# Insights from GAN Training with Kernel Discriminators
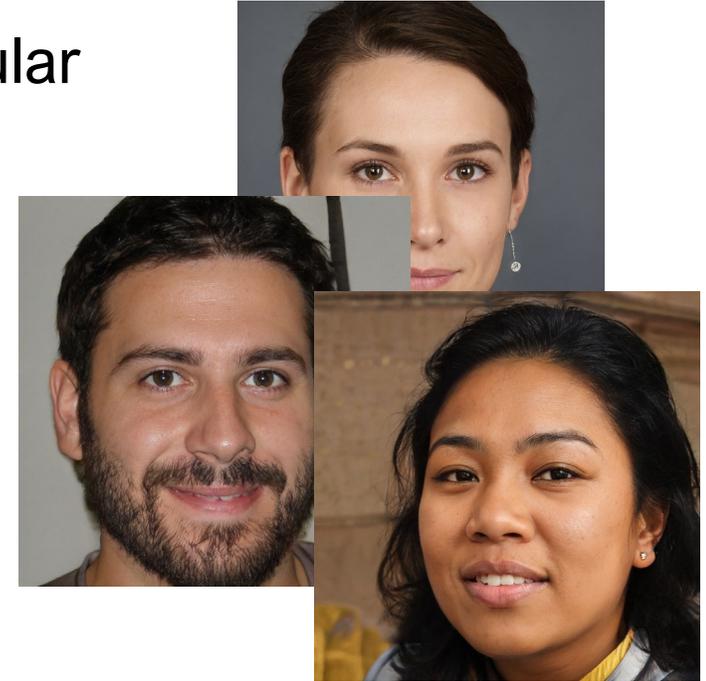
Evan Becker

01/25/2024

# Motivation

**Problem:** Generative Adversarial Networks (GANs) are popular but hard to train

- Alternating min-max optimization for a non-convex objective is not fully understood

**Contribution:** Simple yet expressive framework to analyze convergence and failure modes

- MMD-GAN objective trained with gradient descent ascent
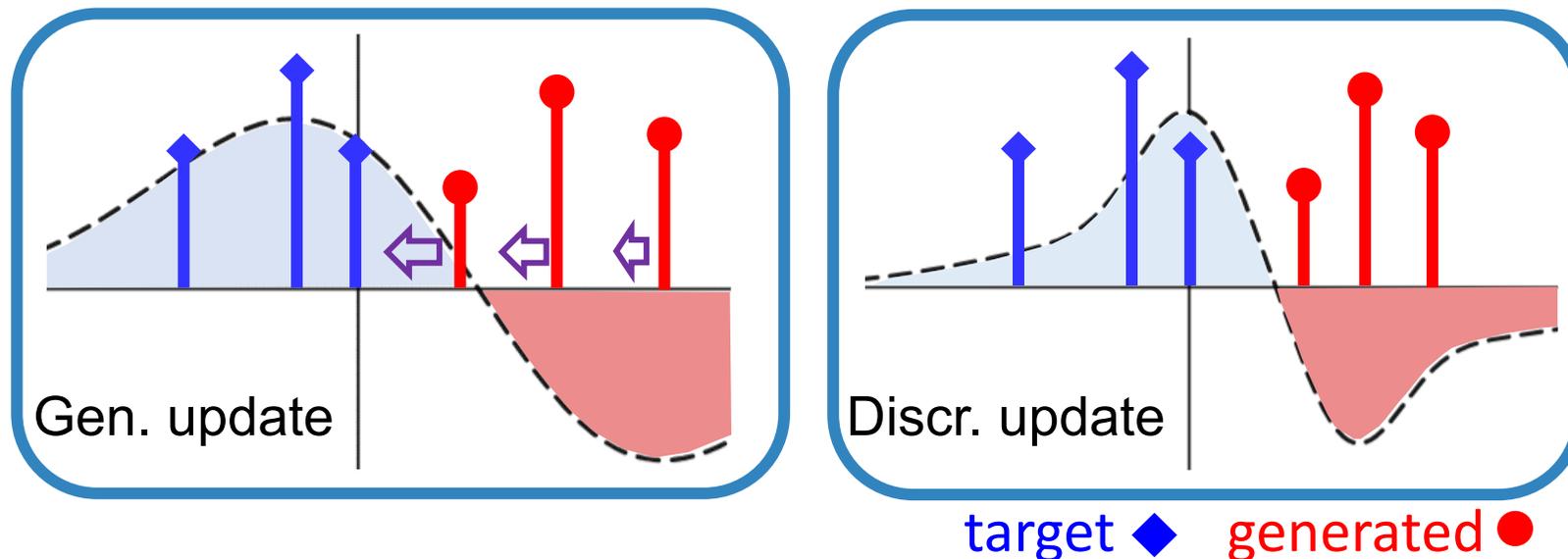- Generator is unconstrained, discriminator is a kernel model

