\partial f(z')}{\partial \theta_D} + \alpha \left( \frac{\partial^2 f(z)}{\partial z \partial \theta_D} \right)^T \frac{\partial f(z')}{\partial z'}, -\frac{\partial f(z')}{\partial \theta_G} - \alpha \left( \frac{\partial^2 f(z)}{\partial z \partial \theta_G} \right)^T \frac{\partial f(z')}{\partial z'} \right]^T \quad (6)$$

In both equations, the first terms represent how  $f(z')$  depends on the parameters directly, which also appear in the gradients from vanilla GANs (eq. 3). However, the second terms are introduced from latent optimisation, accounting for how  $f(z')$  depends on the parameters via the change  $\Delta z$ . For the second equality, we substitute  $\Delta z = \alpha \frac{\partial f(z)}{\partial z}$  as the gradient-based update of  $z$  and use  $\frac{\partial f(z')}{\partial \Delta z} = \frac{\partial f(z')}{\partial z'}$ . Further differentiating  $\Delta z$  results in the second-order terms  $\left( \frac{\partial^2 f(z)}{\partial z \partial \theta_D} \right)^T$  and  $\left( \frac{\partial^2 f(z)}{\partial z \partial \theta_G} \right)^T$ . The original GAN's gradient (eq. 3) does not include any second-order term, since  $\Delta z = 0$  without latent optimisation. LOGAN computes these extra terms by automatic differentiation when back-propagating through the latent optimisation process (see Algorithm 1).

#### 3.1 RELATION WITH SGA

Balduzzi et al. (2018); Gemp & Mahadevan (2018) proposed Symplectic Gradient Adjustment (SGA) to improve the dynamics of gradient-based methods in adversarial games. SGA addresses an important problem with gradient-based optimisation in GANs: the vector-field generated by the losses of the discriminator and generator is not a gradient vector field. It follows that gradient descent is not guaranteed to find a local optimum and can cycle, which can slow down convergence or lead to phenomena like mode collapse and mode hopping.

For a game with gradient  $g$  (eq. 3), the Hessian is the second order derivatives with respect to the parameters,  $H = \nabla_{\theta} g$ . SGA uses the adjusted gradient

$$g^* = g + \lambda A^T g \quad (7)$$

where  $\lambda$  is a positive constant and  $A = \frac{1}{2}(H - H^T)$  is the anti-symmetric component of the Hessian. Applying SGA to GANs yields the adjusted updates (see Appendix A.1 for details):

$$g^* = \left[ \frac{\partial f(z)}{\partial \theta_D} + \lambda \left( \frac{\partial^2 f(z)}{\partial \theta_G \partial \theta_D} \right)^T \frac{\partial f(z)}{\partial \theta_G}, -\frac{\partial f(z)}{\partial \theta_G} + \lambda \left( \frac{\partial^2 f(z)}{\partial \theta_D \partial \theta_G} \right)^T \frac{\partial f(z)}{\partial \theta_D} \right]^T \quad (8)$$

Compared with  $g$  in eq. 3, the adjusted gradient  $g^*$  has second-order terms reflecting the interactions between  $D$  and  $G$ . SGA significantly improves GAN training in simple examples (Balduzzi et al., 2018), allowing faster and more robust convergence to stable fixed points (local Nash equilibria). Unfortunately, SGA is expensive to scale because computing the second-order derivatives with

respect to all parameters is expensive. It remains unclear whether SGA can be incorporated into very large scale models using more efficient implementation (e.g., Hessian-vector products from modified back propagation Pearlmutter (1994)).

The SGA updates in eq. 8 and the LOGAN updates in eq. 6 are strikingly similar, suggesting that the latent step used by LOGAN reduces the negative effects of cycling by introducing a symplectic gradient adjustment into the optimisation procedure. The role of the latent step can be formalised in terms of a third player, whose goal is to help the generator (see appendix A for details). Crucially, latent optimisation approximates SGA using only second-order derivatives with respect to the latent  $z$  and parameters of the discriminator and generator *separately*. The second order terms involving parameters of both the discriminator and the generator – which are expensive to compute – are not used. In short, with a simple modification of the original GAN training algorithm, latent optimisation couples the gradients of the discriminator and generator in a way similar to SGA.

### 3.2 RELATION WITH UNROLLED GANS

In addition, latent optimisation can be seen as unrolling GANs (Metz et al., 2016) in the space of the latent source, rather than the parameters. Unrolling in the latent space has the advantages that:

1. LOGAN is more scalable than Unrolled GANs because it avoids unrolling the parameter updating process, which is prohibitively expensive for models with a large number of parameters.
2. While unrolling the update of  $D$  only affects the parameters of  $G$  (as in Metz et al. 2016), latent optimisation effects both  $D$  and  $G$  as shown in eq. 6.

We next formally present this connection by first showing that SGA can be seen as approximating Unrolled GANs (Metz et al., 2016). For the update  $\theta'_D = \theta_D + \Delta\theta_D$ , we have the Taylor expansion approximation at  $\theta_D$ :

$$f(z; \theta_D + \Delta\theta_D, \theta_G) \approx f(z; \theta_D, \theta_G) + \left( \frac{\partial f(z; \theta_D, \theta_G)}{\partial \theta_D} \right)^T \Delta\theta_D \quad (9)$$

Substitute the gradient descent parameter update  $\Delta\theta_D = -\alpha \frac{\partial f(z; \theta_D, \theta_G)}{\partial \theta_D}$ , and take the derivatives with respect to  $\theta_G$  on both sides:

$$\frac{\partial f(z; \theta_D + \Delta\theta_D, \theta_G)}{\partial \theta_G} \approx \frac{\partial f(z; \theta_D, \theta_G)}{\partial \theta_G} - 2\alpha \left( \frac{\partial^2 f(z; \theta_D, \theta_G)}{\partial \theta_D \partial \theta_G} \right)^T \frac{\partial f(z; \theta_D, \theta_G)}{\partial \theta_D} \quad (10)$$

which has the same form as eq. 8 (taking the negative sign). Compared with the exact gradient from the unroll:

$$\frac{\partial f(z; \theta_D + \Delta\theta_D, \theta_G)}{\partial \theta_G} = \frac{\partial f(z; \theta'_D, \theta_G)}{\partial \theta_G} - 2\alpha \left( \frac{\partial^2 f(z; \theta_D, \theta_G)}{\partial \theta_D \partial \theta_G} \right)^T \frac{\partial f(z; \theta'_D, \theta_G)}{\partial (\theta'_D)} \quad (11)$$

The approximation in eq. 10 comes from using  $\frac{\partial f(z; \theta_D, \theta_G)}{\partial \theta_D} \approx \frac{\partial f(z; \theta'_D, \theta_G)}{\partial \theta'_D}$  and  $\frac{\partial f(z; \theta_D, \theta_G)}{\partial \theta_G} \approx \frac{\partial f(z; \theta'_D, \theta_G)}{\partial \theta_G}$  as a result of additional linear approximation.

At this point, unrolling  $D$  update only affects  $\theta_D$ . Although it is expensive to unroll both  $D$  and  $G$ , in principle, we can unroll  $G$  update and compute the gradient of  $\theta_D$  similarly using  $\Delta\theta_G = \alpha \frac{\partial f(z; \theta_D, \theta_G)}{\partial \theta_G}$ :

$$\frac{\partial f(z; \theta_D, \theta_G + \Delta\theta_G)}{\partial \theta_D} \approx \frac{\partial f(z; \theta_D, \theta_G)}{\partial \theta_D} + 2\alpha \left( \frac{\partial^2 f(z; \theta_D, \theta_G)}{\partial \theta_G \partial \theta_D} \right)^T \frac{\partial f(z; \theta_D, \theta_G)}{\partial \theta_G} \quad (12)$$

which gives us the same update rule as SGA (eq. 8). This correspondence based on first order Taylor expansion is unsurprising, as SGA is based on linearising the adversarial dynamics (Balduzzi et al., 2018).

Therefore, given the previous section, we can view LOGAN as further approximating Unrolled GAN, by unrolling the update of latent source  $z$  instead of the parameters. Although the information from  $z$  is limited compared with all the parameters, the intuition from Unrolled GANs applies here: unrolling the update of  $z$  gives  $D$  and  $G$  extra information to react to their opponents, thus avoiding the circular behaviour.

## 4 LOGAN WITH NATURAL GRADIENT DESCENT

Our analysis explains why latent optimisation may help GAN training. In practice, we expect more benefit from latent optimisation from *stronger* optimiser for  $z$ , which can better capture the coupling between  $D$  and  $G$ . Wu et al. (2019) only used basic gradient descent (GD) with a fixed step-size. This choice limits the size  $\Delta z$  can take: in order not to overshoot when the curvature is large, the step size would be too conservative when the curvature is small. We hypothesise that GD is more detrimental for larger models, which have complex loss surfaces with highly varying curvatures. Consistent with this hypothesis, we observed only marginal improvement over the baseline using GD (section 5.2.3, Table 1, Figure 3 b).

