

Maximum Mean Discrepancy Gradient Flow

Michael Arbel¹, Anna Korba¹, Adil Salim² and Arthur Gretton¹

¹Gatsby Computational Neuroscience Unit, University College London ²KAUST

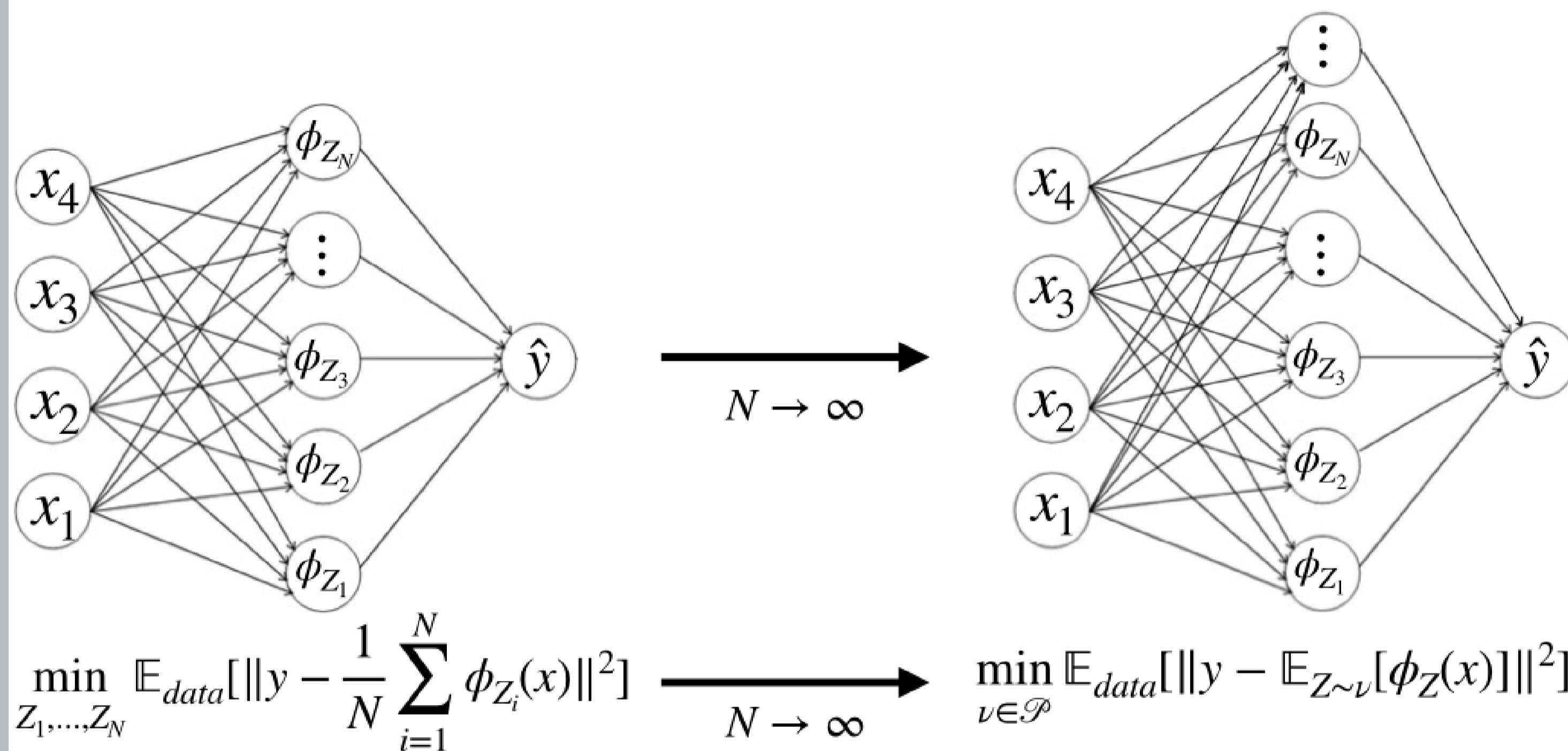


Overview

- General setting:**
- ✓ Non-convex optimization in probability space with the Maximum Mean Discrepancy as a cost function.
 - ✓ Interested in Gradient descent dynamics in the limit of large samples $N \rightarrow \infty$.
- Goals:**
- ✓ Criterion for global convergence of gradient descent when N approaches infinity.
 - ✓ New algorithm based on noise-injection to improve convergence.
 - ✓ Application 1: Optimization of neural networks.
 - ✓ Application 2: Criterion to characterise convergence in Implicit Generative models.