*thispersondoesnotexist.com*

# Problem setup

**Target and Generated Distributions:** Directly parameterized by set of points ($\delta$ Dirac delta function):

$$\mathbb{P}_r(\boldsymbol{x}) = \sum_{i=1}^{N_r} p_i \delta(\boldsymbol{x} - \boldsymbol{x}_i), \quad \mathbb{P}_g(\boldsymbol{x}) = \sum_{j=1}^{N_g} \widetilde{p}_j \delta(\boldsymbol{x} - \widetilde{\boldsymbol{x}}_j),$$

**Discriminator:** $f : \mathcal{X} \to \mathbb{R}$ from a reproducing kernel hilbert space (RKHS) $\mathcal{H}_K$ with positive definite kernel function $K : \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$.
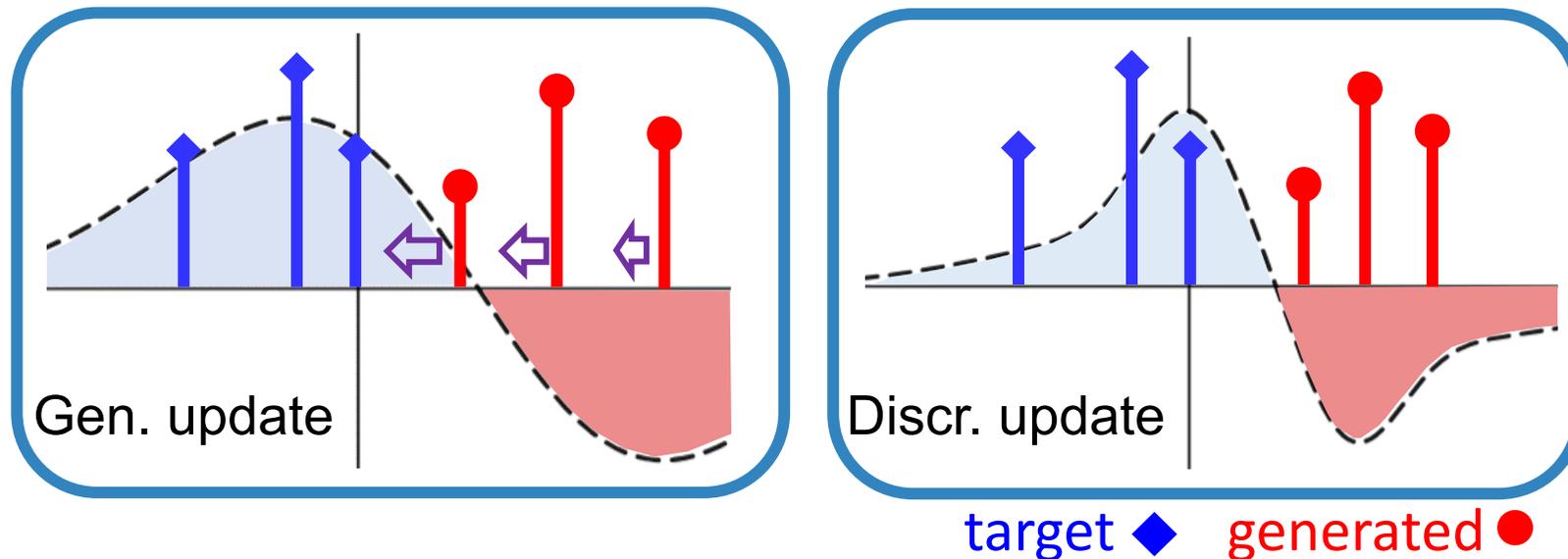


