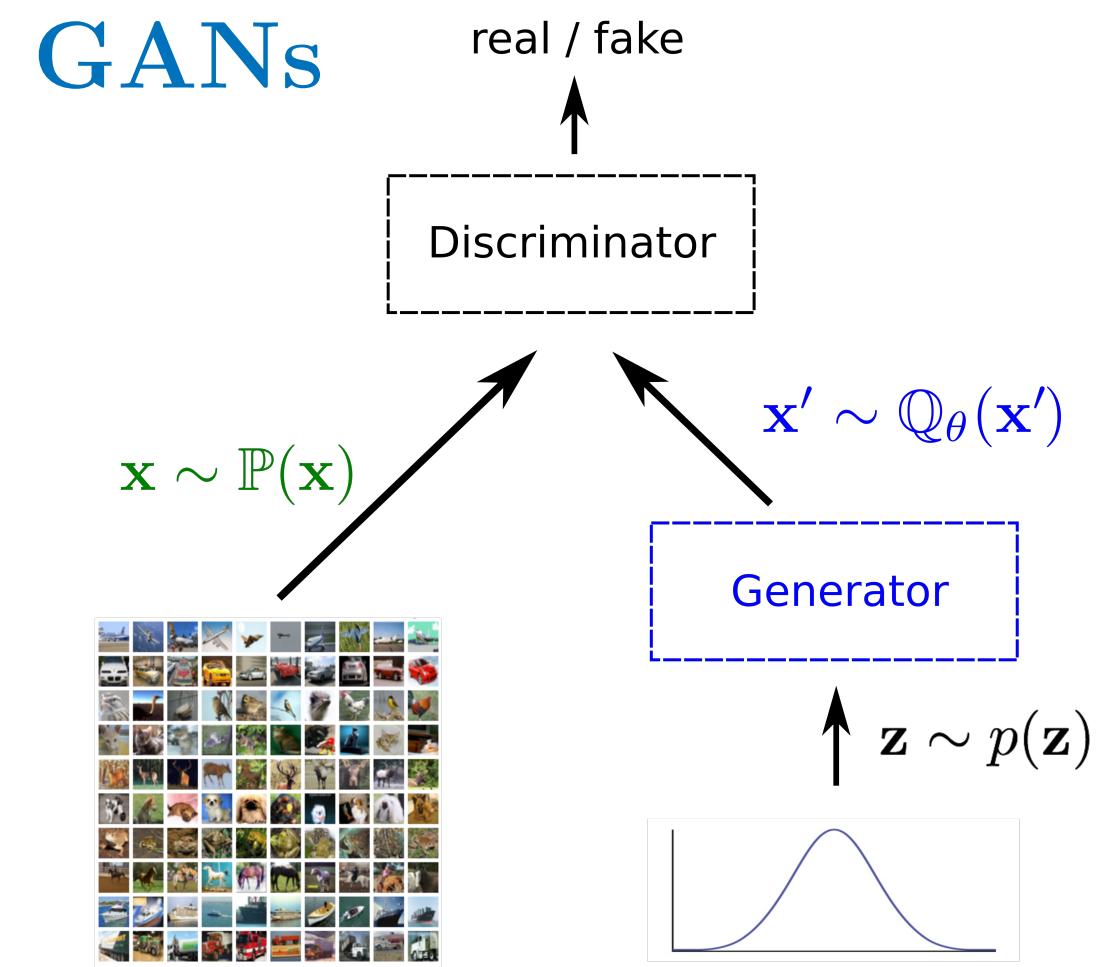


# Stabilizing Training of Generative Adversarial Networks through Regularization

Kevin Roth Aurélien Lucchi Sebastian Nowozin Thomas Hofmann

Institute of Machine Learning, ETH Zürich, Switzerland  
Microsoft Research, Cambridge, UK