In this work, we propose using natural gradient descent (NGD, Amari 1998) for latent optimisation. NGD is an approximate second-order optimisation method, and has been applied successfully in many domains (Pascanu & Bengio, 2013; Martens, 2014). By using the positive semi-definite (PSD) Gauss-Newton matrix to approximate the (possibly negative definite) Hessian, NGD often works even better than exact second-order methods. NGD is expensive in high dimensional parameter spaces, even with approximations (Martens, 2014). However, we demonstrate that it is efficient for latent optimisation, even in very large models.

Given the gradient of  $z$ ,  $g = \frac{\partial f(z)}{\partial z}$ , NGD computes the update as

$$\Delta z = \alpha F^{-1} g \quad (13)$$

where the Fisher information matrix  $F$  is defined as

$$F = \mathbb{E}_{p(t|z)} [\nabla \ln p(t|z) \nabla \ln p(t|z)^T] \quad (14)$$

The log-likelihood function  $\ln p(t|z)$  typically corresponds to commonly used error functions such as the cross entropy loss. This correspondence is not necessary when we interpret NGD as an approximate second-order method, as has long been done (Martens, 2014). Nevertheless, Appendix D provides a Poisson log-likelihood interpretation for the hinge loss commonly used in GANs (Lim & Ye, 2017; Tran et al., 2017). An important difference between latent optimisation and commonly seen scenarios using NGD is that the expectation over the condition ( $z$ ) is absent. Since each  $z$  is only responsible for generating one image, it only minimises the loss  $L_G(z)$  for this particular instance.

More specifically, we use the *empirical* Fisher  $F'$  with Tikhonov damping, as in TONGA (Roux et al., 2008)

$$F' = g \cdot g^T + \beta I \quad (15)$$

$F'$  is cheaper to compute compared with the full Fisher, since  $g$  is already available. The *damping factor*  $\beta$  regularises the step size, which is important when  $F'$  only poorly approximates the Hessian or when the Hessian changes too much across the step. Using the Sherman-Morrison formula, the NGD update can be simplified into the following closed form:

$$\Delta z = \alpha \left( \frac{I}{\beta} - \frac{g g^T}{\beta^2 + \beta g^T g} \right) g = \frac{\alpha}{\beta + \|g\|^2} g \quad (16)$$

which does not involve any matrix inversion. Thus, NGD adapts the step size according to the curvature estimate  $c = \frac{1}{\beta + \|g\|^2}$ . When  $\beta$  is small, NGD normalises the gradient by its squared L2-norm. NGD automatically smooths the scale of updates by down-scaling the gradients as their norm grows, which also contributes to the smoothed norms of updates (Appendix A.2). Since the NGD update remains proportional to  $g$ , our analysis based on gradient descent in section 3 still holds.

### ADDITIONAL REGULARISATION

Various regularisation techniques are often necessary to ensure the stable training of GANs. Here we highlight two of them that we found particularly useful in combination with LOGAN. First, we found regularising the Euclidean norm of optimisation step,

$$R_z = w_r \cdot \|\Delta z\|_2^2 \quad (17)$$

where the scalar weight  $w_r$  is a parameter, as introduced by Wu et al. (2019) is necessary, especially for large models. This term is added to both the generator loss and discriminator loss in training.

Wu et al. (2019) suggested this term is related to optimal transport; more recently, Tanaka (2019) formalised this connection in Discriminator Optimal Transport (DOT). We left the exact connection between our work and DOT to future investigation, but here note that while DOT improves evaluation performance, our method mainly focuses on training. Consequently, although the regulariser  $R_z$  shares the same form as that in DOT, they function differently: it regularises the update of parameters here, but the latent code  $z$  in DOT.

In addition, we found it is more stable to optimise only a portion  $c$  of  $z$ , leaving some of its elements completely random, which can be seen as an additional damping mechanism while preserve more randomness from the latent source.

## 5 EXPERIMENTS AND ANALYSIS

We tested our algorithm for both medium (DCGAN, Radford et al. 2015; Miyato et al. 2018) and large scale (BigGAN, Brock et al. 2018) models. We use the standard hyper-parameter settings for each GAN model, without further optimising them with LOGAN. We performed grid-search over the four parameters introduced in LOGAN: the latent step size  $\alpha$ , damping factor  $\beta$ , the regularisation weight  $w_r$ , and the portion of  $z$  being optimised as  $c$ . Details of the grid search are summarised in Appendix E. Additional empirical analysis of latent optimisation is presented in Appendix B.

### 5.1 EXPERIMENTS WITH DCGAN ON CIFAR

To verify if our proposed NGD optimiser works well for latent optimisation, we first test LOGAN at more moderate scales for direct comparison with Wu et al. (2019) using basic GD. Here we apply latent optimisation on Spectral Normalised GANs (SN-GANs, Miyato et al. 2018).

The experiments follows the same basic setup and hyper-parameter settings as the CS-GAN in Wu et al. (2019). There is no class conditioning in this model. With NGD, we use a large step size of  $\alpha = 0.9$  and the damping factor  $\beta = 0.1$  for optimising  $z$ . We found the weight of 0.1 for the regulariser  $R_z$  (eq. 17), and optimising 80% of the latent source worked best for SN-GANs. All other parameters are same as in Wu et al. 2019.

In addition, we found running extra latent optimisation steps benefited evaluation, so we use ten steps of latent optimisation in evaluation for results in this section, although the models were still trained with a single optimisation step. This is different from in larger models, where optimisation is unnecessary in evaluation (see section 5.2.2 for more details).

Table 2 shows the FID and IS alongside SN-GAN and CS-CAN which used the same architecture. The scores are computed based on 10,000 samples following the same procedure as in Wu et al. (2019). We observe that NGD brought significant improvement over CS-GAN (i.e., LOGAN with GD for optimising  $z$ ). Compared with the baseline SN-GAN model without employing any latent optimisation, there is an improvement of 16.8% in IS and 39.6% in FID. Figure 4 compares random samples from these two models. Overall, samples from LOGAN (NGD) have higher contrasts and sharper contours.

Table 2: Comparison of Scores. The first and second columns are reproduced from Miyato et al. (2018) and Wu et al. (2019) respectively. We report the Inception Score (IS, higher is better, Salimans et al. 2016) and Fréchet Inception Distance (FID, lower is better, Heusel et al. 2017).

	SN-GAN	CS-GAN	LOGAN (NGD)
FID	29.3	$23.1 \pm 0.5$	<b><math>17.7 \pm 0.4</math></b>
IS	$7.42 \pm 0.08$	$7.80 \pm 0.05$	<b><math>8.67 \pm 0.05</math></b>

### 5.2 EXPERIMENTS WITH BIGGAN ON IMAGENET

To illustrate the scalability of our algorithm, we next focus on large scale models based on BigGAN-deep (Brock et al., 2018) trained on  $128 \times 128$  size images from the ImageNet dataset (Deng et al., 2009).



Figure 4: (a) Samples from SN-GAN. (b) Samples from LOGAN.

### 5.2.1 MODEL CONFIGURATION

We used the standard BigGAN-deep architecture with three minor modifications: 1. We increased the size of the latent source from 128 to 256, to compensate the randomness of the source lost when optimising  $z$ . 2. We use the uniform distribution  $\mathcal{U}(-1, 1)$  instead of the standard normal distribution  $\mathcal{N}(0, 1)$  for  $p(z)$  to be consistent with the clipping operation (Algorithm 1). 3. We use leaky ReLU (with the slope of 0.2 for the negative part) instead of ReLU as the non-linearity for smoother gradient flow for  $\frac{\partial f(z)}{\partial z}$ .

Consistent with the detailed findings in Brock et al. (2018), our experiment with this baseline model obtains only slightly better scores compared with those in Brock et al. (2018) (Table 1, see also Figure 12 in Appendix G). We computed the FID and IS as in Brock et al. (2018), and computed IS values from checkpoints with the lowest FIDs. Finally, we computed the means and standard deviations for both measures from 5 models with different random seeds.

To apply latent optimisation with NGD, we use the same large step size of  $\alpha = 0.9$  as in SN-GAN (section 5.1). However, we found much heavier damping is essential for BigGAN, so we use the damping factor  $\beta = 5.0$ , and only optimise 50% of  $z$ 's elements. Consistent with Tanaka (2019), we found a much larger weight of 300.0 for the regulariser  $R_z$  (eq. 17) works best, since deeper models generally have larger Lipschitz constants. All other hyper-parameters, including learning rates and a large batch size of 2048, remain the same as in BigGAN-deep. We call this model LOGAN (NGD).

### 5.2.2 BASIC RESULTS

Employing the same architecture and number of parameters as the BigGAN-deep baseline, LOGAN (NGD) achieved better FID and IS (Table 1). As observed by Brock et al. (2018), BigGAN training eventually collapsed in every experiment. Training with LOGAN also collapsed, perhaps due to higher-order dynamics beyond the scope we have analysed, but it took significantly longer (600k steps versus 300k steps with BigGAN-deep).

