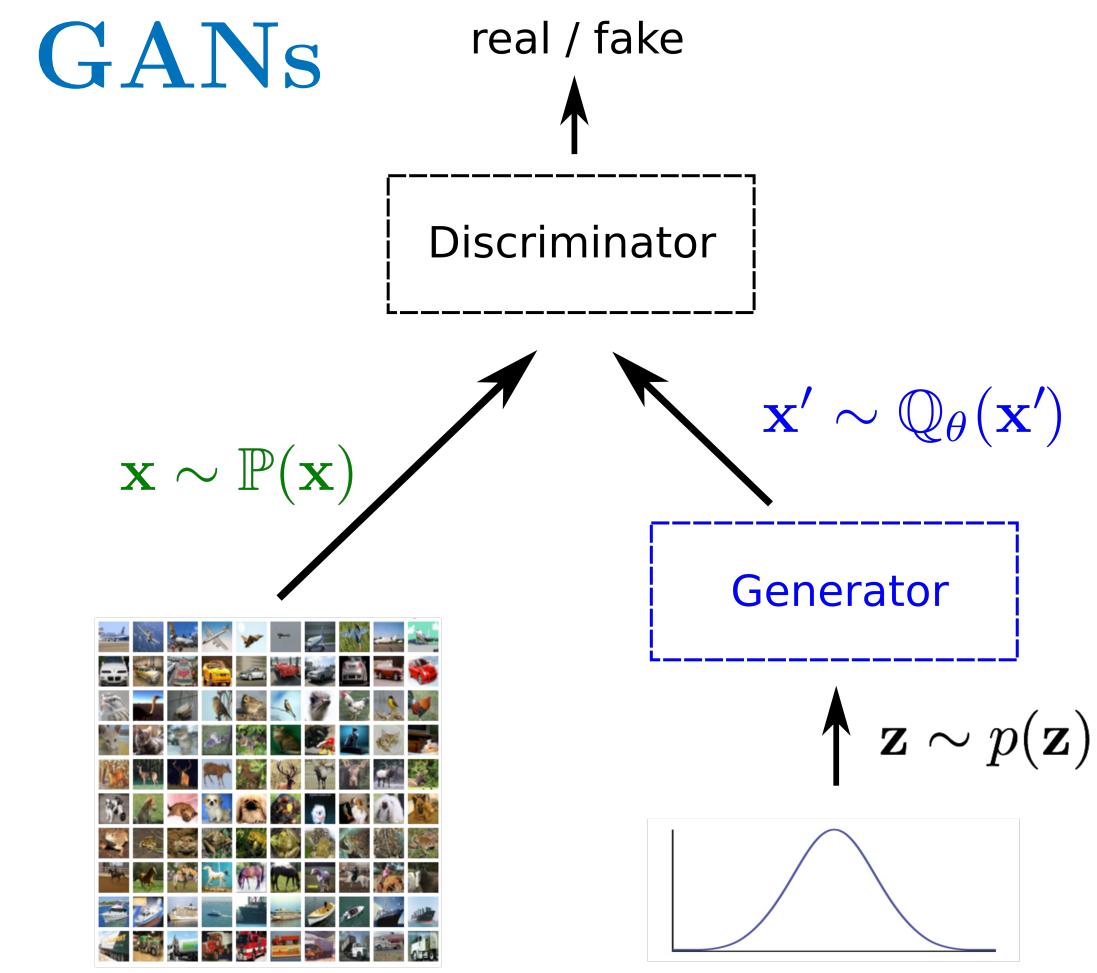


Stabilizing Training of Generative Adversarial Networks through Regularization

Kevin Roth Aurélien Lucchi Sebastian Nowozin Thomas Hofmann

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



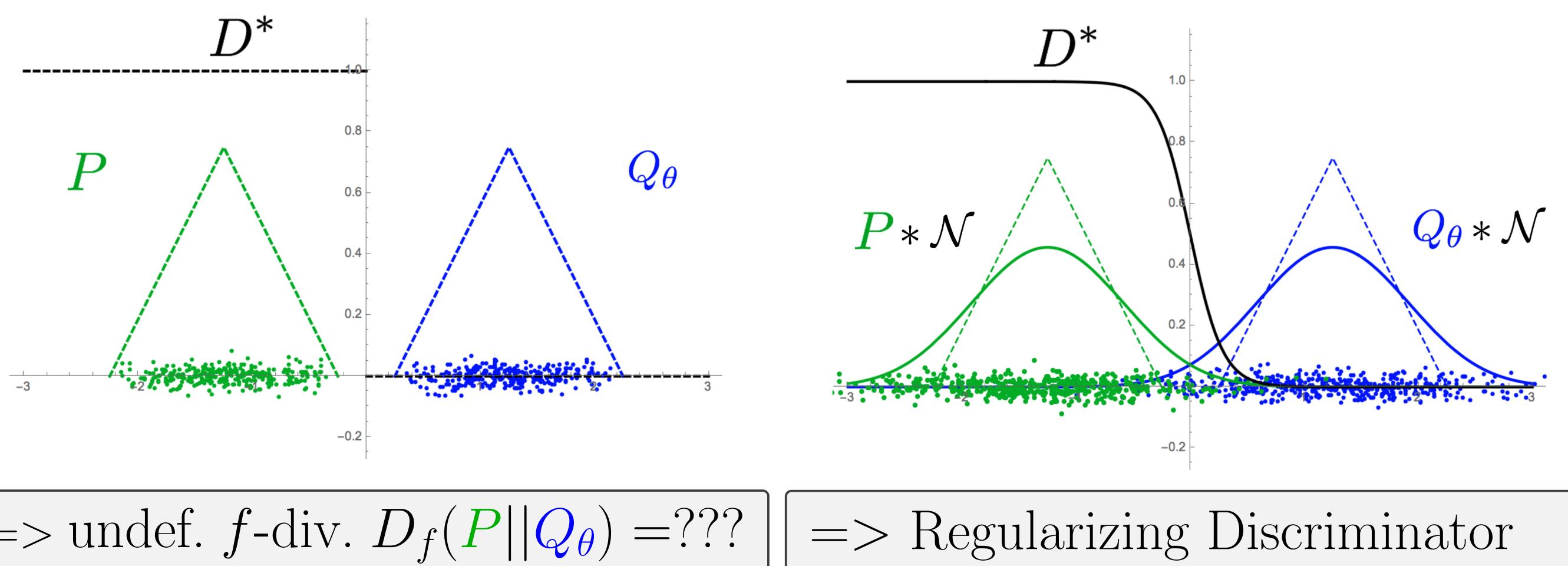
GANs are exciting, but ...
