

# GANs Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium

Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter

LIT AI Lab & Institute of Bioinformatics, Johannes Kepler Universität Linz, Austria

**Abstract** Generative Adversarial Networks (GANs) excel at creating realistic images with complex models for which maximum likelihood is infeasible. However, the convergence of GAN training has still not been proved. We propose a two time-scale update rule (TTUR) for training GANs with stochastic gradient descent on arbitrary GAN loss functions. TTUR has an individual learning rate for both the discriminator and the generator. Using the theory of stochastic approximation, we prove that the TTUR converges under mild assumptions to a stationary local Nash equilibrium. The convergence carries over to the popular Adam optimization, for which we prove that it follows the dynamics of a heavy ball with friction and thus prefers flat minima in the objective landscape. For the evaluation of the performance of GANs at image generation, we introduce the “Fréchet Inception Distance” (FID) which captures the similarity of generated images to real ones better than the Inception Score. In experiments, TTUR improves learning for DCGANs and Improved Wasserstein GANs (WGAN-GP) outperforming conventional GAN training on CelebA, CIFAR-10, SVHN, LSUN Bedrooms, and the One Billion Word Benchmark.

## General GAN Update Settings

**Generator Updates:**

$$\theta_{n+1} = \theta_n + a(n) \left( h(\theta_n, w_n) + M_n^{(\theta)} \right)$$

learning rate:  $a(n)$

true gradient:  $h(\theta, w)$

gradient noise:  $M^{(\theta)}$

stochastic gradient:  $\hat{h}(\theta, w) = h(\theta, w) + M^{(\theta)}$

**Discriminator Updates:**

$$w_{n+1} = w_n + b(n) \left( g(\theta_n, w_n) + M_n^{(w)} \right)$$

learning rate:  $b(n)$

true gradient:  $g(\theta, w)$

gradient noise:  $M^{(w)}$

stochastic gradient:  $\hat{g}(\theta, w) = g(\theta, w) + M^{(w)}$

## Two Time-Scale Update Rule Converges

**Theorem 1** (Borkar 1997). *If the assumptions are satisfied, then the updates converge to  $(\theta^*, \lambda(\theta^*))$  a.s.*

$\theta^*$ ,  $\lambda(\theta^*)$  are local asymptotically stable attractors,  $g(\theta^*, \lambda(\theta^*)) = 0$ ,  $h(\theta^*, \lambda(\theta^*)) = 0$ , the solution  $(\theta^*, \lambda(\theta^*))$  is a **stationary local Nash equilibrium**. The assumptions are:

(A1) The gradients  $h$  and  $g$  are Lipschitz.

(A2)  $\sum_n a(n) = \infty$ ,  $\sum_n a^2(n) < \infty$ ,  $\sum_n b(n) = \infty$ ,  $\sum_n b^2(n) < \infty$ ,  $a(n) = o(b(n))$ .

(A3) The stochastic gradient errors  $\{M_n^{(\theta)}\}$  and  $\{M_n^{(w)}\}$  are martingale difference sequences w.r.t. the increasing  $\sigma$ -field  $\mathcal{F}_n = \sigma(\theta_l, w_l, M_l^{(\theta)}, M_l^{(w)}, l \leq n)$ ,  $n \geq 0$  with  $E[\|M_n^{(\theta)}\|^2 | \mathcal{F}_n^{(\theta)}] \leq B_1$  and  $E[\|M_n^{(w)}\|^2 | \mathcal{F}_n^{(w)}] \leq B_2$ , where  $B_1$  and  $B_2$  are positive deterministic constants.

(A4) For each  $\theta$ , the ODE  $\dot{w}(t) = g(\theta, w(t))$  has a local asymptotically stable attractor  $\lambda(\theta)$  within a domain of attraction  $G_\theta$  such that  $\lambda$  is Lipschitz. The ODE  $\dot{\theta}(t) = h(\theta(t), \lambda(\theta(t)))$  has a local asymptotically stable attractor  $\theta^*$  within a domain of attraction.

(A5)  $\sup_n \|\theta_n\| < \infty$  and  $\sup_n \|w_n\| < \infty$ .