During training, LOGAN was about 3 times slower per step compared with BigGAN-deep because of the additional forward and backward passes. In contrast to experiments with smaller models (section 5.1), we found that optimising  $z$  during evaluation did not improve sample scores (even up to 10 steps), so we do not optimise  $z$  for evaluation. Therefore, LOGAN has the same evaluation cost as original BigGAN-deep. To help understand this behaviour, we plot the change from  $\Delta z$  during training in Figure 5 a. Although the movement in Euclidean space  $\|\Delta z\|$  grew until training collapsed, the movement in  $D$ 's output space, measured as  $\|f(z + \Delta z) - f(z)\|$ , remained unchanged (see Appendix F for details). As shown in our analysis, optimising  $z$  improves the training dynamics, so LOGANs work well after training without requiring latent optimisation. We reckon that smaller



models might not be “over-parametrised” enough to fully amortise the computation from optimising  $z$ , which can then further exploit the architecture in evaluation time. Appendix B further illustrates these different behaviours. We aim to further investigate this difference in future studies.

Given the criticism of FID and IS as heuristics metrics for sample distributions, we further measure how these samples directly contribute to downstream classification task via the recently proposed Classification Accuracy Score (CAS, Ravuri & Vinyals 2019). Unlike FID and IS, this metric favours likelihood-based models, which are more likely cover all modes representing different classes. The CAS from LOGAN nearly halved the gap between the state-of-the-art GANs and VQ-VAE2 (Razavi et al., 2019). See Appendix C for more details.

### 5.2.3 ABLATION STUDIES

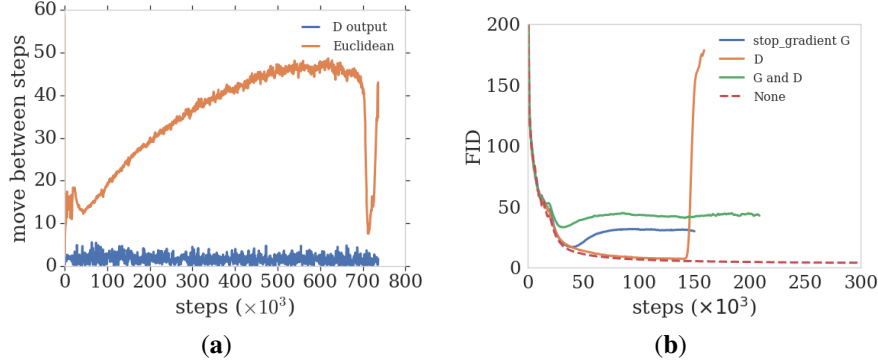


Figure 5: **(a)** The change from  $\Delta z$  across training, in  $D$ ’s output space and  $z$ ’s Euclidean space. The distances are normalised by their standard derivations computed from a moving window of size 20 (1007 data points in total). **(b)** Training curves from models with different “stop\_gradient” operations. For reference, the training curve from an unablated model is plotted as the dashed line. All instances with stop\_gradient collapsed (FID went up) early in training.

We verify our theoretical analysis in section 3 by examining key components of Algorithm 1 via ablation studies. First, we experiment with using basic GD to optimising  $z$ , as in Wu et al. (2019), and call this model LOGAN (GD). A smaller step size of  $\alpha = 0.0001$  was required; larger values were unstable and led to premature collapse of training. As shown in Table 1, the scores from LOGAN (GD) were worse than LOGAN (NGD) and similar to the baseline model.

We then evaluate the effects of removing those terms depending on  $\frac{\partial f(z)}{\partial z}$  in eq. 6, which are not in the ordinary gradient (eq. 3). Since we computed these terms by back-propagating through the latent optimisation procedure, we removed them by selectively blocking back-propagation with “stop\_gradient” operations (e.g., in TensorFlow, Abadi et al. 2016). Figure 5 b shows the change of FIDs for the three models corresponding to removing  $\left(\frac{\partial \Delta z}{\partial \theta_G}\right)^T \frac{\partial f(z')}{\partial z'}$ , removing  $\left(\frac{\partial \Delta z}{\partial \theta_D}\right)^T \frac{\partial f(z')}{\partial z'}$  and removing both terms. As predicted by our analysis (section 3), both terms help stabilise training; training diverged early for all three ablations.

### 5.2.4 TRUNCATION AND SAMPLES

Truncation is a technique introduced by Brock et al. (2018) to illustrate the trade-off between the FID and IS in a trained model. For a model trained with  $z \sim p(z)$  from a source distribution symmetric around 0, such as the standard normal distribution  $\mathcal{N}(0, 1)$  and the uniform distribution  $\mathcal{U}(-1, 1)$ , down-scaling (truncating) the source  $\tilde{z} = s \cdot z$  with  $0 \leq s < 1$  gives samples with higher visual quality but reduced diversity. We see this quantified in higher IS scores and lower FID when evaluating samples from truncated distributions.

Figure 3 b plots the truncation curves for the baseline BigGAN-deep model, LOGAN (GD) and LOGAN (NGD), obtained by varying the truncation (value of  $s$ ) from 1.0 (no truncation, upper-left

ends of the curves) to 0.02 (extreme truncation, bottom-right ends). Each curve shows the trade-off between FID and IS for an individual model; curves towards the upper-right corner indicate better overall sample quality. The relative positions of curves in figure 3 (b) shows LOGAN (NGD) has the best sample quality. Interestingly, although LOGAN (GD) and the baseline model have similar scores without truncation (upper-left ends of the curves, see also Table 1), LOGAN (GD) was better behaved with increasing truncation, suggesting LOGAN (GD) still converged to a better equilibrium. For further reference, we plot truncation curves from additional baseline models in Figure 12 (Appendix G).

Figure 1 and Figure 2 show samples from selected points along the truncation curves. In the high IS regime, C and D on the truncation curves both have similarly high IS of near 260. Samples from batches with such high IS have almost photo-realistic image quality. Figure 1 shows that while the baseline model produced nearly uniform samples, LOGAN (NGD) could still generate highly diverse samples. On the other hand, A and B from Figure 3 b have similarly low FID of near 5, indicating high sample diversity. Samples in Figure 2 b show higher quality compared with those in a (e.g., the interfaces between the elephants and ground, the contours around the pandas).

## 6 CONCLUSION

In this work, we present LOGAN, which significantly improves the state of the art in large scale GAN training for image generation by optimising the latent source  $z$ . Our results illustrate improvements in quantitative evaluation and samples with higher quality and diversity. Moreover, our analysis suggests that LOGAN fundamentally improves adversarial training dynamics. LOGAN is related to the energy-based formulation of a GAN’s discriminator (Dai et al., 2017; Kumar et al., 2019; Du & Mordatch, 2019), when latent optimisation is viewed as descending the energy function defined by the discriminator. From this view, sampling from the distribution implicitly defined by this energy function, via, e.g., Langevin sampling (Welling & Teh, 2011), may bring further benefits. Another class of approaches regularises the entropy of the generator outputs to reduce mode collapse (Belghazi et al., 2018; Dieng et al., 2019). Such techniques could be combined with LOGAN to further improve coverage of the underlying data distribution. Moreover, we expect our method to be useful in other tasks that involve adversarial training, including representation learning and inference (Donahue et al., 2017; Dumoulin et al., 2017; Donahue & Simonyan, 2019), text generation (Zhang et al., 2019), style learning (Zhu et al., 2017; Karras et al., 2019), audio generation (Donahue et al., 2018) and video generation (Vondrick et al., 2016; Clark et al., 2019).

## ACKNOWLEDGMENTS

We thank Mihaela Rosca, Suman Ravuri and James Martens for comments on the draft and insightful discussions.

## REFERENCES

- Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. Tensorflow: A system for large-scale machine learning. In *12th USENIX Symposium on Operating Systems Design and Implementation OSDI*, pp. 265–283, 2016.
- Shun-Ichi Amari. Natural gradient works efficiently in learning. *Neural computation*, 10(2):251–276, 1998.
- Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein GAN. *arXiv preprint arXiv:1701.07875*, 2017.
- Sanjeev Arora, Rong Ge, Yingyu Liang, Tengyu Ma, and Yi Zhang. Generalization and equilibrium in generative adversarial nets. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 224–232. JMLR. org, 2017.
- David Balduzzi, Sebastien Racaniere, James Martens, Jakob Foerster, Karl Tuyls, and Thore Graepel. The mechanics of n-player differentiable games. *arXiv preprint arXiv:1802.05642*, 2018.