Motivation 1: Gradient Descent dynamics in neural networks

Easier to characterize the gradient descent dynamics in the Mean-field limit [3, 7]

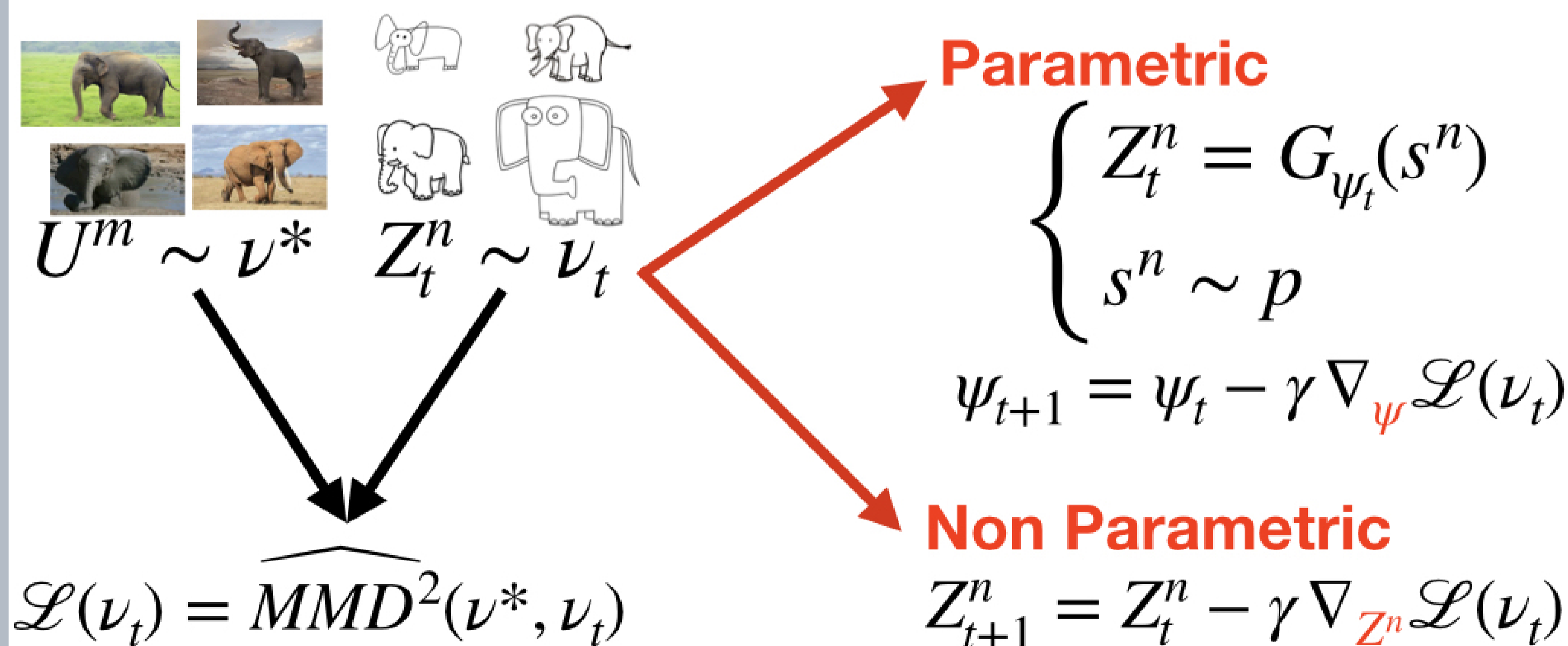


In the well-specified case, i.e.: $\mathbb{E}_{data}[y|x] = \mathbb{E}_{U \sim \nu^*} [\phi_U(x)]$, equivalent to minimizing the MMD with a random feature kernel k :

$$\min_{\nu \in \mathcal{P}} MMD_k^2(\nu^*, \nu), \quad k(Z, Z') = \mathbb{E}_{data} [\phi_Z(x)^\top \phi_{Z'}(x)]$$

Motivation 2: Implicit Generative models

- Good performance for Implicit generative models using the MMD as a loss [5, 2, 1].
- Hard to characterize global convergence.