Gen. update

Discr. update

target ◆   generated ●

# Problem setup

**Optimization:** $\min\limits_{\widetilde{X}} \max\limits_{f \in \mathcal{H}} \mathcal{L}(f, \widetilde{X})$ with loss defined as
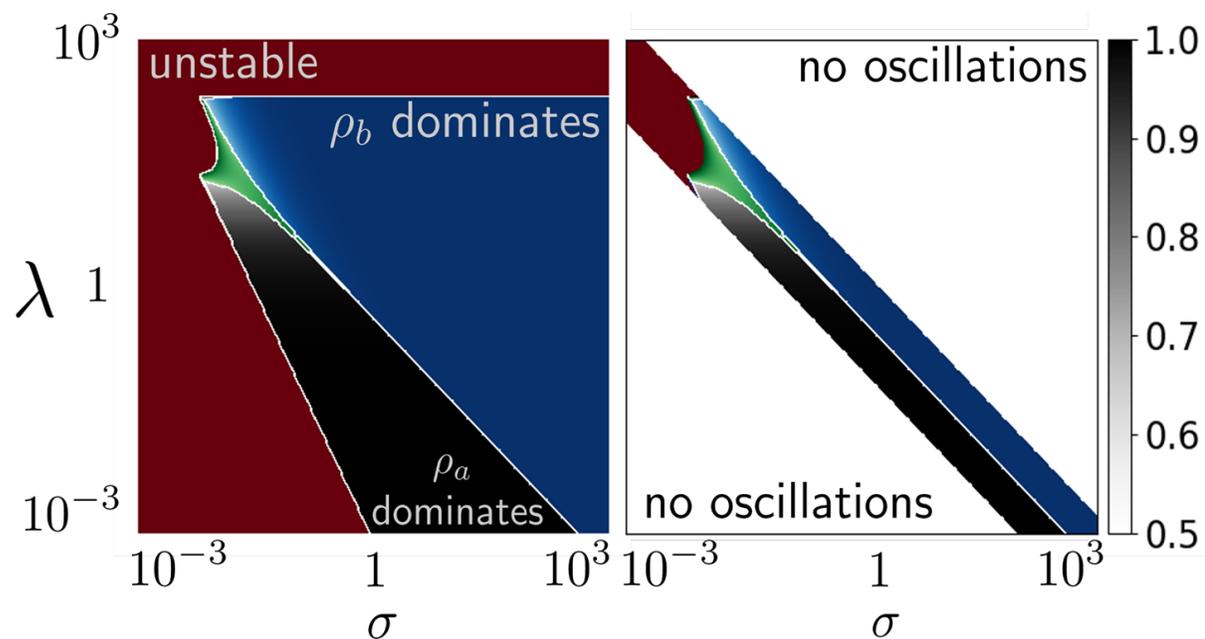
$$\mathcal{L}(f, \widetilde{X}) := \sum_{i=1}^{N_r} p_i f(x_i) - \sum_{i=1}^{N_g} \widetilde{p}_i f(\widetilde{x}_j) - \tfrac{\lambda}{2} \|f\|_{\mathcal{H}}^2 \,.$$

*Note*: Just maximizing the objective over $f \in \mathcal{H}$ results in the maximum mean discrepancy (MMD) between distributions, which is a common two-sample test statistic

# First-Order Analysis

**Idea:** Look at some local region around each true point $x_i$



*Dominating eigenvalues when using an RBF kernel discriminator*

**Insight 1** (Becker et al. 22): Existence of good local minima (exact convergence) and bad (mode collapse).
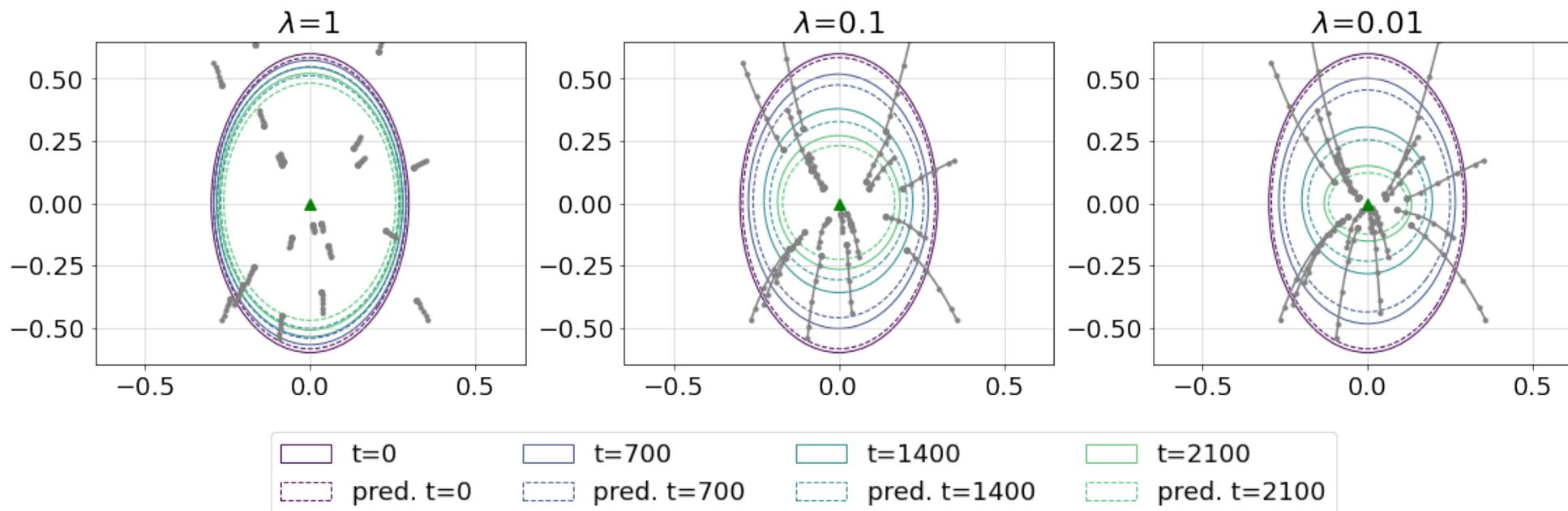
**Insight 2** (Becker et al. 23): Analyze the *rate* of convergence by looking at eigenvalue functions. Can prescribe hyperparameters to achieve fastest local convergence.