they're also notoriously hard to train!

GAN objective [3]

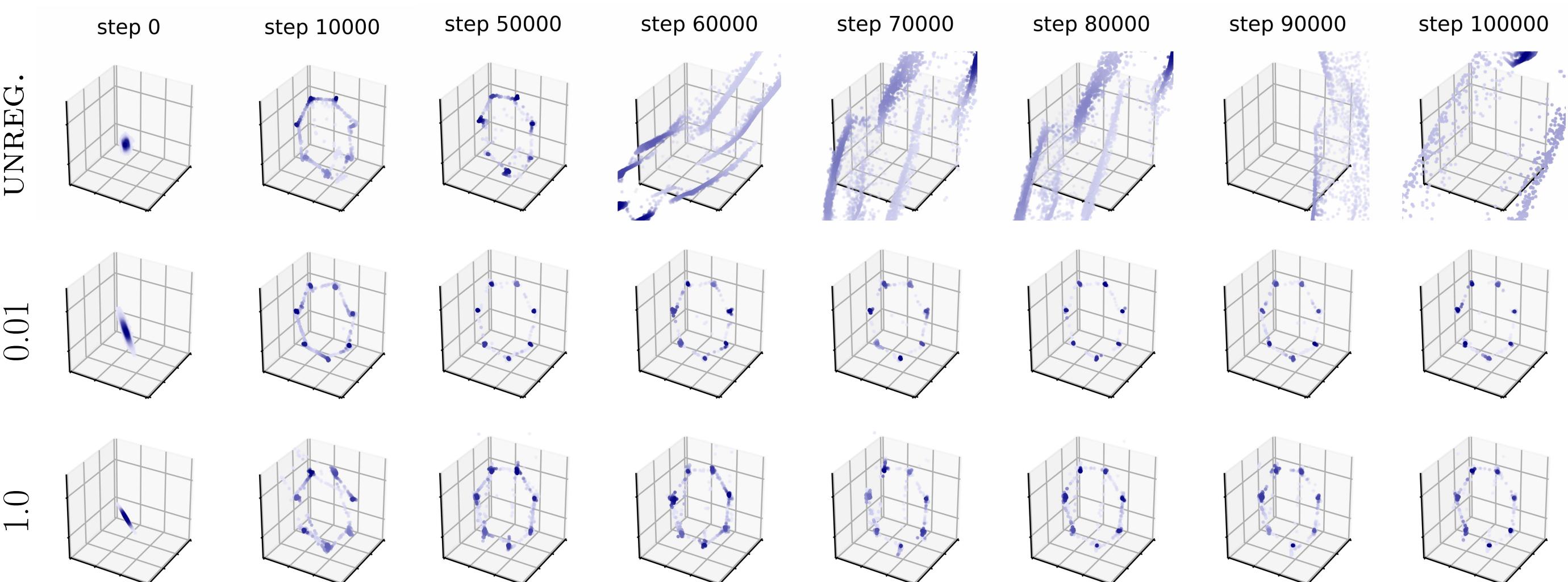
$$\begin{aligned} & \min_G \max_D \mathbb{E}_P[\log D(x)] + \mathbb{E}_{p(z)}[\log(1 - D(G(z)))] \\ & = \min_G D_{JS}(P||Q_\theta) \text{ for Bayes-optimal } D^*(x) \end{aligned}$$

