

Reducing Noise in GAN Training with Variance Reduced Extragradient

Tatjana Chavdarova^{*,1,2}, Gauthier Gidel^{*,1,3},
 François Fleuret² and Simon Lacoste-Julien^{1,3}

^{*}Equal contribution; ¹Mila, Université de Montréal; ²EPFL, Idiap; ³Element AI

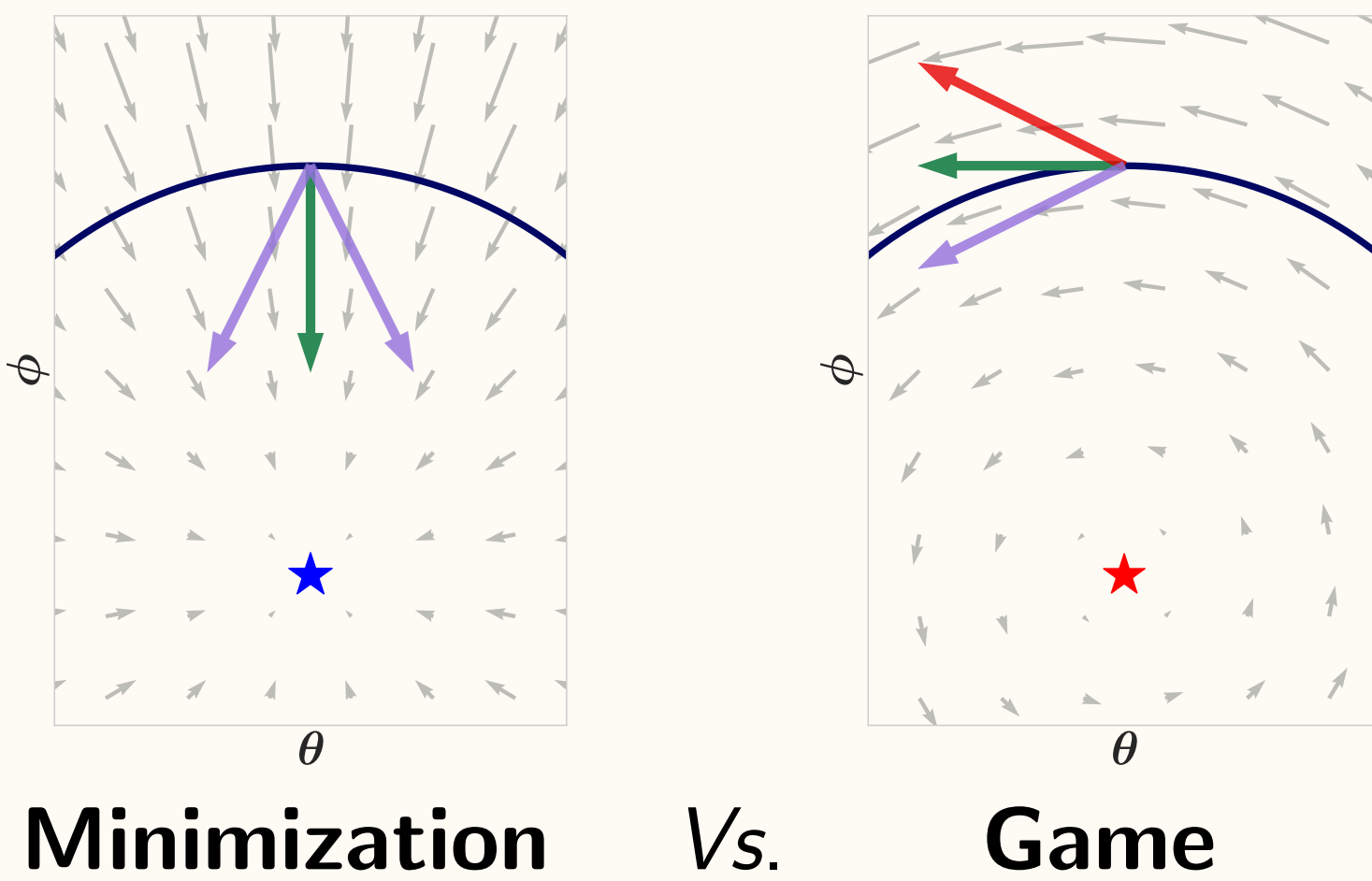
Overview

Takeaways

- Games harder to optimize → **Extragradient** damps the oscillations of the game.
- Not having the full gradient (e.g. stochastic gradient) breaks Extragradient.
- We propose “**SVRE**” that uses **variance reduction** to fix stochastic Extragradient
- Theoretically the fastest method (under some standard assumptions).
- Empirically much more stable than the baseline and yields improvements late in learning.

Motivation: VR for game optimization

- BigGAN: 8-fold increased batch size yields 46% relative improvement of Inception Score on ImageNet
- In practice (often) only the variance component of *Adam* is used: $\beta_1 = 0$ (tuned)
- Intuition on why noisy game vector field is more problematic than noisy minimization:



Background

Two-player games Equilibrium