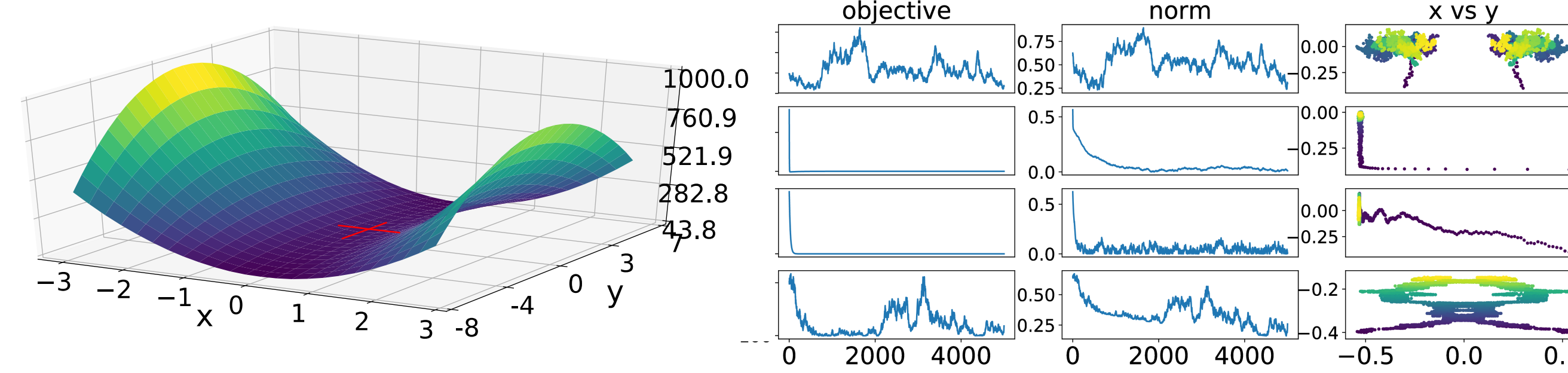
## Fréchet Inception Distance to Evaluate GANs

For  $p(\cdot)$  the distribution of model samples and  $p_w(\cdot)$  the distribution of the real world samples the equality  $p(\cdot) = p_w(\cdot)$  holds if and only if  $\int p(\cdot) f(x) dx = \int p_w(\cdot) f(x) dx$  where  $f(x)$  is the first and second order polynomial of the coding layer of an Inception model. From the resulting moments the vector of mean and covariance  $(m, C)$  can be computed.

The Fréchet distance  $d(\cdot, \cdot)$  between the Gaussian with mean and covariance  $(m, C)$  obtained from  $p(\cdot)$  and the Gaussian  $(m_w, C_w)$  obtained from  $p_w(\cdot)$  is called the “Fréchet Inception Distance” (FID), given by:

$$d^2((m, C), (m_w, C_w)) = \|m - m_w\|_2^2 + \text{Tr}(C + C_w - 2(CC_w)^{1/2}).$$