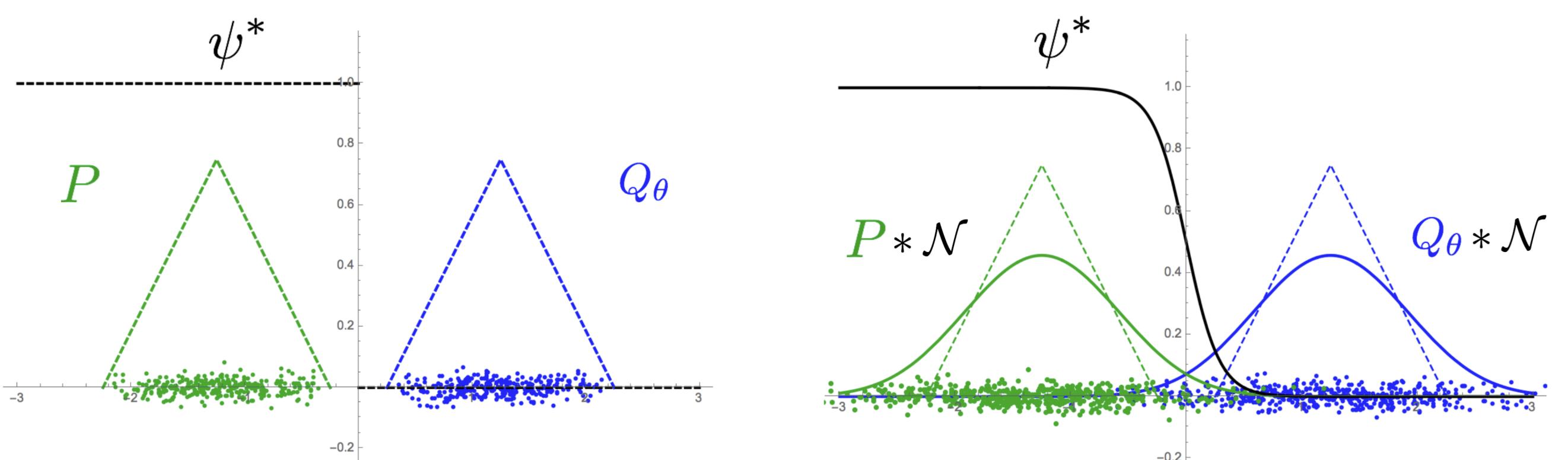


GANs are exciting, but ...  
they're also notoriously hard to train!

GAN objective [3]

$$\begin{aligned} & \min_G \max_{\varphi} \mathbf{E}_{\mathbb{P}}[\ln \varphi(x)] + \mathbf{E}_{p(z)}[\ln(1 - \varphi(G(z)))] \\ & = \min_G D_{\text{JS}}(\mathbb{P} \| Q_\theta) \text{ for Bayes-optimal } \varphi^*(x) \end{aligned}$$

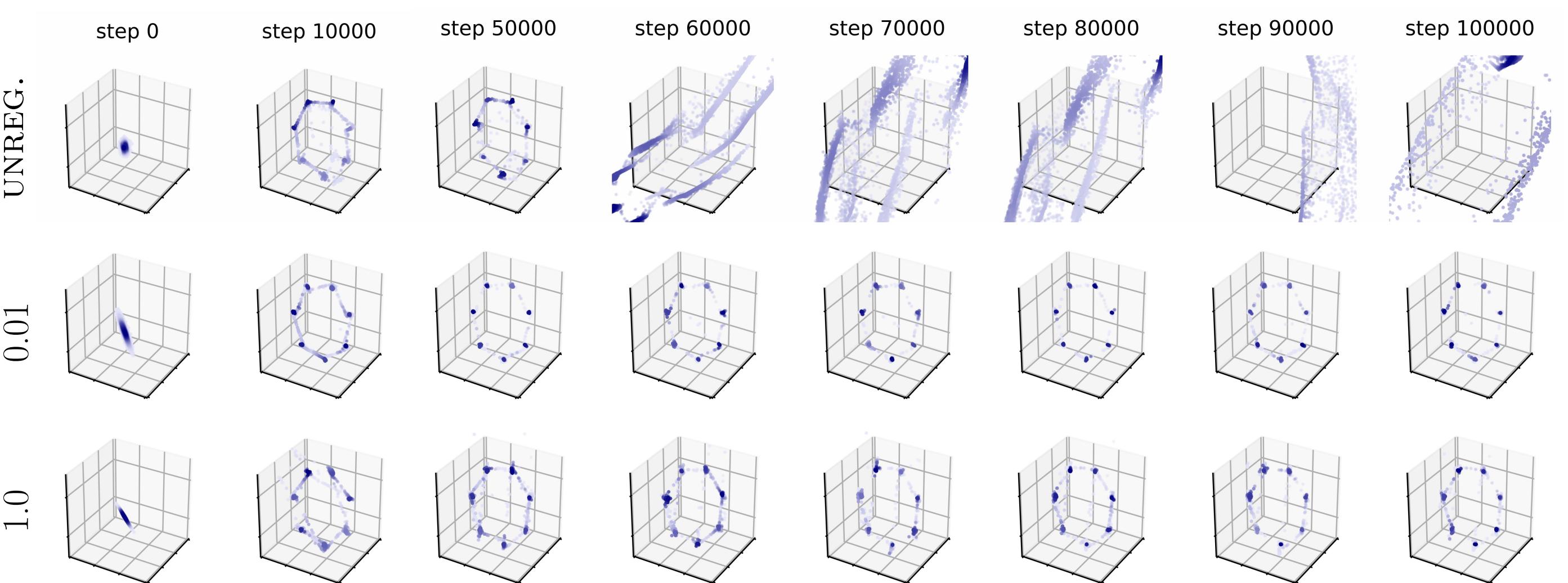
**Problem:** Dim. Mismatch<sup>[1, 6]</sup> **Solution:** Adding Noise  
or  $\text{supp}(\mathbb{P}) \cap \text{supp}(Q_\theta) = \emptyset$  (Convolving Densities)



=> Undef. f-div.  $D_f(\mathbb{P} || Q_\theta) = ???$

=> Regularizing Discriminator

## Dimensionally Misspecified Submanifold Mixture



Unstable unregularized GAN vs. stable regularized GANs for different levels of  $\gamma$ .  
The regularized GAN can essentially be trained indefinitely without collapse.