Generalizes mini-max formulation:

$$\theta^* \in \arg \min_{\theta \in \Theta} \mathcal{L}^G(\theta, \varphi^*),$$

$$\varphi^* \in \arg \min_{\varphi \in \Phi} \mathcal{L}^D(\theta^*, \varphi).$$

Stationary conditions

Objective: point with zero gradient.

$$\|\nabla_{\theta} \mathcal{L}^G(\theta^*, \varphi^*)\| = \|\nabla_{\varphi} \mathcal{L}^D(\theta^*, \varphi^*)\| = 0.$$

$\omega \stackrel{\text{def}}{=} (\theta, \varphi)$, $\omega^* \stackrel{\text{def}}{=} (\theta^*, \varphi^*)$, $\Omega \stackrel{\text{def}}{=} \Theta \times \Phi$,
 Can be reformulated as $F(\omega^*) = 0$ where,

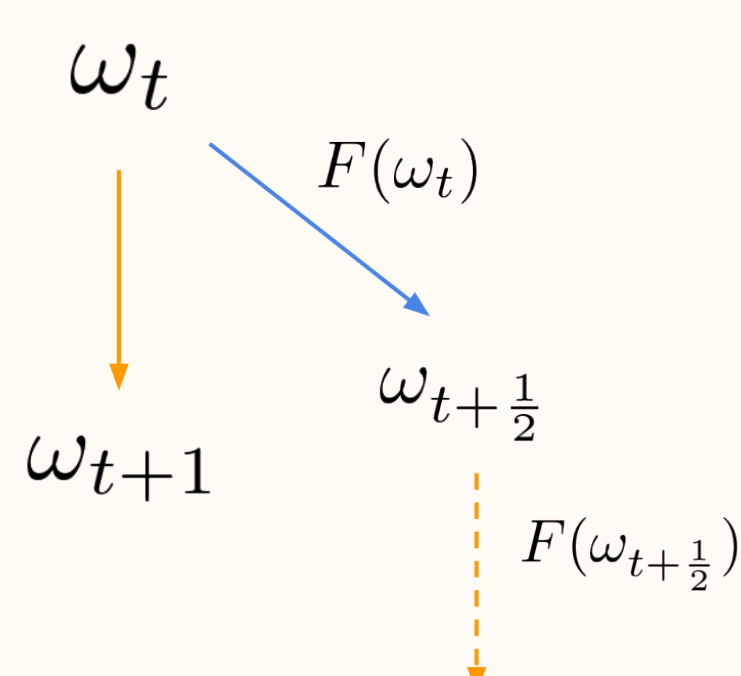
$$F(\omega) \stackrel{\text{def}}{=} (\nabla_{\theta} \mathcal{L}^G(\theta, \varphi), \nabla_{\varphi} \mathcal{L}^D(\theta, \varphi)).$$

Extragradient

$$\omega_{t+\frac{1}{2}} = \omega_t - \gamma_t F(\omega_t) \quad (\text{extrapolation})$$

$$\omega_{t+1} = \omega_t - \gamma_t F(\omega_{t+\frac{1}{2}}) \quad (\text{update})$$

Intuition: Look one step in the future and anticipate the next move of the adversary.
 Close to *implicit* method.