- 
- Mohamed Ishmael Belghazi, Aristide Baratin, Sai Rajeswar, Sherjil Ozair, Yoshua Bengio, Aaron Courville, and R Devon Hjelm. MINE: Mutual information neural estimation. In *International Conference on Machine Learning*, 2018.
- Vivek S Borkar. Stochastic approximation with two time scales. *Systems & Control Letters*, 29(5): 291–294, 1997.
- Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale GAN training for high fidelity natural image synthesis. *arXiv preprint arXiv:1809.11096*, 2018.
- Emmanuel J Candes, Justin K Romberg, and Terence Tao. Stable signal recovery from incomplete and inaccurate measurements. *Communications on Pure and Applied Mathematics: A Journal Issued by the Courant Institute of Mathematical Sciences*, 59(8):1207–1223, 2006.
- Aidan Clark, Jeff Donahue, and Karen Simonyan. Efficient video generation on complex datasets. *arXiv preprint arXiv:1907.06571*, 2019.
- Zihang Dai, Amjad Almahairi, Philip Bachman, Eduard Hovy, and Aaron Courville. Calibrating energy-based generative adversarial networks. *ICLR*, 2017.
- J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR*, 2009.
- Adji B. Dieng, Francisco J. R. Ruiz, David M. Blei, and Michalis K. Titsias. Prescribed generative adversarial networks. *arXiv preprint arXiv:1910.04302*, 2019.
- Chris Donahue, Julian McAuley, and Miller Puckette. Adversarial audio synthesis. *arXiv preprint arXiv:1802.04208*, 2018.
- Jeff Donahue and Karen Simonyan. Large scale adversarial representation learning. In *NeurIPS*, 2019.
- Jeff Donahue, Philipp Krähenbühl, and Trevor Darrell. Adversarial feature learning. In *ICLR*, 2017.
- David L Donoho. Compressed sensing. *IEEE Transactions on information theory*, 52(4):1289–1306, 2006.
- Yilun Du and Igor Mordatch. Implicit generation and generalization in energy-based models. *Advances in Neural Information Processing Systems*, 2019.
- Vincent Dumoulin, Ishmael Belghazi, Ben Poole, Olivier Mastropietro, Alex Lamb, Martin Arjovsky, and Aaron Courville. Adversarially learned inference. In *ICLR*, 2017.
- Ian Gemp and Sridhar Mahadevan. Global Convergence to the Equilibrium of GANs using Variational Inequalities. In *Arxiv:1808.01531*, 2018.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pp. 2672–2680, 2014.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Advances in Neural Information Processing Systems*, pp. 6626–6637, 2017.
- Morris W Hirsch. Convergent activation dynamics in continuous time networks. *Neural networks*, 2(5):331–349, 1989.
- Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4401–4410, 2019.

- 
- Vijaymohan R Konda and Vivek S Borkar. Actor-critic-type learning algorithms for markov decision processes. *SIAM Journal on control and Optimization*, 38(1):94–123, 1999.
- Rithesh Kumar, Sherjil Ozair, Anirudh Goyal, Aaron Courville, and Yoshua Bengio. Maximum entropy generators for energy-based models. *arXiv preprint arXiv:1901.08508*, 2019.
- Alistair Letcher, David Balduzzi, Sébastien Racanière, James Martens, Jakob N Foerster, Karl Tuyls, and Thore Graepel. Differentiable game mechanics. *Journal of Machine Learning Research*, 20(84):1–40, 2019.
- Jae Hyun Lim and Jong Chul Ye. Geometric GAN. *arXiv preprint arXiv:1705.02894*, 2017.
- James Martens. New insights and perspectives on the natural gradient method. *arXiv preprint arXiv:1412.1193*, 2014.
- Lars Mescheder, Sebastian Nowozin, and Andreas Geiger. The numerics of gans. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems 30*, pp. 1825–1835. Curran Associates, Inc., 2017. URL <http://papers.nips.cc/paper/6779-the-numerics-of-gans.pdf>.
- Lars Mescheder, Andreas Geiger, and Sebastian Nowozin. Which training methods for gans do actually converge? In *International Conference on Machine Learning*, pp. 3478–3487, 2018.
- Luke Metz, Ben Poole, David Pfau, and Jascha Sohl-Dickstein. Unrolled generative adversarial networks. *CoRR*, abs/1611.02163, 2016. URL <http://arxiv.org/abs/1611.02163>.
- Takeru Miyato and Masanori Koyama. cGANs with projection discriminator. *arXiv preprint arXiv:1802.05637*, 2018.
- Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative adversarial networks. *arXiv preprint arXiv:1802.05957*, 2018.
- Sebastian Nowozin, Botond Cseke, and Ryota Tomioka.  $f$ -GAN: Training generative neural samplers using variational divergence minimization. In *Advances in neural information processing systems*, pp. 271–279, 2016.
- Razvan Pascanu and Yoshua Bengio. Revisiting natural gradient for deep networks. *arXiv preprint arXiv:1301.3584*, 2013.
- Barak A Pearlmutter. Fast exact multiplication by the hessian. *Neural computation*, 6(1):147–160, 1994.
- Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.
- Suman Ravuri and Oriol Vinyals. Classification accuracy score for conditional generative models. In *Advances in Neural Information Processing Systems*, pp. 12268–12279, 2019.
- Ali Razavi, Aaron van den Oord, and Oriol Vinyals. Generating diverse high-fidelity images with vq-vae-2. In *Advances in Neural Information Processing Systems*, pp. 14866–14876, 2019.
- Nicolas L Roux, Pierre-Antoine Manzagol, and Yoshua Bengio. Topmoumoute online natural gradient algorithm. In *Advances in neural information processing systems*, pp. 849–856, 2008.
- Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. In *Advances in neural information processing systems*, pp. 2234–2242, 2016.
- Andrew M Saxe, James L McClelland, and Surya Ganguli. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. *ICLR*, 2014.
- Akinori Tanaka. Discriminator optimal transport. In *Advances in Neural Information Processing Systems*, pp. 6813–6823, 2019.

Dustin Tran, Rajesh Ranganath, and David Blei. Hierarchical implicit models and likelihood-free variational inference. In *Advances in Neural Information Processing Systems*, pp. 5523–5533, 2017.

Carl Vondrick, Hamed Pirsiavash, and Antonio Torralba. Generating videos with scene dynamics. In *Advances In Neural Information Processing Systems*, pp. 613–621, 2016.

Max Welling and Yee W Teh. Bayesian learning via stochastic gradient langevin dynamics. In *Proceedings of the 28th international conference on machine learning (ICML-11)*, pp. 681–688, 2011.

Yan Wu, Mihaela Rosca, and Timothy Lillicrap. Deep compressed sensing. In *International Conference on Machine Learning*, pp. 6850–6860, 2019.

Han Zhang, Ian Goodfellow, Dimitris Metaxas, and Augustus Odena. Self-attention generative adversarial networks. In *International Conference on Machine Learning*, pp. 7354–7363, 2019.

Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pp. 2223–2232, 2017.

## A DETAILED ANALYSIS OF LATENT OPTIMISATION

In this section we present complementary of LOGAN. In particular, we show how the algorithm brings together ideas from symplectic gradient adjustment and stochastic approximation with two time scales.

### A.1 APPROXIMATE SYMPLECTIC GRADIENT ADJUSTMENT

To analyse LOGAN as a differentiable game we treat the latent step  $\Delta z$  as adding a *third player* to the original game played by the discriminator and generator. The third player’s parameter,  $\Delta z$ , is optimised online for each  $z \sim p(z)$ . Together the three players (latent player, discriminator, and generator) have losses *averaged over a batch of samples*:

$$L = [\eta L_G, L_D, L_G]^T \quad (18)$$

where  $\eta = \frac{1}{N}$  ( $N$  is the batch size) reflects the fact that each  $\Delta z$  is only optimised for a single sample  $z$ , so its contribution to the total loss across a batch is small compared with  $\theta_D$  and  $\theta_G$  which are directly optimised for batch losses. This choice of  $\eta$  is essential for the following derivation, and has important practical implication. It means that the per-sample loss  $L_G(z')$ , instead of the loss summed over a batch  $\sum_{n=1}^N L_G(z'_n)$ , should be the only loss function guiding latent optimisation. Therefore, when using natural gradient descent (Section 4), the Fisher information matrix should only be computed using the current sample  $z$ .