Problem: dim. mismatch
or $\text{supp}(P) \cap \text{supp}(Q_\theta) = \emptyset$

Solution: Adding Noise
(Convolving Densities)



Dimensionally Misspecified Submanifold Mixture



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, volume 2016, 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. Sønderby, 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]

$$D_f(P||Q) \geq \sup_{\psi \in \Psi} (\mathbb{E}_P[\psi] - \mathbb{E}_Q[f^c \circ \psi]) \quad (1)$$

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

• Theory: Convolving Distributions

$$\mathbb{E}_P \mathbb{E}_\Lambda[h(\psi(x + \xi))] = \int h(\psi(x)) \int p(x - \xi) \lambda(\xi) d\xi dx = \mathbb{E}_{P * \Lambda}[h \circ \psi] \quad (2)$$

=> Convolved f -GAN Objective

$$F(\mathbb{P} * \Lambda_\gamma, \mathbb{Q} * \Lambda_\gamma; \psi) = \mathbb{E}_{P * \Lambda_\gamma}[\psi] - \mathbb{E}_{Q * \Lambda_\gamma}[f^c \circ \psi] \quad (3)$$

Analytic Approximation

• For small noise variance γ we can Taylor expand ψ around $\xi = 0$ [2]:

$$\psi(x + \xi) = \psi(x) + [\nabla \psi(x)] \xi + \frac{1}{2} \xi^T [\nabla^2 \psi(x)] \xi + \mathcal{O}(\xi^3) \quad (4)$$

• Third-order approximation

$$F_\gamma = F_0 + \frac{\gamma}{2} \{ \mathbb{E}_P[\Delta \psi] - \mathbb{E}_Q[\Delta(f^c \circ \psi)] \} + \mathcal{O}(\gamma^2) \quad (5)$$

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

Simplification

• Chain-rule:

$$\Delta(f^c \circ \psi) = (f^{c''} \circ \psi) \cdot \|\nabla \psi\|^2 + (f^{c'} \circ \psi) \Delta \psi \quad (6)$$

• Property of optimal discriminant ψ^* [4]

$$(f^{c'} \circ \psi^*) dQ = dP \quad (7)$$

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

$$\mathbb{E}_P[\Delta \psi^*] - \mathbb{E}_Q[\Delta(f^c \circ \psi^*)] = -\mathbb{E}_Q[(f^{c''} \circ \psi^*) \cdot \|\nabla \psi^*\|^2] \quad (8)$$

Regularized f -GAN

$$\begin{aligned} F_\gamma(\mathbb{P}, \mathbb{Q}_\theta; \psi) &= \mathbb{E}_P[\psi] - \mathbb{E}_{Q_\theta}[f^c \circ \psi] - \frac{\gamma}{2} \Omega_f(Q_\theta; \psi) \\ \Omega_f(Q_\theta; \psi) &:= \mathbb{E}_{Q_\theta}[(f^{c''} \circ \psi) \cdot \|\nabla \psi\|^2] \end{aligned} \quad (9)$$

- “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 γ_0 , annealing decay factor α , number of discriminator update steps n_ψ per generator iteration, minibatch size m , number of training iterations T

Require: Initial discriminator parameters ω_0 , initial generator parameters θ_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 $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\} \sim \mathbb{P}$.