1 minute 5 minutes 10 minutes
 Reducing Noise in GANs

Methods for solving bilinear games

$$\min_{\theta \in \mathbb{R}^d} \max_{\varphi \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \theta^\top A_i \varphi$$

Method	Gradient method	Extragradient
Batch	$\ \omega_t - \omega^*\ \rightarrow \infty$	$\ \omega_t - \omega^*\ \rightarrow 0$
Stochastic	No hope for convergence	$\ \omega_t - \omega^*\ \rightarrow \infty$

Stochasticity breaks extragradient !

Stochastic variance reduced gradient

Finite sum assumption:

$$\mathcal{L}^G(\omega) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i^G(\omega), \quad \mathcal{L}^D(\omega) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i^D(\omega)$$

Unbiased estimates of the gradient are:

$$d_i^G(\omega) := \nabla \mathcal{L}_i^G(\omega) - \nabla \mathcal{L}_i^G(\omega^S) + \mu_\theta^S$$

$$d_i^D(\omega) := \nabla \mathcal{L}_i^D(\omega) - \nabla \mathcal{L}_i^D(\omega^S) + \mu_\varphi^S$$

- μ : Full-batch gradient at the snapshot ω^S .
- Index i sampled uniformly over $\{1, \dots, n\}$.
- $\mathbb{E}[d_i^G(\omega)] = \frac{1}{n} \sum_{i=1}^n \nabla \mathcal{L}_i^G(\omega) = \nabla \mathcal{L}^G(\omega)$.
- If ω^S is close to $\omega \rightarrow$ small variance.

SVRE: Stochastic Variance Reduced Extragradient

SVRE combines **Extragradient**:

$$\omega_{t+\frac{1}{2}} = \omega_t - \gamma_t F_i(\omega_t) \quad (\text{extrapolation})$$

$$\omega_{t+1} = \omega_t - \gamma_t F_i(\omega_{t+\frac{1}{2}}) \quad (\text{update})$$

with **Variance Reduction**:

$$F_i(\omega) := \begin{pmatrix} \nabla \mathcal{L}_i^G(\omega) - \nabla \mathcal{L}_i^G(\omega^S) + \mu_\theta^S \\ \nabla \mathcal{L}_i^D(\omega) - \nabla \mathcal{L}_i^D(\omega^S) + \mu_\varphi^S \end{pmatrix}$$