The resulting simultaneous gradient is

$$\begin{aligned} g &= \left[ \eta \frac{\partial L_G(z')}{\partial \Delta z}, \frac{\partial L_D(z')}{\partial \theta_D}, \frac{\partial L_G(z')}{\partial \theta_G} \right]^T \\ &= \left[ -\eta \frac{\partial f(z')}{\partial \Delta z}, \frac{\partial f(z')}{\partial \theta_D}, -\frac{\partial f(z')}{\partial \theta_G} \right]^T \end{aligned} \quad (19)$$

Following Balduzzi et al. (2018), we can write the Hessian of the game as:

$$H = \begin{bmatrix} -\eta \frac{\partial^2 f(z')}{\partial \Delta z^2} & -\eta \frac{\partial^2 f(z')}{\partial \Delta z \partial \theta_D} & -\eta \frac{\partial^2 f(z')}{\partial \Delta z \partial \theta_G} \\ \frac{\partial^2 f(z')}{\partial \theta_D \partial \Delta z} & \frac{\partial^2 f(z')}{\partial \theta_D^2} & \frac{\partial^2 f(z')}{\partial \theta_D \partial \theta_G} \\ -\frac{\partial^2 f(z')}{\partial \theta_G \partial \Delta z} & -\frac{\partial^2 f(z')}{\partial \theta_G \partial \theta_D} & -\frac{\partial^2 f(z')}{\partial \theta_G^2} \end{bmatrix} \quad (20)$$



The presence of a non-zero anti-symmetric component in the Hessian

$$A = \frac{1}{2}(H - H^T)$$

$$= \begin{bmatrix} 0 & -\frac{1+\eta}{2} \frac{\partial^2 f(z')}{\partial \Delta z \partial \theta_D} & \frac{1-\eta}{2} \frac{\partial^2 f(z')}{\partial \Delta z \partial \theta_G} \\ \frac{1+\eta}{2} \frac{\partial^2 f(z')}{\partial \theta_D \partial \Delta z} & 0 & \frac{\partial^2 f(z')}{\partial \theta_D \partial \theta_G} \\ -\frac{1-\eta}{2} \frac{\partial^2 f(z')}{\partial \theta_G \partial \Delta z} & -\frac{\partial^2 f(z')}{\partial \theta_G \partial \theta_D} & 0 \end{bmatrix} \quad (21)$$

implies the dynamics have a rotational component which can cause cycling or slow down convergence. Since  $\eta \ll 1$  for typical batch sizes (e.g.,  $\frac{1}{64}$  for DCGAN and  $\frac{1}{2048}$  for BigGAN-deep), we abbreviate  $\gamma = \frac{1+\eta}{2} \approx \frac{1-\eta}{2}$  to simplify notations.

Symplectic gradient adjustment (SGA) counteracts the rotational force by adding an adjustment term to the gradient to obtain  $g^* \leftarrow g + \lambda A^T g$ , which for the discriminator and generator has the form:

$$g_D^* = \frac{\partial f(z')}{\partial \theta_D} + \lambda \gamma \left( \frac{\partial^2 f(z')}{\partial \Delta z \partial \theta_D} \right)^T \frac{\partial f(z')}{\partial \Delta z} + \lambda \left( \frac{\partial^2 f(z')}{\partial \theta_G \partial \theta_D} \right)^T \frac{\partial f(z')}{\partial \theta_G} \quad (22)$$

$$g_G^* = -\frac{\partial f(z')}{\partial \theta_G} - \lambda \gamma \left( \frac{\partial^2 f(z')}{\partial \Delta z \partial \theta_G} \right)^T \frac{\partial f(z')}{\partial \Delta z} + \lambda \left( \frac{\partial^2 f(z')}{\partial \theta_D \partial \theta_G} \right)^T \frac{\partial f(z')}{\partial \theta_D} \quad (23)$$

The gradient with respect to  $\Delta z$  is ignored since the convergence of training only depends on  $\theta_D$  and  $\theta_G$ .

If we drop the last terms in eq.22 and 23, which are expensive to compute for large models with high-dimensional  $\theta_D$  and  $\theta_G$ , and use  $\frac{\partial f(z')}{\partial \Delta z} = \frac{\partial f(z')}{\partial z'}$ , the adjusted updates can be rewritten as

$$g_D^* \approx \frac{\partial f(z')}{\partial \theta_D} + \lambda \gamma \left( \frac{\partial^2 f(z')}{\partial z' \partial \theta_D} \right)^T \frac{\partial f(z')}{\partial z'} \quad (24)$$

$$g_G^* \approx -\frac{\partial f(z')}{\partial \theta_G} - \lambda \gamma \left( \frac{\partial^2 f(z')}{\partial z' \partial \theta_G} \right)^T \frac{\partial f(z')}{\partial z'} \quad (25)$$

Because of the third player, there are still the terms depend on  $\frac{\partial f(z')}{\partial z'}$  to adjust the gradients. Efficiently computing  $\frac{\partial^2 f(z')}{\partial z' \partial \theta_D}$  and  $\frac{\partial^2 f(z')}{\partial z' \partial \theta_G}$  is non-trivial (e.g., Pearlmutter 1994). However, if we introduce the local approximation

$$\frac{\partial^2 f(z')}{\partial z' \partial \theta_D} \approx \frac{\partial^2 f(z)}{\partial z \partial \theta_D} \quad \frac{\partial^2 f(z')}{\partial z' \partial \theta_G} \approx \frac{\partial^2 f(z)}{\partial z \partial \theta_G} \quad (26)$$

then the adjusted gradient becomes identical to eq. 6 from latent optimisation.

In other words, automatic differentiation by commonly used machine learning packages can compute the adjusted gradient for  $\theta_D$  and  $\theta_G$  when back-propagating through the latent optimisation process. Despite the approximation involved in this analysis, both our experiments in section 5 and the results from Wu et al. (2019) verified that latent optimisation can significantly improve GAN training.

## A.2 STOCHASTIC APPROXIMATION WITH TWO TIME SCALES

This section shows that latent optimisation accelerates the speed of updating  $D$  relative to the speed of updating  $G$ , facilitating convergence according to Heusel et al. (2017) (see also Figure 6 b). Intuitively, the generator requires less updating compared with  $D$  to achieve the same reduction of loss because latent optimisation “helps”  $G$ .

Heusel et al. (2017) used the theory of stochastic approximation to analyse GAN training. Viewing the training process as stochastic approximation with two time scales (Borkar, 1997; Konda & Borkar, 1999), they suggest that the update of  $D$  should be fast enough compared with that of  $G$ . Under mild assumptions, Heusel et al. (2017) proved that such two time-scale update converges to local Nash equilibrium. Their analysis follows the idea of  $(\tau, \delta)$  perturbation (Hirsch, 1989), where

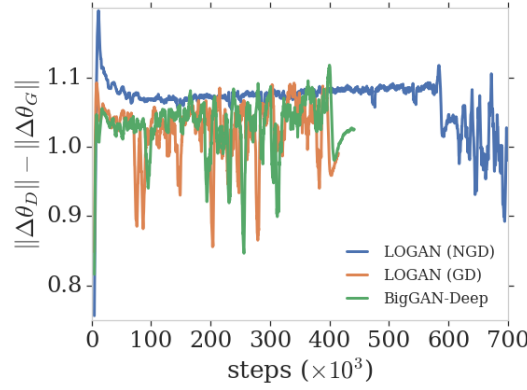


Figure 6: The update speed of the discriminator relative to the generator shown as the difference  $\|\Delta\theta_D\| - \|\Delta\theta_G\|$  after each update step. Lines are smoothed with a moving average using window size 20 (in total, there are 3007, 1659 and 1768 data points for each curve). All curves oscillated strongly after training collapsed.

the slow updates ( $G$ ) are interpreted as a small perturbation over the ODE describing the fast update ( $D$ ). Importantly, the size of perturbation  $\delta$  is measured by the magnitude of parameter change, which is affected by both the learning rate and gradients.

Here we show, in accordance with Heusel et al. (2017), that LOGAN accelerates discriminator updates and slows down generator updates, thus helping the convergence of discriminator. We start from analysing the change of  $\theta_G$ . We assume that, without LO, it takes  $\Delta\theta_G = \theta'_G - \theta_G$  to make a small constant amount of reduction in loss  $L_G$ :

$$\rho = -f(z; \theta_D, \theta_G + \Delta\theta_G) + f(z; \theta_D, \theta_G) \quad (27)$$

Now using the optimised  $z' = z + \Delta z$ , we assess the change  $\delta\theta_G$  required to achieve the same amount of reduction:

$$\rho = -f(z + \Delta z; \theta_D, \theta_G + \delta\theta_G) + f(z; \theta_D, \theta_G) \quad (28)$$

Intuitively, when  $z$  “helps”  $\theta_G$  to achieve the same goal of increasing  $f(z; \theta_D, \theta_G)$  by  $\rho$ , the responsible of  $\theta_G$  becomes smaller, so it does not need to change as much as  $\Delta\theta_G$ , thus  $\|\delta\theta_G\| < \|\Delta\theta_G\|$ .