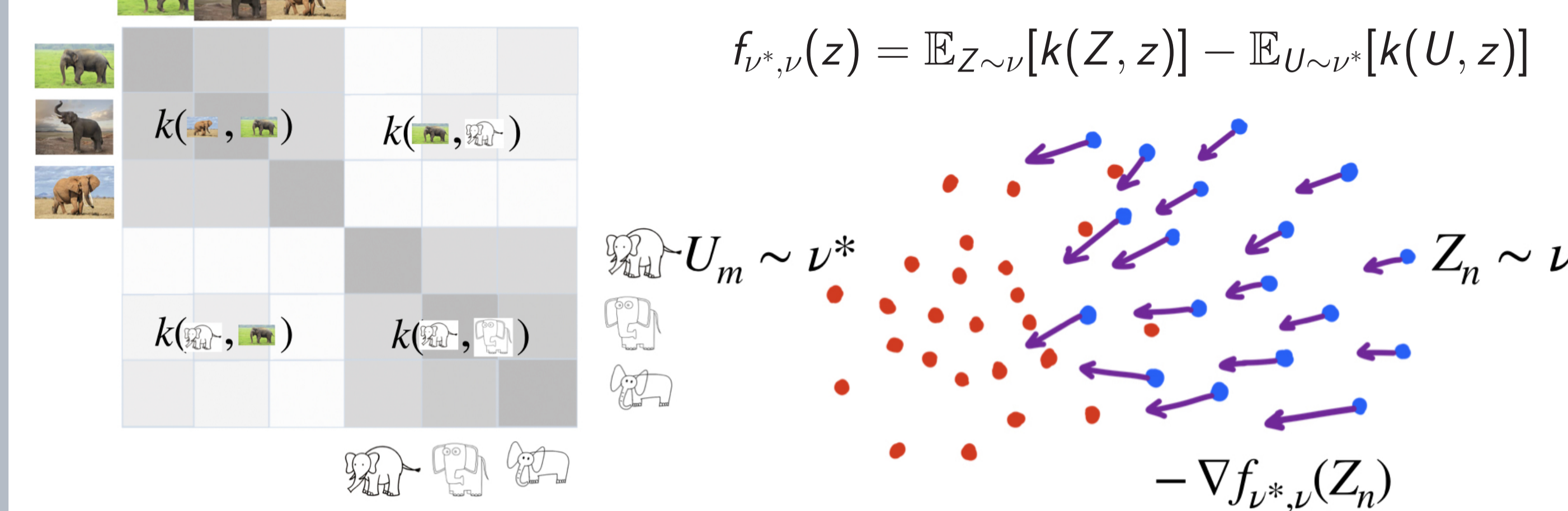
- Easier to analyse the non-parametric setting.

Maximum Mean Discrepancy (MMD)