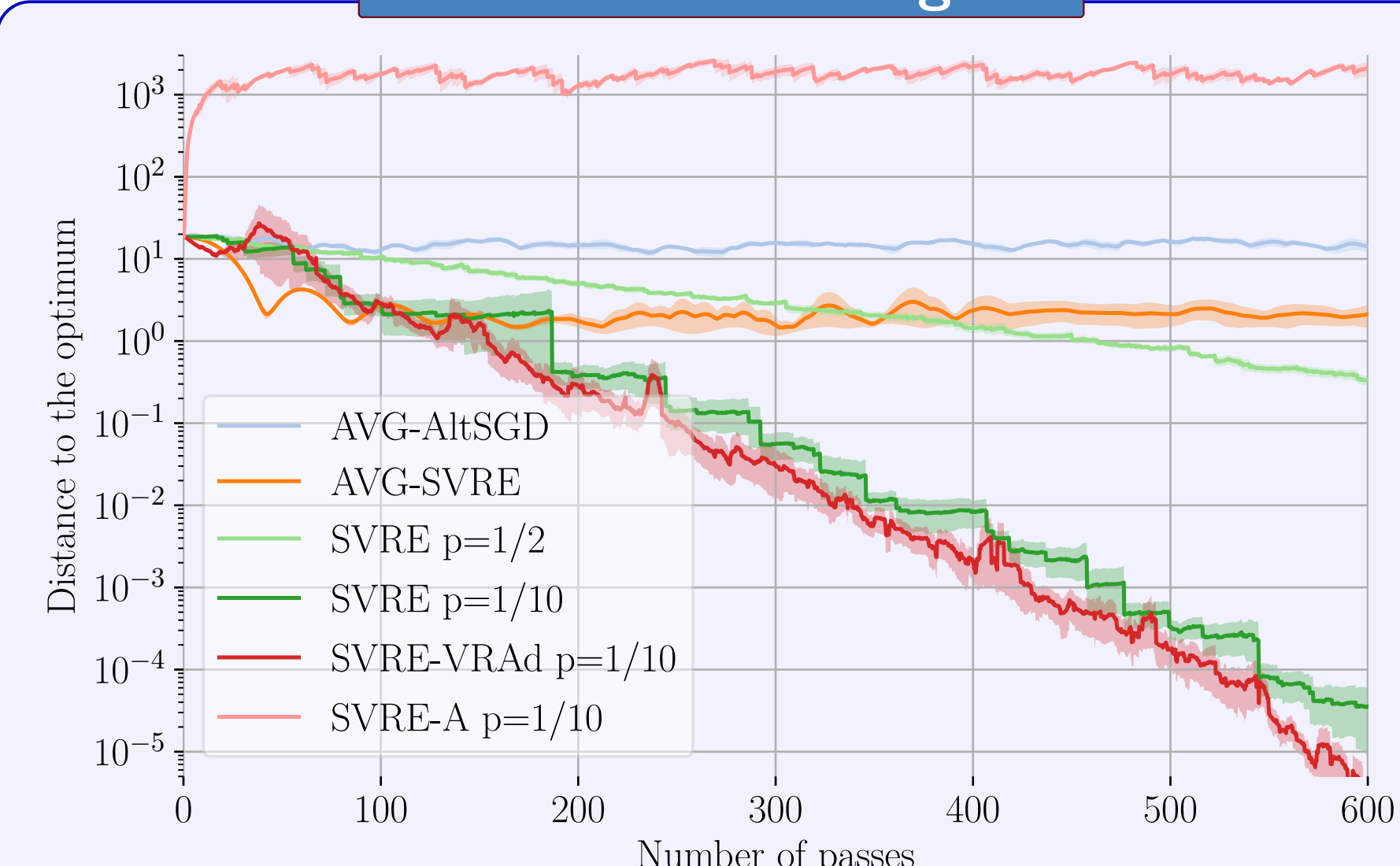
Recall $\omega := (\theta, \varphi)$ is the *joint* parameter.

Comparison of variance reduced methods for games

Method	Complexity	μ -adaptivity
SVRG	$\ln(1/\epsilon) \times (n + \frac{\bar{L}^2}{\mu^2})$	\times
A. SVRG	$\ln(1/\epsilon) \times (n + \sqrt{n} \frac{\bar{L}}{\mu})$	\times
SVRE	$\ln(1/\epsilon) \times (n + \frac{\bar{L}}{\mu})$	sometimes

- μ strong monotonicity.
- \bar{L} average Lipschitz constant of the gradient.
- $\bar{\ell}$ average cocoercivity: $\bar{L} \leq \bar{\ell} \leq \bar{L}^2/\mu$