Formally,  $f(z; \theta_D, \theta_G)$  and  $f(z + \Delta z; \theta_D, \theta_G + \delta\theta_G)$  have the following Taylor expansions around  $z$  and  $\theta_G$ :

$$f(z; \theta_D, \theta_G + \delta\theta_G) = f(z; \theta_D, \theta_G) + \left( \frac{\partial f(z; \theta_D, \theta_G)}{\partial \theta_G} \right)^T \Delta\theta_G + \epsilon(\Delta\theta_G) \quad (29)$$

$$f(z + \Delta z; \theta_D, \theta_G + \delta\theta_G) = f(z; \theta_D, \theta_G) + \left( \frac{\partial f(z; \theta_D, \theta_G)}{\partial z} \right)^T \Delta z \quad (30)$$

$$+ \left( \frac{\partial f(z + \Delta z; \theta_D, \theta_G)}{\partial \theta_G} \right)^T \delta\theta_G + \epsilon(\Delta z, \delta\theta_G) \quad (31)$$

Where  $\epsilon(\cdot)$ ’s are higher order terms of the increments. Using the assumption of eq. 27 and 28, we can combine eq. 29 and 31:

$$\left( \frac{\partial f(z; \theta_D, \theta_G)}{\partial \theta_G} \right)^T \Delta\theta_G = \left( \frac{\partial f(z; \theta_D, \theta_G)}{\partial z} \right)^T \Delta z + \left( \frac{\partial f(z + \Delta z; \theta_D, \theta_G)}{\partial \theta_G} \right)^T \delta\theta_G + \epsilon \quad (32)$$

where  $\epsilon = \epsilon(\Delta z, \delta\theta_G) - \epsilon(\Delta\theta_G)$ . Since  $\Delta z \propto \frac{\partial f(z; \theta_D, \theta_G)}{\partial z}$  in gradient descent (eq. 3),

$$\frac{\partial f(z; \theta_D, \theta_G)}{\partial z} \Delta z > 0 \quad (33)$$

Therefore, we have the inequality

$$\left( \frac{\partial f(z; \theta_D, \theta_G)}{\partial \theta_G} \right)^T \Delta\theta_G < \left( \frac{\partial f(z + \Delta z; \theta_D, \theta_G)}{\partial \theta_G} \right)^T \delta\theta_G + \epsilon \quad (34)$$

If we further assume  $\Delta\theta_G$  and  $\delta\theta_G$  are obtained from stochastic gradient descent with identical learning rate,

$$\Delta\theta_G = \alpha \frac{\partial f(z; \theta_D, \theta_G)}{\partial \theta_G} \quad \delta\theta_G = \alpha \frac{\partial f(z; \theta_D, \theta_G)}{\partial \theta_G} \quad (35)$$

substituting eq. 35 into eq. 34 gives

$$\|\Delta\theta_G\| < \|\delta\theta_G\| + \epsilon \quad (36)$$

The same analysis applies to the discriminator. The similar intuition is that it takes the discriminator additional effort to compensate the exploitation from the optimised  $z'$ . We then obtain

$$\left( \frac{\partial f(z; \theta_D, \theta_G)}{\partial \theta_D} \right)^T \Delta\theta_D = \left( \frac{\partial f(z; \theta_D, \theta_G)}{\partial z} \right)^T \Delta z + \left( \frac{\partial f(z + \Delta z; \theta_D, \theta_G)}{\partial \theta_D} \right)^T \delta\theta_D + \epsilon \quad (37)$$

However, since the adversarial loss  $L_D = -L_G$ , we have  $\Delta\theta_D = -\alpha \frac{\partial f(z; \theta_D, \theta_G)}{\partial \theta_D}$  and  $\delta\theta_D = -\alpha \frac{\partial f(z; \theta_D, \theta_G)}{\partial \theta_D}$  taking the opposite signs of eq.35. For sufficiently small  $\Delta z$ ,  $\Delta\theta_G$  and  $\delta\theta_G$ ,  $\epsilon$  is close to zero, so  $\|\Delta\theta_D\| < \|\delta\theta_D\|$  under our assumptions of small  $\Delta z$ ,  $\Delta\theta_G$  and  $\delta\theta_G$ .

Importantly, the bigger the product  $\frac{\partial f(z)}{\partial z} \Delta z$  is, the more robust the inequality is to the error from  $\epsilon$ . Moreover, bigger steps increase the speed gap between updating D and G, further facilitating convergence (in accordance with Heusel et al. (2017)). Overall, our analysis suggests:

1. More than one gradient descent step may not be helpful, since  $\Delta z$  from multiple GD steps may deviate from the direction of  $\frac{\partial f(z)}{\partial z}$ .
2. A large step of  $\Delta z$  is helpful in facilitating convergence by widening the gap between D and G updates (Heusel et al., 2017).
3. However, the step of  $\Delta z$  cannot be too large. In addition to the linear approximation we used throughout our analysis, the approximate SGA breaks down when eq.26 is strongly violated when ‘‘overshoot’’ brings the gradients at  $\frac{\partial f(z')}{\partial z'}$  to the opposite sign of  $\frac{\partial f(z)}{\partial z}$ .

## B ADDITIONAL ANALYSIS OF LATENT OPTIMISATION

Here we analyse the relationship between the number of latent optimisation steps during evaluation and the final FIDs and inception scores. As in the main paper, we train the SN-GAN model with only 1 latent optimisation step, but test them with  $\{0, 1, 5, 10, 20, 30\}$  steps during evaluation. For lower variance in computing the scores, we use 10,000 samples for evaluation (as oppose to 5000 samples used in the main paper for direct comparison with other baselines). The inception scores are taken from checkpoints with the best (lowest) FIDs, and the error bars indicate standard deviations obtained from 3 different random seeds.

Figure 7 shows that the scores can be substantially improved with extra optimisation steps in evaluation. Although only 1 step was used in training, up to around 20 steps at evaluation could still improve sample quality. Beyond that, the return from extra computation became diminishing.

We did not observe similar improvement with BigGANs in evaluation. To contrast the difference between them, Figure 8 and 9 illustrate the change of samples made by latent optimisation. In both cases, 10 latent optimisation steps were applied in evaluation, after the models were trained with 1 optimisation step. While the effect of improvement is clear in Figure 8 (from SN-GAN), the changes in Figure 9 are barely observable from inspecting the samples alone before and after latent optimisation.

## C EVALUATING CLASSIFICATION ACCURACY SCORE

To compute CAS, a ResNet classifier (He et al., 2016) is trained on samples from our model instead of the ImageNet dataset. Each data sample is replaced by a sample from the model conditioned on the same class. The trained model is then evaluated on the ImageNet dataset as in a standard classification task. We use the same schedule as in Ravuri & Vinyals (2019) for training the ResNet classifier, but stopped earlier at about 10k steps, where the classification accuracy peaked. See

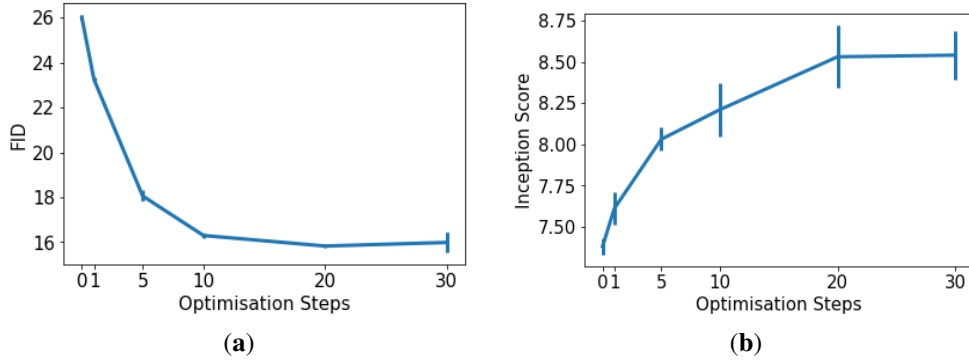


Figure 7: FIDs (a) and Inception Scores (b) obtained with different latent optimisation steps *at evaluation*.

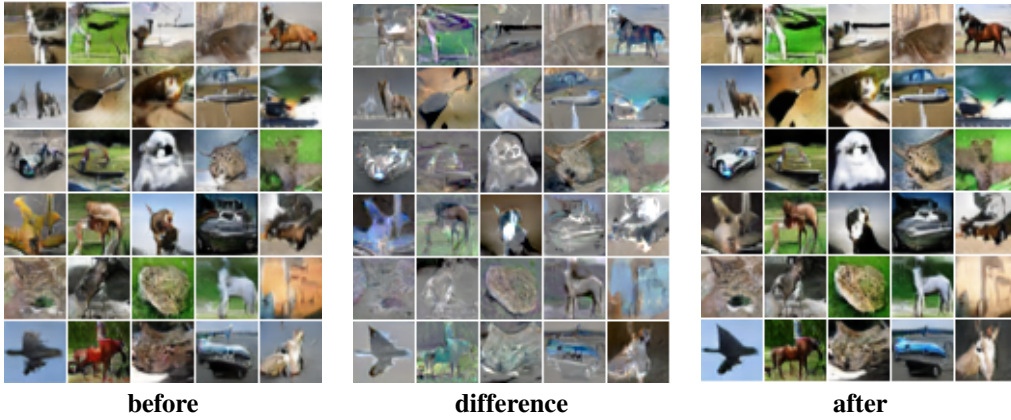


Figure 8: CIFAR samples from SN-GAN before, after latent optimisation, and the differences between them.