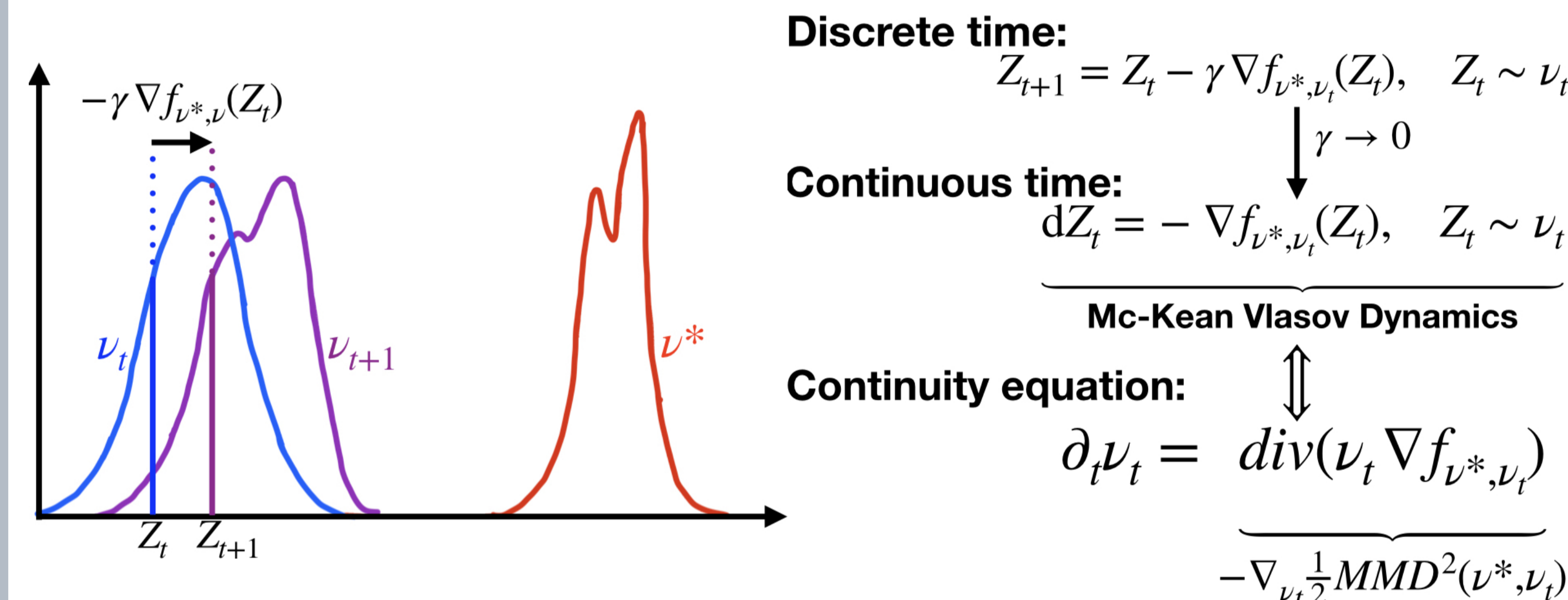
The MMD is a distance between probability distributions defined using a positive semi-definite kernel k [4]:

$$\frac{1}{2} MMD_k^2(\nu^*, \nu) = \frac{1}{2} \mathbb{E}_{Z, Z' \sim \nu} [k(Z, Z')] - \mathbb{E}_{Z \sim \nu^*} [k(Z, U)] + \frac{1}{2} \mathbb{E}_{U, U' \sim \nu^*} [k(U, U')]$$

- ✓ Easy to estimate from samples: ✓ Can be interpreted as the energy relative to a potential function $f_{\nu^*, \nu} := \frac{\delta}{\delta \nu} \frac{1}{2} MMD_k^2(\nu^*, \nu)$