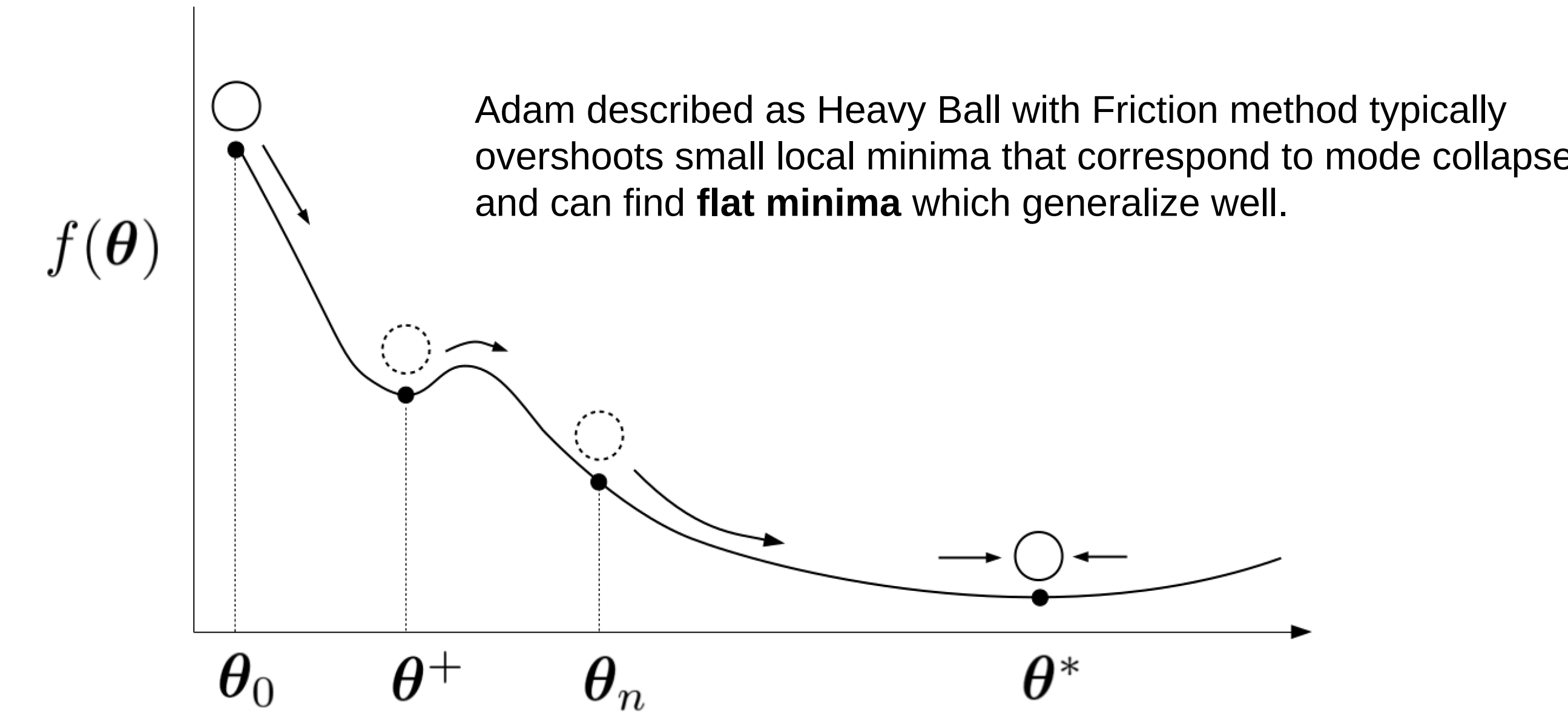
## Saddle Point Example



TTUR on a simple toy min/max saddle point problem. The objective  $\min_x \max_y f(x, y) = (1 + x^2)(100 - y^2)$  with a saddle point at  $(x, y) = (0, 0)$  fulfills assumption A4. The norm  $\|(x, y)\|$  measures the distance of the parameter vector to the saddle point.  $(x, y)$  is updated by *gradient descent* in  $x$  and *gradient ascent* in  $y$  with additive Gaussian noise to simulate a stochastic update. The updates should converge to the saddle point  $(x, y) = (0, 0)$  with objective value  $f(0, 0) = 100$  and norm 0.

- First row: One time-scale update rule with learning rate 0.01 diverges.
- Second row: One time scale update rule with learning rate 0.001 converges.
- Third row: TTUR with learning rates 0.0001 for  $x$  and 0.01 for  $y$  converges faster than with equal time-scale updates. Moves directly to the saddle point.
- Fourth row: TTUR with learning rates 0.01 and 0.0001 converges but slower.

## Adam Follows a Heavy Ball with Friction ODE



**Theorem 2.** *If Adam is used with  $\beta_1 = 1 - a(n+1)r(n)$ ,  $\beta_2 = 1 - \alpha a(n+1)r(n)$  and with  $\nabla f$  as the full gradient of the lower bounded, continuously differentiable objective  $f$ , then for stationary second moments of the gradient, Adam follows the differential equation for Heavy Ball with Friction (HBF):*

$$\ddot{\theta}_t + a(t) \dot{\theta}_t + \nabla f(\theta_t) = 0.$$

*Adam converges for gradients  $\nabla f$  that are  $L$ -Lipschitz.*

Adam described by a differential equation with the Lyapunov function:

$$E(t) = \frac{1}{2} |\dot{\theta}(t)|^2 + f(\theta(t)) \quad \text{with}$$

$$\dot{E}(t) = -a(t) |\dot{\theta}(t)|^2 < 0.$$

Gadat et al. derived a discrete and stochastic version of the HBF:

$$\theta_{n+1} = \theta_n - a(n+1) m_n$$

$$m_{n+1} = m_n + a(n+1) r(n) (\nabla f(\theta_n) - m_n) + a(n+1) r(n) M_{n+1}.$$

The recursion can be rewritten as

$$\theta_{n+1} = \theta_n - a(n+1) m_n$$

$$m_{n+1} = (1 - a(n+1) r(n)) m_n + a(n+1) r(n) (\nabla f(\theta_n) + M_{n+1}).$$