Ravuri & Vinyals (2019) for the motivation of CAS and more details of the training and evaluation procedure. We report both top-5 and top-1 classification accuracy in Table 3. Although higher resolution generally brings better CAS, we use the  $128 \times 128$  model as in the main paper due to limited computational resource. Despite this, the CAS from LOGAN nearly halved the gaps between BigGAN-deep and VQ-VAE2 at a higher resolution of  $256 \times 256$ .

Table 3: CAS for different models. Except LOGAN, numbers from all other models are reproduced from Razavi et al. (2019).

	Top-5 Accuracy	Top-1 Accuracy
BigGAN-deep ( $128 \times 128$ )	64.44%	40.64%
BigGAN-deep ( $256 \times 256$ )	65.92%	42.65%
LOGAN ( $128 \times 128$ )	71.97%	47.86%
VQ-VAE2 ( $256 \times 256$ )	77.59%	54.83%
Real Data ( $256 \times 256$ )	88.79%	68.82%

## D POISSON LIKELIHOOD FROM HINGE LOSS

Here we provide a probabilistic interpretation of the hinge loss for the generator, which leads naturally to the scenario of a family of discriminators. Although this interpretation is not necessary for our current algorithm, it may provides useful guidance for incorporating multiple discriminators.

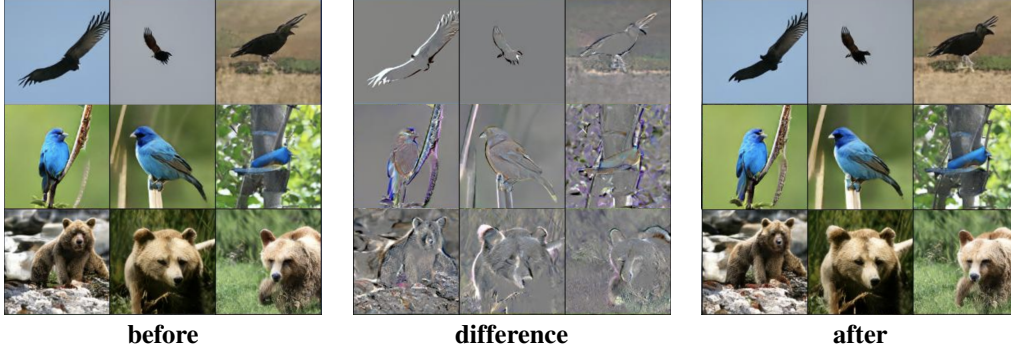


Figure 9: ImageNet samples from BigGAN-deep before, after latent optimisation, and the differences between them.

We introduce the label  $t = 1$  for real data and  $t = 0$  fake samples. This section shows that the generator hinge loss

$$L_G = -D(G(z)) \quad (38)$$

can be interpreted as a negative log-likelihood function:

$$L_G = -\ln p(t = 1; D, G(z)) \quad (39)$$

Here  $p(t = 1; z, D, G)$  is the probability that the generated image  $G(z)$  can fool the discriminator  $D$ .

The original GAN’s discriminator can be interpreted as outputting a Bernoulli distribution  $p(t; \beta_G) = \beta_G^t \cdot (1 - \beta_G)^{1-t}$ . In this case, if we parameterise  $\beta_G = D(G(z))$ , the generator loss is the negative log-likelihood

$$\begin{aligned} -\ln P(t = 1; D, G(z)) &= -\ln p(t = 1; \beta_G) \\ &= -\ln \beta_G = -\ln D(G(z)) \end{aligned} \quad (40)$$

Bernoulli, however, is not the only valid choice as the discriminator’s output distribution. Instead of sampling “1” or “0”, we assume that there are *many* identical discriminators that can independently vote to reject an input sample as fake. The number of votes  $k$  in a given interval can be described by a Poisson distribution with parameter  $\lambda$  with the following PMF:

$$p(k; \lambda) = \frac{\lambda^k e^{-\lambda}}{k!} \quad (41)$$

The probability that a generated image can fool *all* the discriminators is the probability of  $G(z)$  receiving no vote for rejection

$$p(k = 0; \lambda) = e^{-\lambda} \quad (42)$$

Therefore, we have the following negative log-likelihood as the generator loss if we parameterise  $\lambda = -D(G(z))$ :

$$-\ln p(k = 0; D, G(z)) = -\ln p(k = 0; \lambda) = -D(G(z)) \quad (43)$$

This interpretation has a caveat that when  $D(G(z)) > 0$  the Poisson distribution is not well defined. However, in general the discriminator’s hinge loss

$$L_D = -\min(0, -1 + D(x)) - \min(0, -1 - D(G(z))) \quad (44)$$

pushes  $D(G(z)) < 0$  via training.

## E HYPER-PARAMETER SEARCH

We first searched the hyper-parameters for the DCGAN (section 5.1) over the following range:

Base on the results from DCGANs, hyper-parameter search on the following grid was performed for the BigGAN-deep (section 5.2):



Table 4: Hyper-parameter grid for the DCGAN. Best values from the grid search are highlighted.

	values
$\alpha$	{0.01, 0.1, 0.3, 0.5, 0.7, 0.8, <b>0.9</b> , 1.0}
$\beta$	{0.01, <b>0.1</b> , 0.3, 0.5, 0.7, 0.9, 1.0}
$w_r$	{0.01, <b>0.1</b> , 0.2, 0.3, 0.5, 0.7, 1.0, 2.0, 3.0}
$c$	{30%, 50%, 70%, <b>80%</b> , 90%, 100%}

Table 5: Hyper-parameter grid for the BigGAN-deep. Best values from the grid search are highlighted.

	values
$\alpha$	{0.7, 0.8, <b>0.9</b> , 1.0}
$\beta$	{0.1, 0.5, 1.0, 2.0, 3.0, 4.0, <b>5.0</b> , 7.0, 9.0}
$w_r$	{0.1, 0.2, 0.5, 1.0, 2.0, 5.0, 10.0, 30.0, 100.0, 200.0, <b>300.0</b> , 400.0, 500.0}
$c$	{30%, <b>50%</b> , 70%, 80%}

## F DETAILS IN COMPUTING DISTANCES IN FIGURE 5 A

For a temporal sequence  $x_1, x_2, \dots, x_T$  (e.g. the changes of  $z$  or  $f(z)$  at each training step in this paper), to normalise its variance while accounting for the non-stationarity, we process it as follows. We first compute the moving average and standard deviation over a window of size  $N$ :

$$\mu_t = \frac{1}{N} \sum_{u=t}^{t+N-1} x_u \quad (45)$$

$$\sigma_t = \sqrt{\frac{1}{N-1} \sum_{u=t}^{t+N-1} (x_u - \mu_u)^2} \quad (46)$$

Then normalise the sequence as:

$$\bar{x}_t = \frac{x_t}{\sigma_t} \quad (47)$$

The result in Figure 5a is robust to the choice of window size. Our experiments with  $N$  from 10 to 50 yielded visually similar plots.

## G ADDITIONAL SAMPLES AND RESULTS

Figure 1 and 2 provide additional samples, organised similar to Figure 1 and 2. Figure 3 shows additional truncation curves.