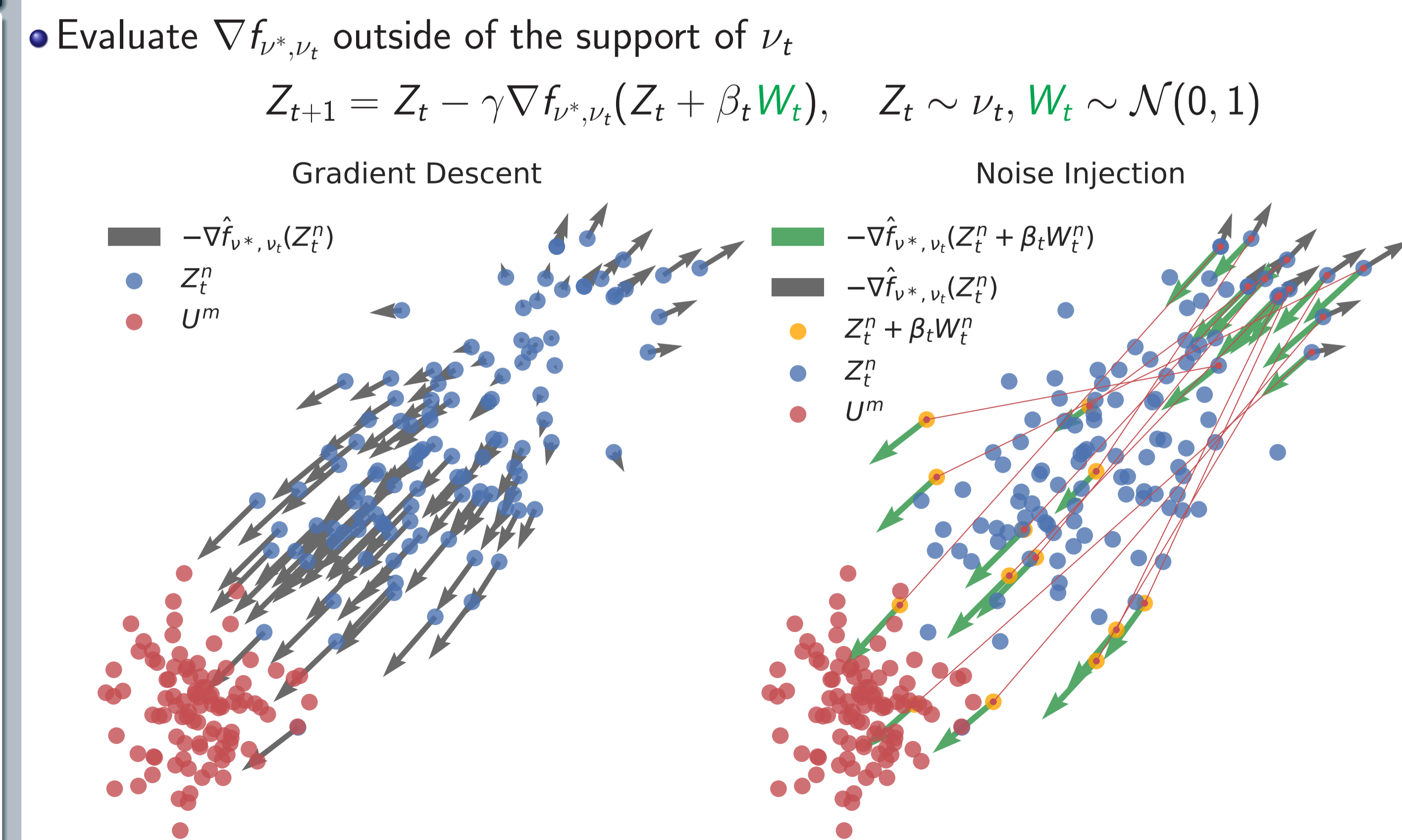


Gradient Flow of the MMD



- Equivalent to Gradient descent when ν_t restricted to the form: $\nu_t = \frac{1}{N} \sum_{n=1}^N \delta_{Z_t^n}$.
- Stationarity distribution ν_∞ satisfies: $\mathbb{E}_{Z \sim \nu_\infty} [\|\nabla f_{\nu^*, \nu_\infty}(Z)\|^2] = 0$

Noise injection (NI)



Theory: Global convergence

1. Criterion for convergence of the gradient flow: *Negative Sobolev Distance*:

$$S(\nu^* | \nu_t) := \sup_{\substack{g \in \mathcal{L}(\nu_t) \\ \|\nabla g\|_{L_2(\nu_t)} \leq 1}} \mathbb{E}_{Z \sim \nu^*} [g(Z)] - \mathbb{E}_{Z \sim \nu_t} [g(Z)]$$

Assume that $S(\nu^* | \nu_t) \leq C$ for all t , and that k is a characteristic kernel, then ν_t converges weakly towards ν^* . Moreover:

$$MMD^2(\nu^*, \nu_t)^2 \leq \frac{1}{MMD^2(\nu^*, \nu_0) + 4\gamma C^{-1}t}$$

2. Criterion for convergence of the noise injection algorithm:

- Decreasing direction:

$$4\gamma^2 \beta_t^2 MMD^2(\nu^*, \nu_t) \leq \mathbb{E}_{\substack{Z \sim \nu_t \\ W \sim \mathcal{N}(0,1)}} [\|\nabla f_{\nu^*, \nu_t}(Z + \beta_t W)\|^2]$$