Stochastic bilinear game



GAN Experiments

Optimization methods

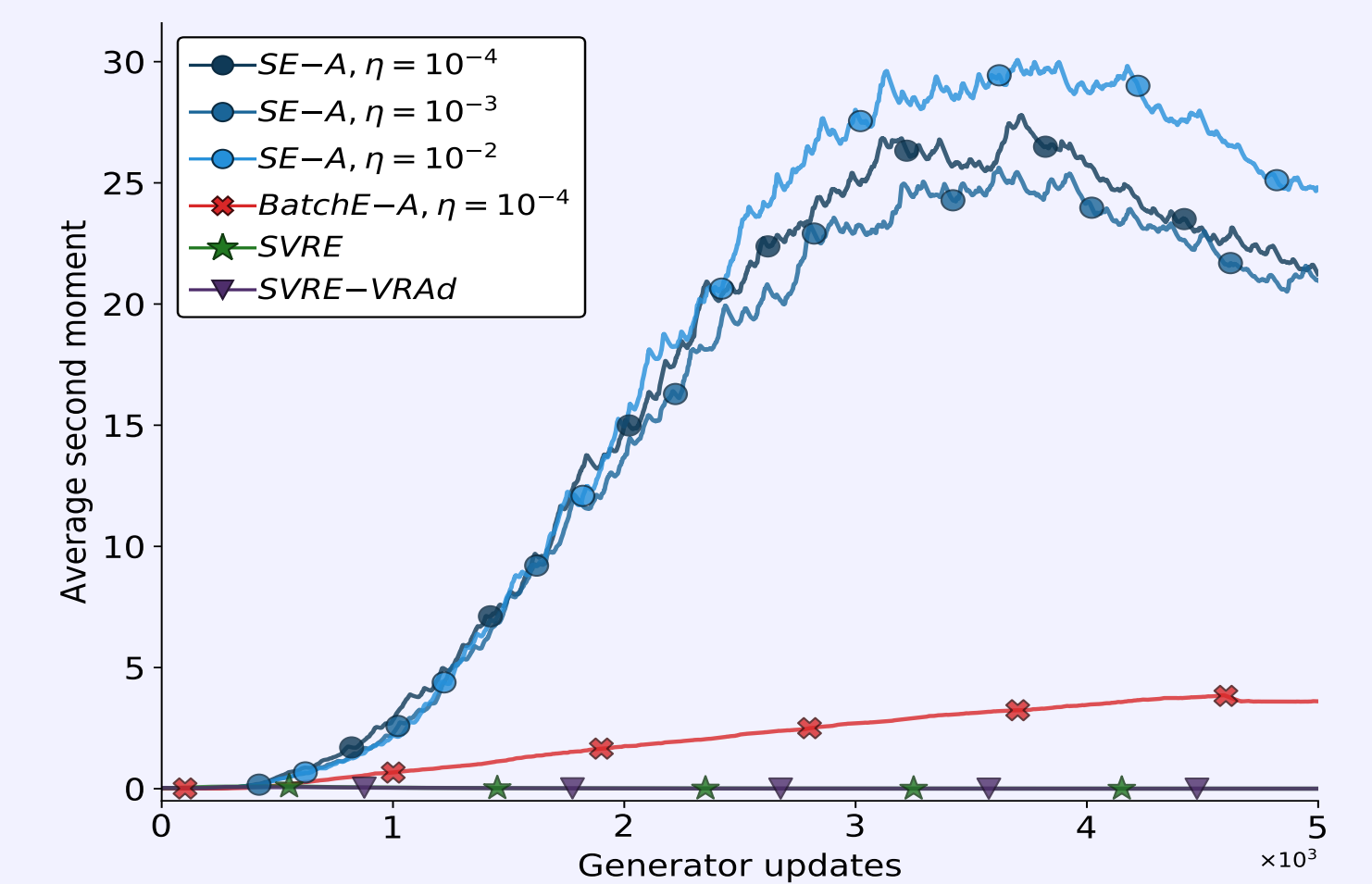
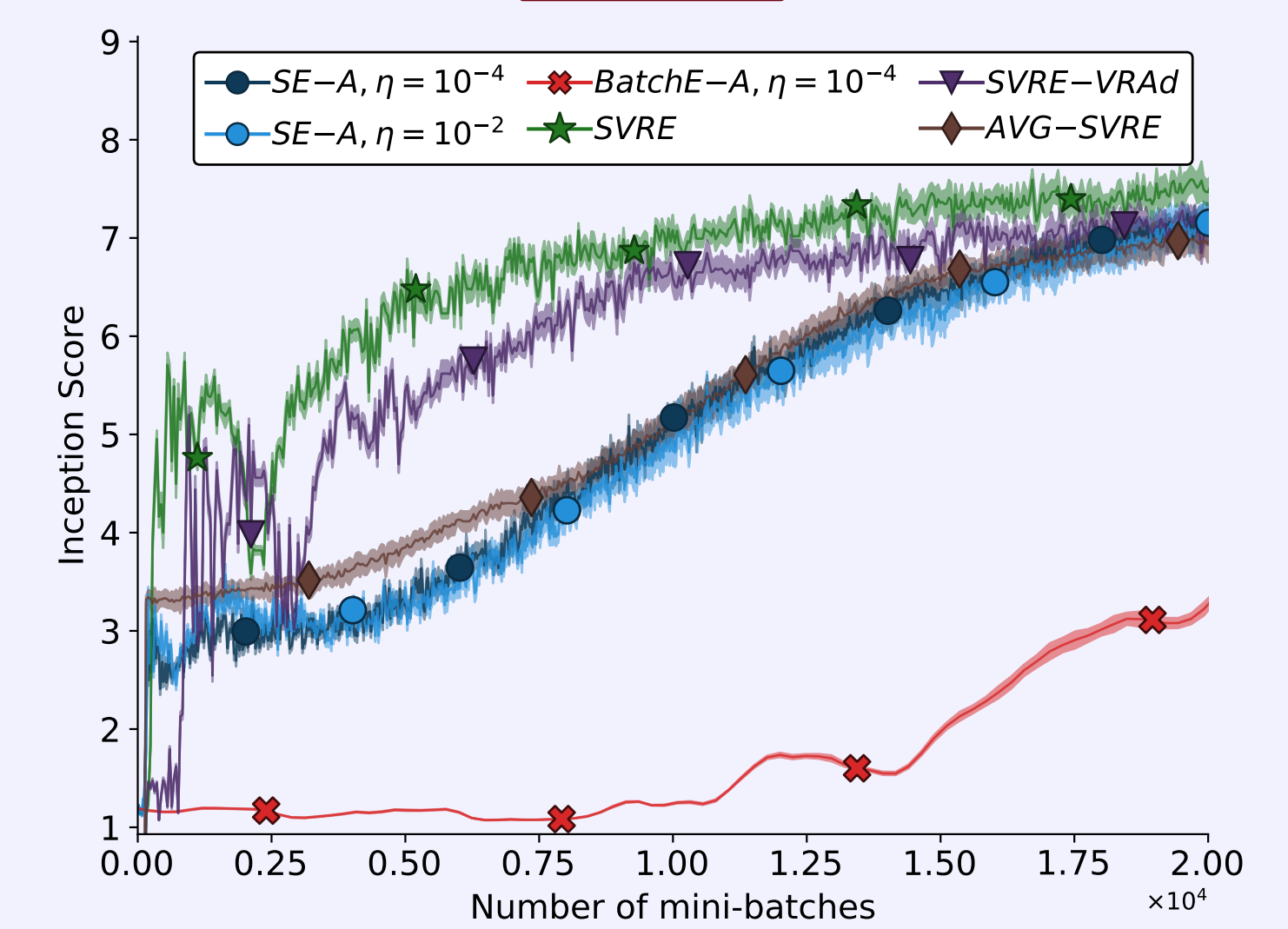
BatchE: full-batch extragradient
SG: stochastic gradient (alternating GAN)
SE: stochastic extragradient
WS-SVRE: warm-start SVRE
A: Adam (adaptive step size method)