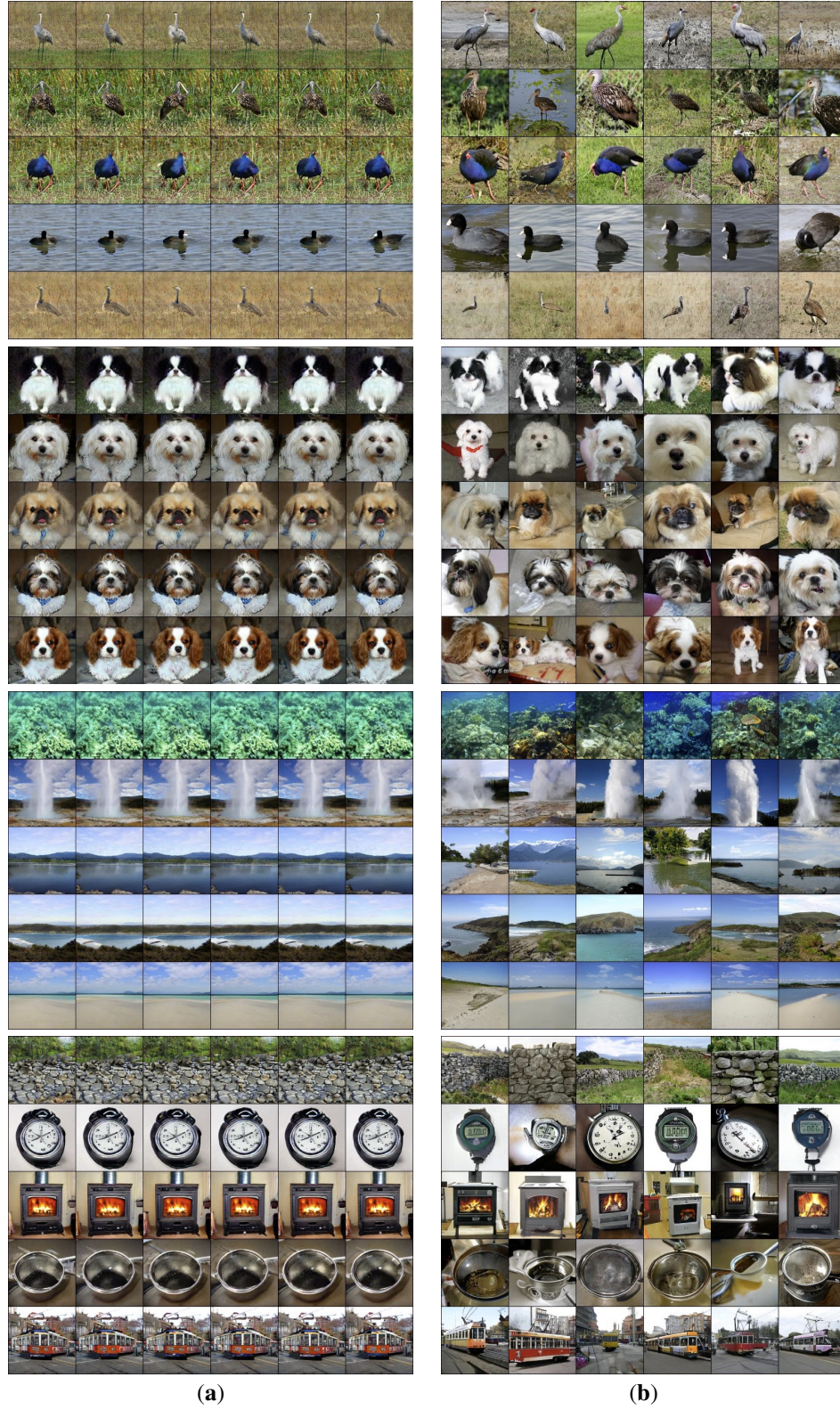


Figure 10: Samples from BigGAN-deep (a) and LOGAN (b) with similarly high Inception scores. Samples from the two panels were drawn from truncations correspond to points C, D in Figure 3. (FID/IS: (a) 27.97/259.4, (b) 8.19/259.9)



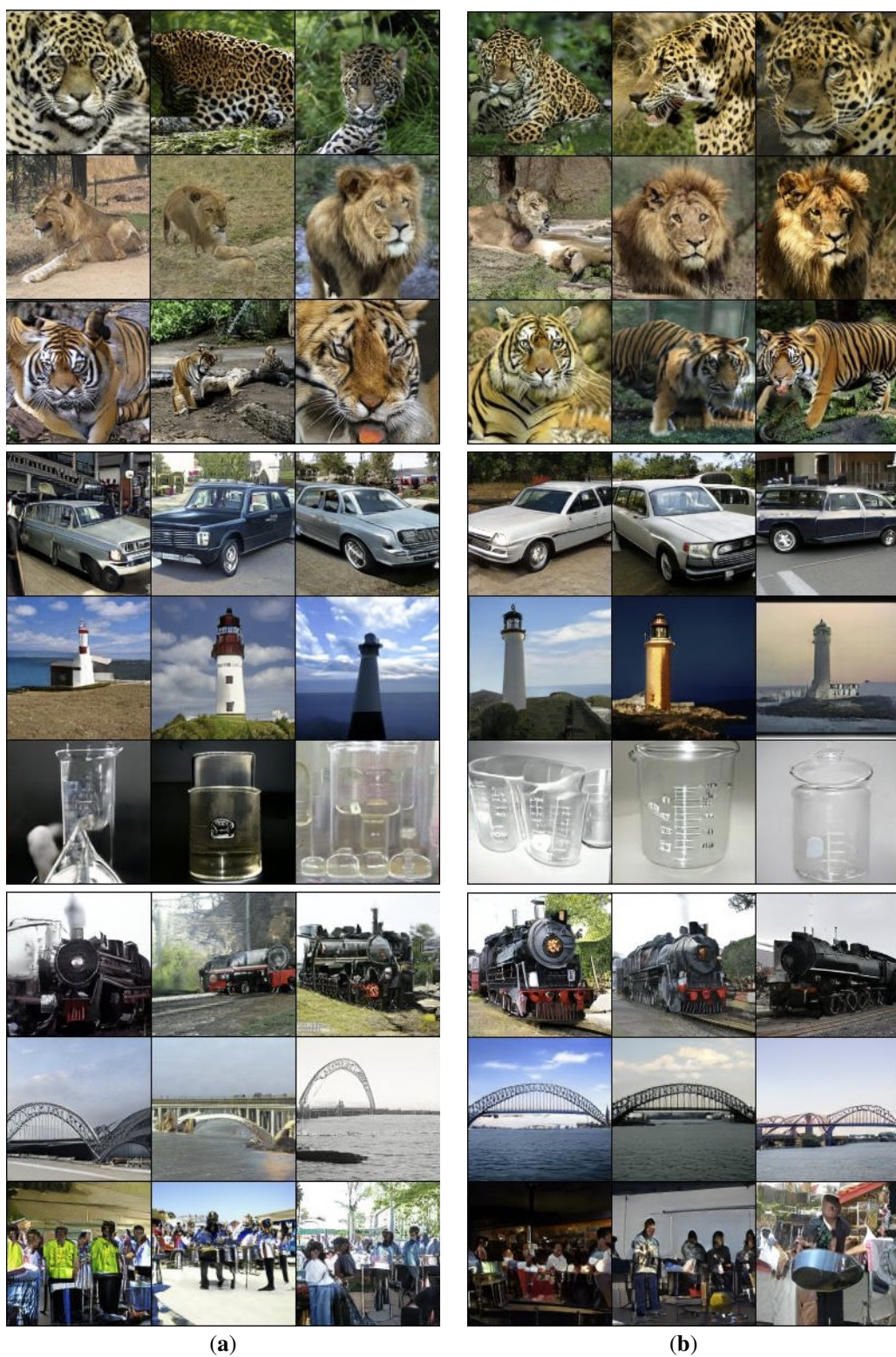


Figure 11: Samples from BigGAN-deep **(a)** and LOGAN **(b)** with similarly low FID. Samples from the two panels were drawn from truncations correspond to points A, B in figure 3b. (FID/IS: **(a)** 5.04/126.8, **(b)** 5.09/217.0)

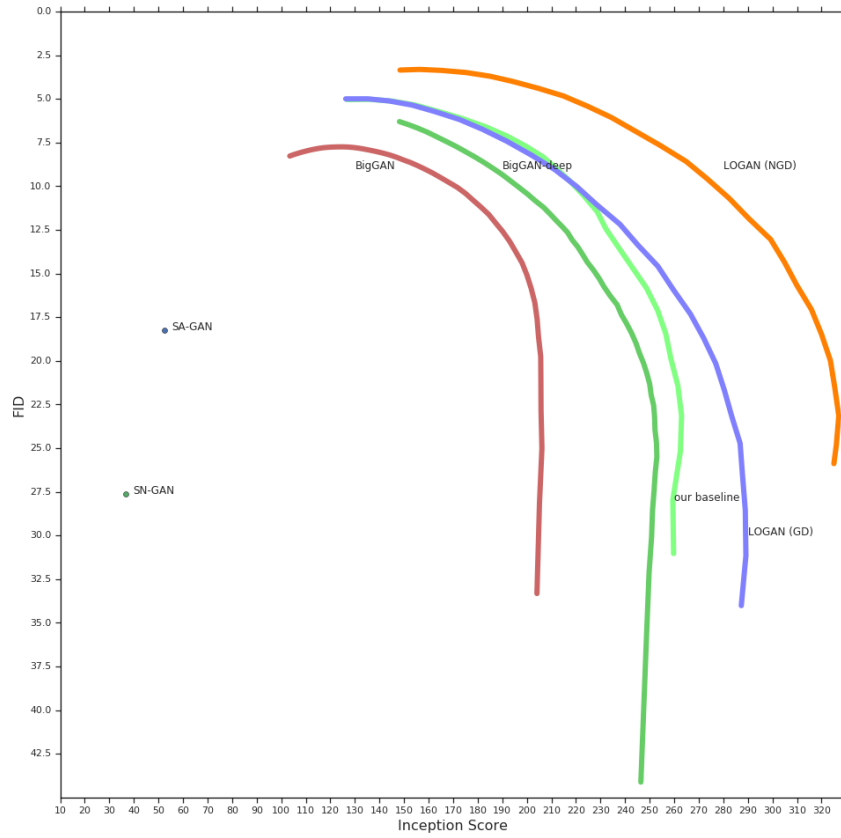


Figure 12: Truncation curves with additional baselines. In addition to the truncation curves reported in Figure 3b, here we also include the Spectral-Normalised GAN (Miyato et al., 2018), Self-Attention GAN (Zhang et al., 2019), original BigGAN and BigGAN-deep as presented in Brock et al. (2018).