## References

- [1] M. Arjovsky and L. Bottou. Towards principled methods for training generative adversarial networks. In *NIPS 2016 Workshop on Adversarial Training*. In review for ICLR, 2017.
- [2] C. M. Bishop. Training with noise is equivalent to tikhonov regularization. *Neural computation*, 7(1):108–116, 1995.
- [3] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative Adversarial Networks. *NIPS*, 2014.
- [4] X. Nguyen, M. J. Wainwright, and M. I. Jordan. Estimating divergence functionals and the likelihood ratio by convex risk minimization. *IEEE Transactions on Information Theory*, 56(11):5847–5861, 2010.
- [5] S. Nowozin, B. Cseke, and R. Tomioka. f-GAN: Training generative neural samplers using variational divergence minimization. In *Advances in Neural Information Processing Systems*, pages 271–279, 2016.
- [6] C. K. Sonderby, J. Caballero, L. Theis, W. Shi, and F. Huszár. Amortised map inference for image super-resolution. *arXiv preprint arXiv:1610.04490*, 2016.

## Training with Noise

From  $f$ -divergences to  $f$ -GAN objectives [5]

$$D_f(\mathbb{P} \| \mathbb{Q}) \geq \sup_{\psi \in \Psi} [F(\mathbb{P}, \mathbb{Q}; \psi) := \mathbf{E}_{\mathbb{P}}[\psi] - \mathbf{E}_{\mathbb{Q}}[f^c \circ \psi]]$$

- Practitioner: Explicitly adding noise  $\xi \sim \Lambda_\gamma \equiv \mathcal{N}(0, \gamma \mathbb{I})$  to  $x \sim \mathbb{P}, Q$
- Theory: Convolving Distributions

$$\mathbf{E}_{\mathbb{P}} \mathbf{E}_\Lambda[h(\psi(x + \xi))] = \int h(\psi(x)) \int p(x - \xi) \lambda(\xi) d\xi dx = \mathbf{E}_{\mathbb{P} * \Lambda}[h \circ \psi]$$

=> Convolved  $f$ -GAN Objective

$$F(\mathbb{P} * \Lambda_\gamma, \mathbb{Q} * \Lambda_\gamma; \psi) = \mathbf{E}_{\mathbb{P} * \Lambda_\gamma}[\psi] - \mathbf{E}_{\mathbb{Q} * \Lambda_\gamma}[f^c \circ \psi]$$

## Analytic Approximation

- For small noise variance  $\gamma$  we can Taylor expand  $\psi$  around  $\xi = 0$  [2]:

$$\psi(x + \xi) = \psi(x) + [\nabla \psi(x)]^T \xi + \frac{1}{2} \xi^T [\nabla^2 \psi(x)] \xi + \mathcal{O}(\xi^3)$$

- Third-order approximation (in  $\xi$ ) of  $F_\gamma$  via  $F = F_0$  plus a correction, i.e.

$$F_\gamma(\mathbb{P}, \mathbb{Q}; \psi) = F(\mathbb{P}, \mathbb{Q}; \psi) + \frac{\gamma}{2} \{ \mathbf{E}_{\mathbb{P}}[\Delta \psi] - \mathbf{E}_{\mathbb{Q}}[\Delta(f^c \circ \psi)] \} + \mathcal{O}(\gamma^2)$$

- Interpretation: Laplace  $\Delta = \text{Tr}(\nabla^2)$  measures how much  $\psi$  and  $f^c \circ \psi$  differ from their local average

## Efficient Gradient-Based Regularization

- Chain-rule:  $\Delta(f^c \circ \psi) = (f^{c''} \circ \psi) \cdot \|\nabla \psi\|^2 + (f^{c'} \circ \psi) \Delta \psi$
- Property of optimal discriminant  $\psi^*$  [4]:  $(f^{c'} \circ \psi^*) d\mathbb{Q} = d\mathbb{P}$

=> Convenient cancellation at  $\psi = \psi^* + \mathcal{O}(\gamma)$  [2]:

$$\mathbf{E}_{\mathbb{P}}[\Delta \psi^*] - \mathbf{E}_{\mathbb{Q}}[\Delta(f^c \circ \psi^*)] = -\mathbf{E}_{\mathbb{Q}}[(f^{c''} \circ \psi^*) \cdot \|\nabla \psi^*\|^2]$$