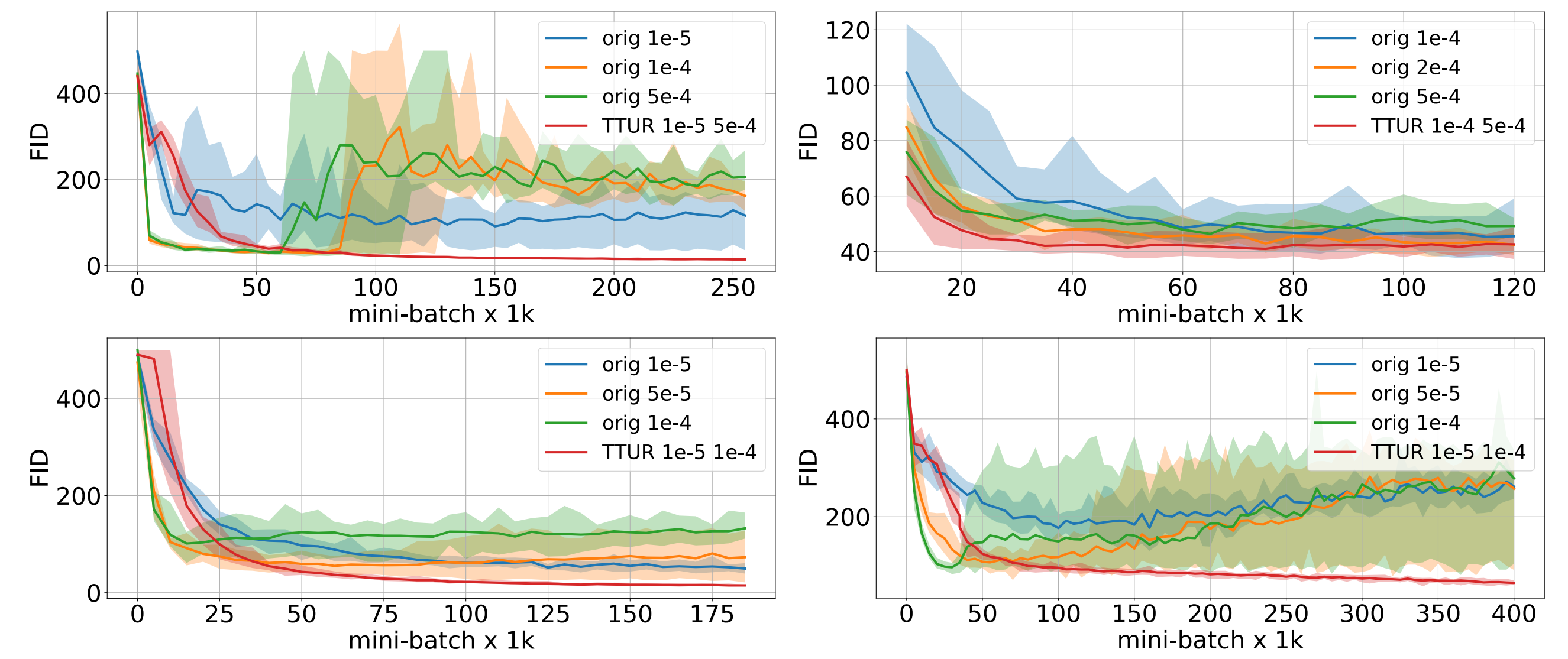
The recursion is the first moment update of Adam.

Learning GANs with TTUR and Adam converges.

## Experiments

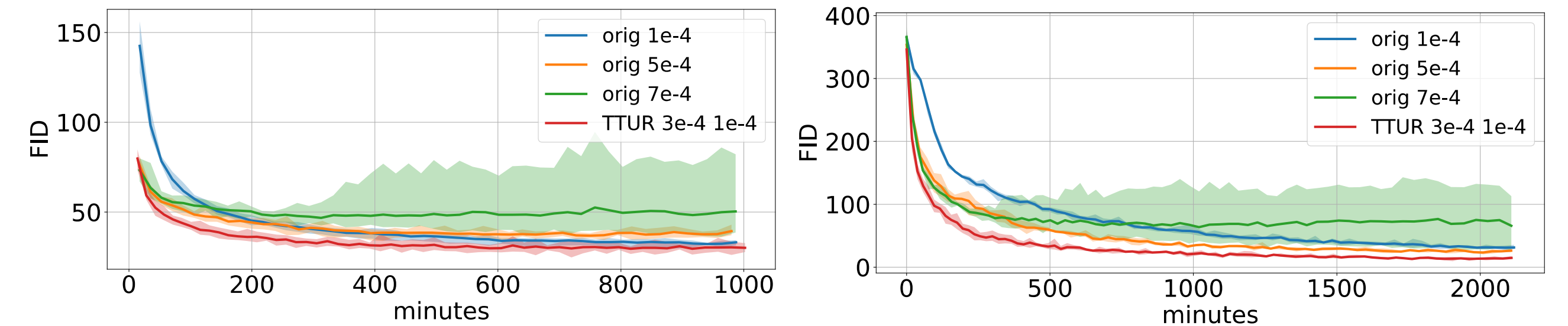
**Model Selection and Evaluation.** Learning rates are optimized to be small to ensure convergence but large enough to allow fast learning. Evaluation metric is the FID for image and the Jensen-Shannon divergence (JSD) for language data. WGAN-GP updates the critic five times for image and ten times for language data per iteration, TTUR updates the critic only once, therefore the training progress is aligned to wall-clock time for better comparison. TTUR learning rates are given for the discriminator  $b$  and generator  $a$  as: “TTUR  $b$   $a$ ”.