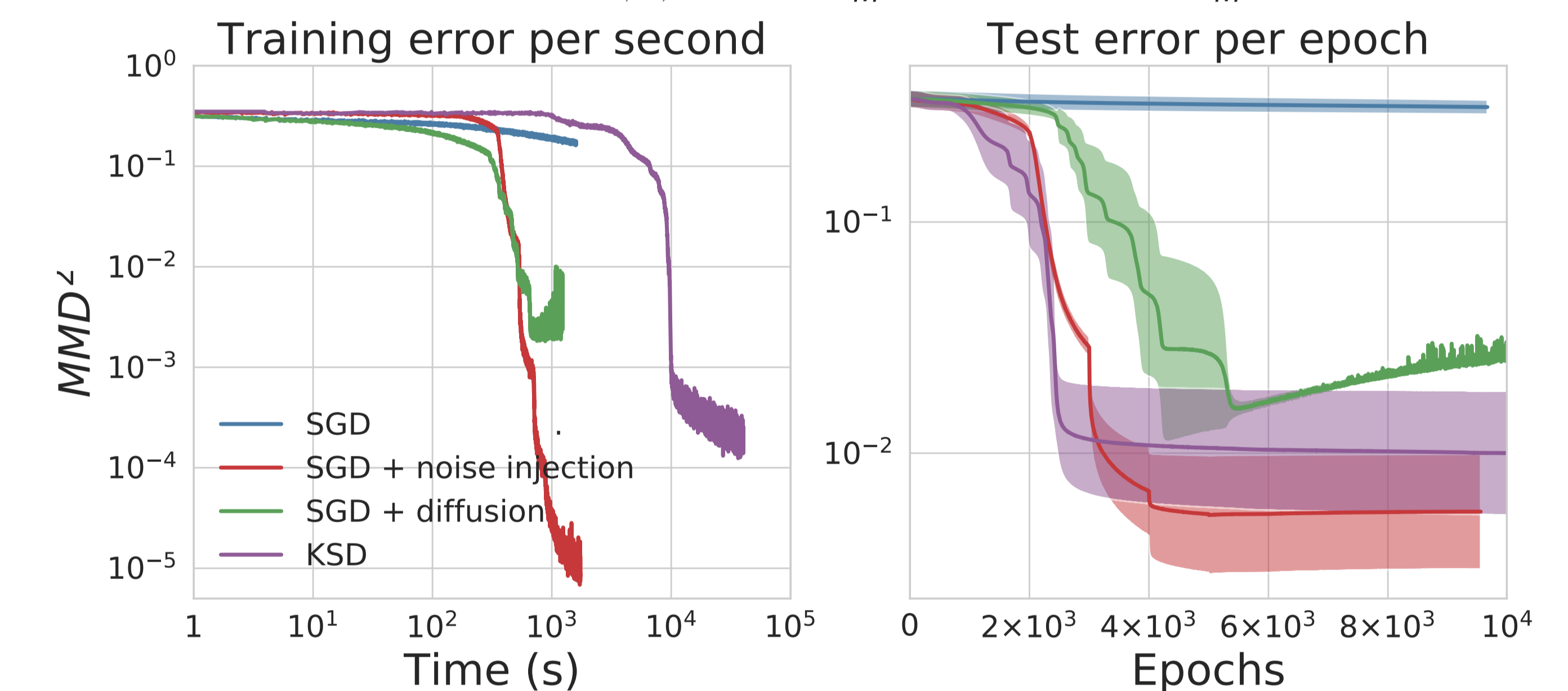
- Large noise: $\sum_{i=0}^t \beta_i^2 \rightarrow \infty$

Then for some constant L :