*hyperparameters*: $\lambda$ regularization, $\sigma$ kernel width, $\eta_g, \eta_d$ learning rates, $\Delta_i := p_i - \sum_{j \in N_i} \widetilde{p}_j$ local probability mass difference

# Second-Order Analysis

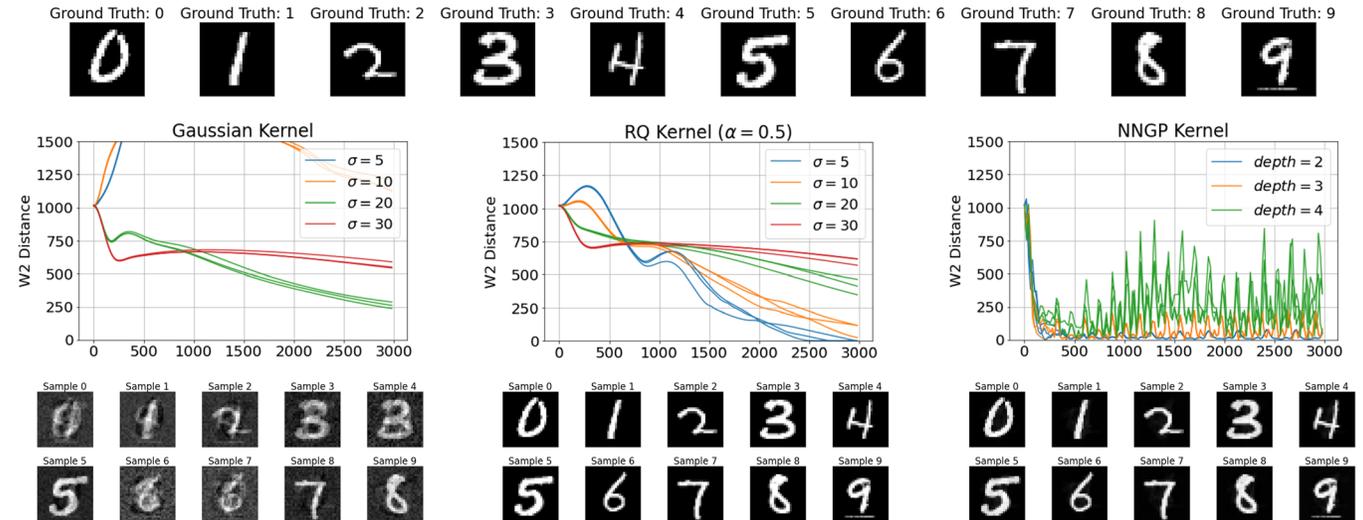**Idea:** Look at the first two moments (mean and variance) of the generated distribution



Predicted vs. Actual Covariance ($\sigma = 1$)

**Insight 3**: Smaller kernel width $\sigma$ increases 'momentum' of convergence in second order analysis (stronger effect than linearized analysis predicts)

*hyperparameters*: $\lambda$ regularization, $\sigma$ kernel width, $\eta_g, \eta_d$ learning rates, $\Delta_i := p_i - \sum_{j \in N_i} \widetilde{p}_j$ local probability mass difference

# Key takeaways:



## Why do we care?

- Neural networks can be thought of as kernel machines whose kernel shape evolves over time (evolves very little in NTK regime)

- Reducing effective kernel width of the discriminator during training promotes fast convergence to good minima, and is already done heuristically (Karras et. al 2018)!

- Min-max optimization is used outside of GANs, and techniques used here could be used to study other systems

# Thank you!

Questions?

Email: evbecker@ucla.edu