Sample minibatch of latent variables from prior $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\} \sim p(\mathbf{z})$.

$\omega \leftarrow \omega + \nabla_\omega (F(\omega, \theta) - \frac{\gamma}{2} \Omega_f(\omega, \theta))$ # gradient ascent

end for

Sample minibatch of latent variables from prior $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\} \sim p(\mathbf{z})$.

$\theta \leftarrow \theta - \nabla_\theta F(\omega, \theta)$ # gradient descent

end for

Cross-Testing Protocol

Regularized $\gamma = 0.1$
True Cond.

	Pos.	Neg.
Pos.	0.9688	0.0002
Neg.	0.0312	0.9998

Cross-testing: FP: 0.0

Unregularized
True Cond.

	Pos.	Neg.
Pos.	1.0	0.0013
Neg.	0.0	0.9987

Cross-testing: FP: 1.0

Cross-Testing:

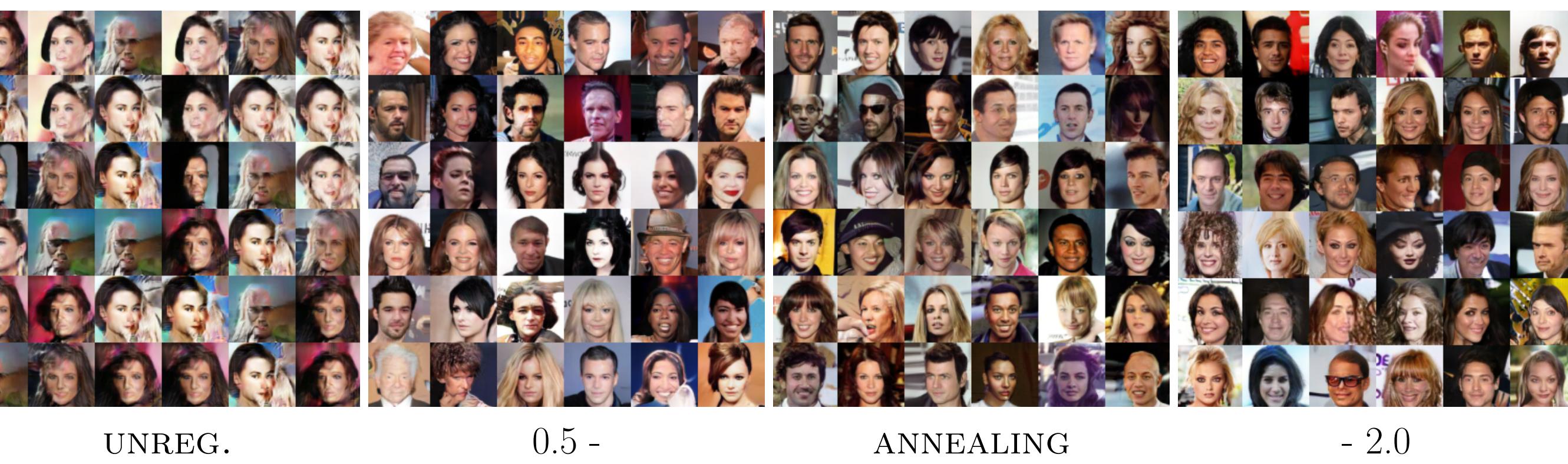
Classify 10k samples generated by the regularized GAN with the discriminator of the unregularized GAN and vice versa.

=> Regularized GAN generalizes better!

Stability across Architectures



Sample quality and diversity



tf code → github.com/rothk

```
# JS-Regularizer
# -----
def Discriminator_Regularizer(self, D1, D1_logits, D1_arg, D2, D2_logits, D2_arg):
    grad_D1_logsits = tf.gradients(D1_logsits, D1_arg)[0]
    grad_D2_logsits = tf.gradients(D2_logsits, D2_arg)[0]
    grad_D1_logsits_norm = tf.norm(tf.reshape(grad_D1_logsits, [self.batch_size,-1]), axis=1, keep_dims=True)
    grad_D2_logsits_norm = tf.norm(tf.reshape(grad_D2_logsits, [self.batch_size,-1]), axis=1, keep_dims=True)
    reg_D1 = tf.multiply(tf.square(1.0-D1), tf.square(grad_D1_logsits_norm))
    reg_D2 = tf.multiply(tf.square(D2), tf.square(grad_D2_logsits_norm))
    return tf.reduce_mean(reg_D1 + reg_D2)
```