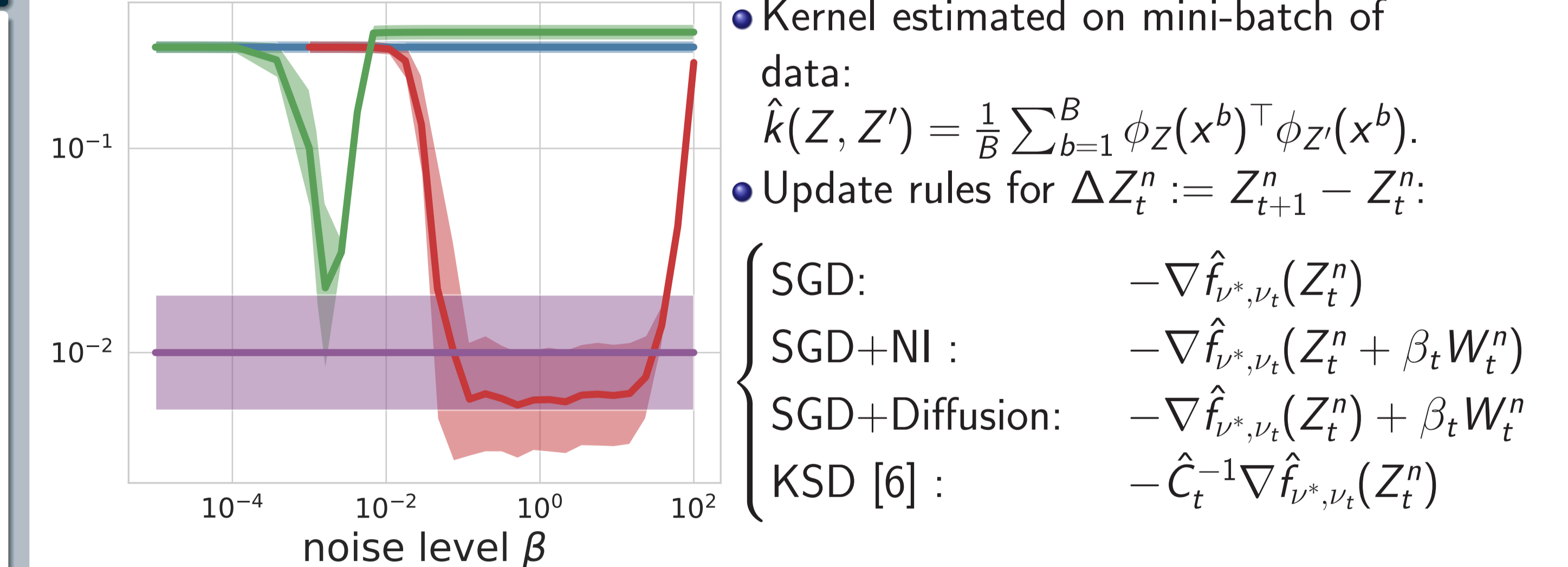
$$MMD^2(\nu^*, \nu_t) \leq MMD(\nu^*, \nu_0) e^{-4\gamma^2(1-3\gamma L) \sum_{i=0}^t \beta_i^2}$$

Experimental Comparison

Student-Teacher networks: $\min_{Z^1, \dots, Z^N} \mathbb{E}_{data} [\|\frac{1}{M} \sum_{m=1}^M \phi_{U^m}(x) - \frac{1}{M} \sum_{n=1}^N \phi_{Z^n}(x)\|^2]$



Sensitivity to noise (Test error)



Bibliography

- M. Arbel, D. J. Sutherland, M. Birnkowski, and A. Gretton. "On gradient regularizers for MMD GANs". *NIPS* (2018).
- M. Birnkowski, D. J. Sutherland, M. Arbel, and A. Gretton. "Demystifying MMD GANs". *ICLR*. 2018.
- L. Chizat and F. Bach. "On the global convergence of gradient descent for over-parameterized models using optimal transport". *NIPS*, 2018.
- A. Gretton, K. M. Borgwardt, M. J. Rasch, B. Schölkopf, and A. Smola. "A kernel two-sample test". *Journal of Machine Learning Research* (2012).
- Y. Li, K. Swersky, and R. Zemel. "Generative Moment Matching Networks". *ICML*. 2015.
- Y. Mroueh, T. Sercu, and A. Raj. *Sobolev Descent: Variational Transport of Distributions via Advection*. Private communication. Apr. 2018.
- G. M. Rotskoff and E. Vanden-Eijnden. "Neural networks as interacting particle systems: Asymptotic convexity of the loss landscape and universal scaling of the approximation error". *arXiv preprint arXiv:1805.00915* (2018).