**DCGAN on Image Data:**



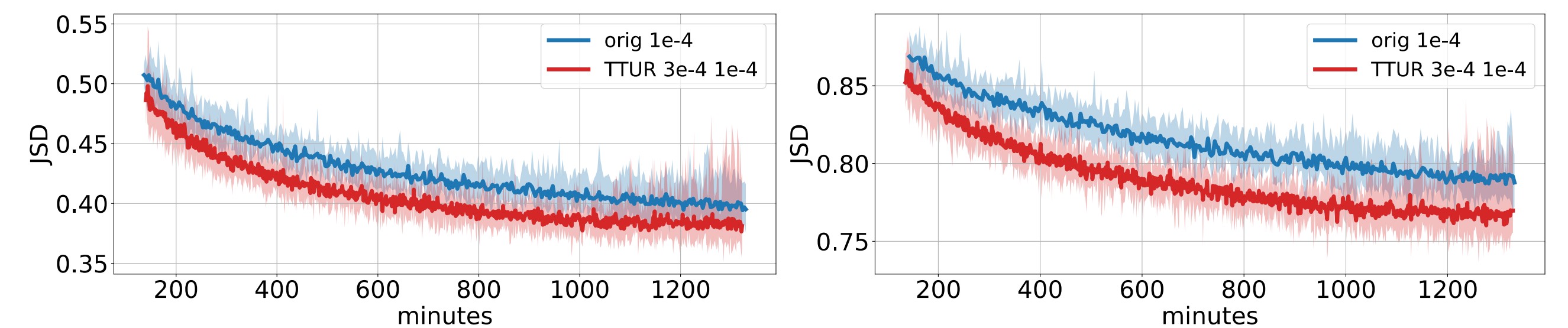
Mean FID (solid line) surrounded by a shaded area bounded by the maximum and minimum over 8 runs. **Top left:** CelebA. **Top right:** CIFAR-10, starting at mini-batch 10k. **Bottom left:** SVHN. **Bottom right:** LSUN Bedrooms.

**WGAN-GP on Image Data:**



Mean FID (solid line) surrounded by a shaded area bounded by the maximum and minimum over 8 runs, aligned to wall-clock time. **Left:** CIFAR-10, starting at minute 20. **Right:** LSUN Bedrooms.

**WGAN-GP on Language Data (One Billion Word Benchmark):**



Mean normalized JSD (solid line) surrounded by a shaded area bounded by the maximum and minimum over 10 runs, aligned to wall-clock time, starting at minute 150. **Left:** 4-gram statistics. **Right:** 6-gram statistics.

**Results:**

DCGAN Image								
dataset	method	$b, a$	updates	FID	method	$b = a$	updates	FID
CelebA	TTUR	1e-5, 5e-4	225k	<b>12.5</b>	orig	5e-4	70k	21.4
CIFAR-10	TTUR	1e-4, 5e-4	75k	<b>36.9</b>	orig	1e-4	100k	37.7
SVHN	TTUR	1e-5, 1e-4	165k	<b>12.5</b>	orig	5e-5	185k	21.4
LSUN	TTUR	1e-5, 1e-4	340k	<b>57.5</b>	orig	5e-5	70k	70.4
WGAN-GP Image								
dataset	method	$b, a$	time(m)	FID	method	$b = a$	time(m)	FID
CIFAR-10	TTUR	3e-4, 1e-4	700	<b>24.8</b>	orig	1e-4	800	29.3
LSUN	TTUR	3e-4, 1e-4	1900	<b>9.5</b>	orig	1e-4	2010	20.5
WGAN-GP Language								
$n$ -gram	method	$b, a$	time(m)	JSD	method	$b = a$	time(m)	JSD
4-gram	TTUR	3e-4, 1e-4	1150	<b>0.35</b>	orig	1e-4	1040	0.38
6-gram	TTUR	3e-4, 1e-4	1120	<b>0.74</b>	orig	1e-4	1070	0.77

Performance with respect to the FID and JSD for optimized number of updates.