Summary

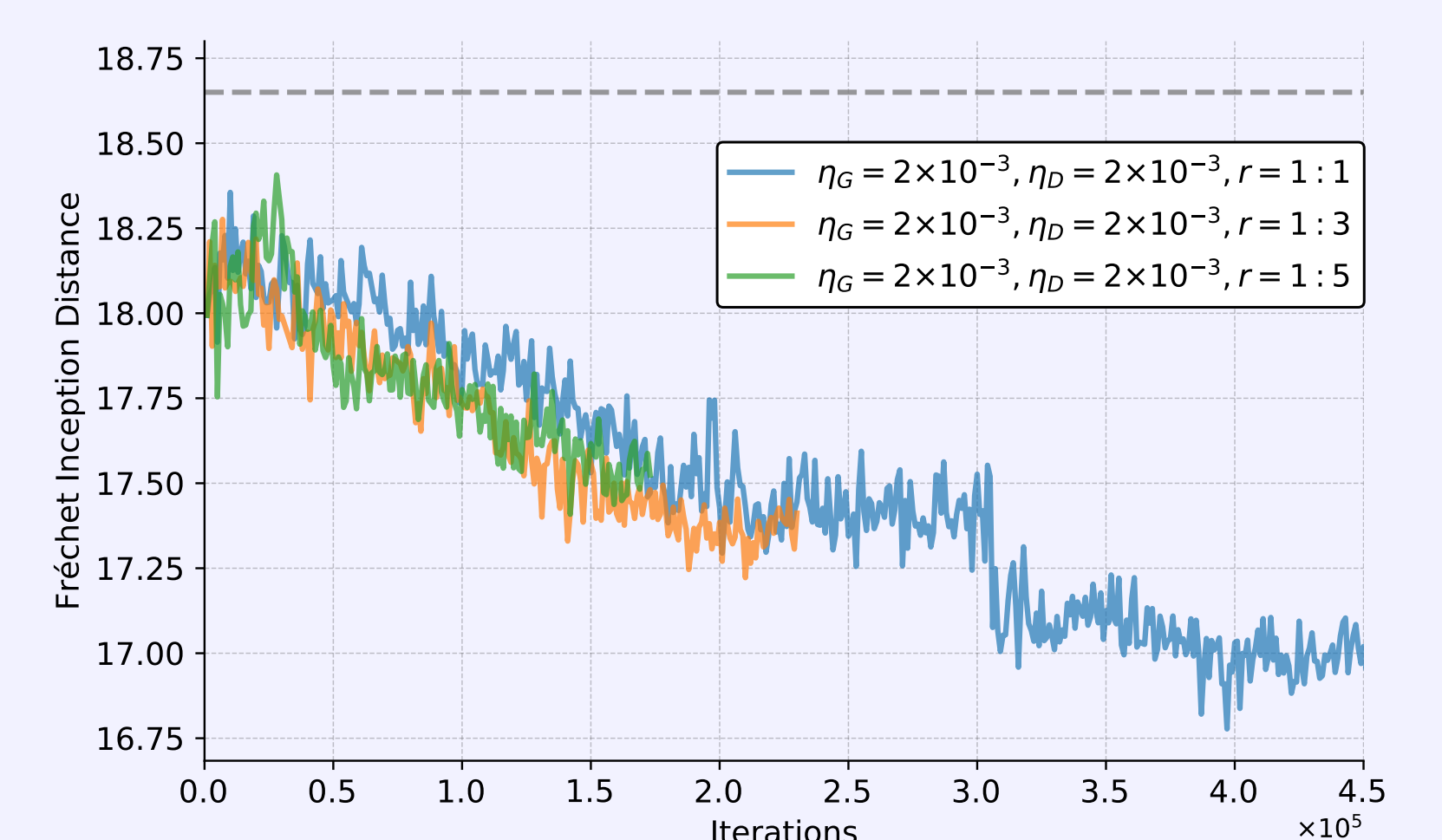
	SG-A	SE-A	SVRE	WS-SVRE
CIFAR-10	21.70	18.65	23.56	16.77
SVHN	5.66	5.14	4.81	4.88

FID (lower is better)

MNIST



WS-SVRE on CIFAR10



SVRE is more robust and stable on SVHN