=> Tractable regularization which avoids (i) detrimental sampling variance, (ii) explicitly convolving the distributions, and (iii) the computation of Laplacians

## Regularized $f$ -GAN

$$F_\gamma(\mathbb{P}, \mathbb{Q}; \psi) = F(\mathbb{P}, \mathbb{Q}; \psi) - \frac{\gamma}{2} \Omega_f(\mathbb{Q}; \psi), \quad \Omega_f(\mathbb{Q}; \psi) := \mathbf{E}_{\mathbb{Q}}[(f^{c''} \circ \psi) \cdot \|\nabla \psi\|^2]$$

$$\begin{aligned} F_{\text{JS}}(\mathbb{P}, \mathbb{Q}; \varphi) &= \mathbf{E}_{\mathbb{P}}[\ln(\varphi)] + \mathbf{E}_{\mathbb{Q}}[\ln(1 - \varphi)] - \frac{\gamma}{2} \Omega_{\text{JS}}(\mathbb{P}, \mathbb{Q}; \varphi) \\ \Omega_{\text{JS}}(\mathbb{P}, \mathbb{Q}; \varphi) &:= \mathbf{E}_{\mathbb{P}}[(1 - \varphi)^2 \|\nabla \varphi\|^2] + \mathbf{E}_{\mathbb{Q}}[\varphi^2 \|\nabla \varphi\|^2] \end{aligned}$$

- “soft Lipschitz” constraint, non-negative weighting function  $f^{c''} \geq 0$
- lower variance compared to explicitly adding noise
- easy to implement & computationally cheap
- can train indefinitely without collapse!

**Algorithm 1** Regularized  $f$ -GAN. Default values:  $\gamma_0 = 2.0$ ,  $\alpha = 0.01$  (with annealing),  $\gamma = 0.1$  (without annealing),  $n_\psi = 1$

**Require:** Initial noise variance  $\gamma_0$ , annealing decay factor  $\alpha$ , number of discriminator update steps  $n_\psi$  per generator iteration, minibatch size  $m$ , number of training iterations  $T$   
**Require:** Initial discriminator parameters  $\omega_0$ , initial generator parameters  $\theta_0$   
**for**  $t = 1, \dots, T$  **do**  
     $\gamma \leftarrow \gamma_0 \cdot \alpha^{t/T}$  # annealing  
    **for**  $1, \dots, n_\psi$  **do**  
        Sample minibatch of real data  $\{x^{(1)}, \dots, x^{(m)}\} \sim \mathbb{P}$ .  
        Sample minibatch of latent variables from prior  $\{z^{(1)}, \dots, z^{(m)}\} \sim p(z)$ .  
         $\omega \leftarrow \omega + \nabla_\omega (F(\omega, \theta) - \frac{\gamma}{2} \Omega_f(\omega, \theta))$  # gradient ascent  
    **end for**  
    Sample minibatch of latent variables from prior  $\{z^{(1)}, \dots, z^{(m)}\} \sim p(z)$ .  
     $\theta \leftarrow \theta - \nabla_\theta F(\omega, \theta)$  # gradient descent  
**end for**

## Cross-Testing Protocol

Regularized  $\gamma = 0.1$

True Cond.	
Pos.	Neg.
0.9688	0.0002
0.0312	0.9998

Cross-testing: FP: 0.0

Unregularized

True Cond.	
Pos.	Neg.
1.0	0.0013
0.0	0.9987

Cross-testing: FP: 1.0

Cross-Testing:  
Classify 10k samples generated by the regularized GAN with the discriminator of the unregularized GAN and vice versa.

=> Regularized GAN generalizes better!

## Stability across Architectures



## Sample Quality and Diversity



tf code → [github.com/rothk](https://github.com/rothk)

```
# JS-Regularizer
# -----
def Discriminator_Regularizer(self, D1, D1_logits, D1_arg, D2, D2_logits, D2_arg):
    grad_D1_logits = tf.gradients(D1_logits, D1_logits, D1_arg)[0]
    grad_D2_logits = tf.gradients(D2_logits, D2_logits, D2_arg)[0]
    grad_D1_logits_norm = tf.norm(tf.reshape(grad_D1_logits, [self.batch_size,-1]), axis=1, keep_dims=True)
    grad_D2_logits_norm = tf.norm(tf.reshape(grad_D2_logits, [self.batch_size,-1]), axis=1, keep_dims=True)
    reg_D1 = tf.multiply(tf.square(1.0-D1), tf.square(grad_D1_logits_norm))
    reg_D2 = tf.multiply(tf.square(D2), tf.square(grad_D2_logits_norm))
    return tf.reduce_mean(reg_D1 + reg_D2)
```