

Flowing Datasets with Wasserstein over Wasserstein Gradient Flows

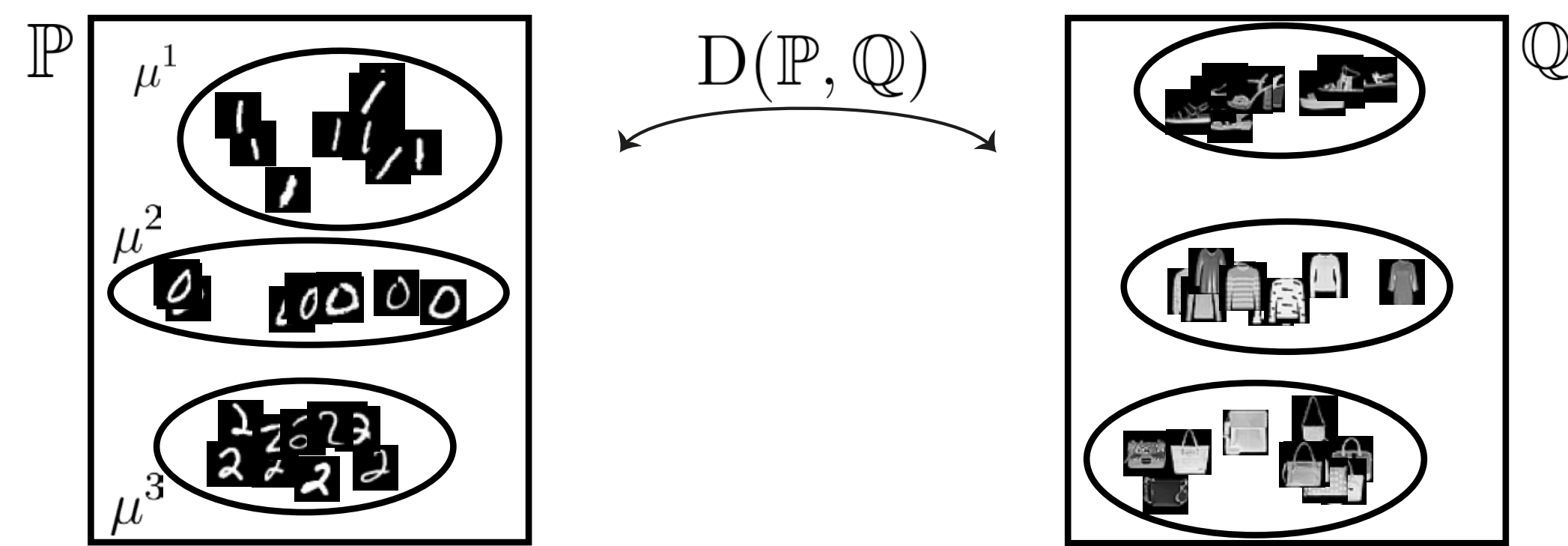
Clément Bonet^{*1}, Christophe Vauthier^{*2}, Anna Korba¹

¹ENSAE, CREST, Institut Polytechnique de Paris; ²Université Paris-Saclay, Laboratoire de Mathématique d'Orsay

Contributions

Goal: move labeled dataset in a coherent way

- Labeled datasets modeled as $\mathbb{P} = \frac{1}{C} \sum_{c=1}^C \delta_{\mu^{c,n}} \in \mathcal{P}_2(\mathcal{P}_2(\mathbb{R}^d))$ where $\mu^{c,n} = \frac{1}{n} \sum_{i=1}^n \delta_{x_i^c}$
- Endow $\mathcal{P}_2(\mathcal{P}_2(\mathbb{R}^d))$ with OT distance WoW
- Minimize $\mathbb{F} : \mathcal{P}_2(\mathcal{P}_2(\mathbb{R}^d)) \rightarrow \mathbb{R}$ using WoW gradient flows
- Application on image datasets



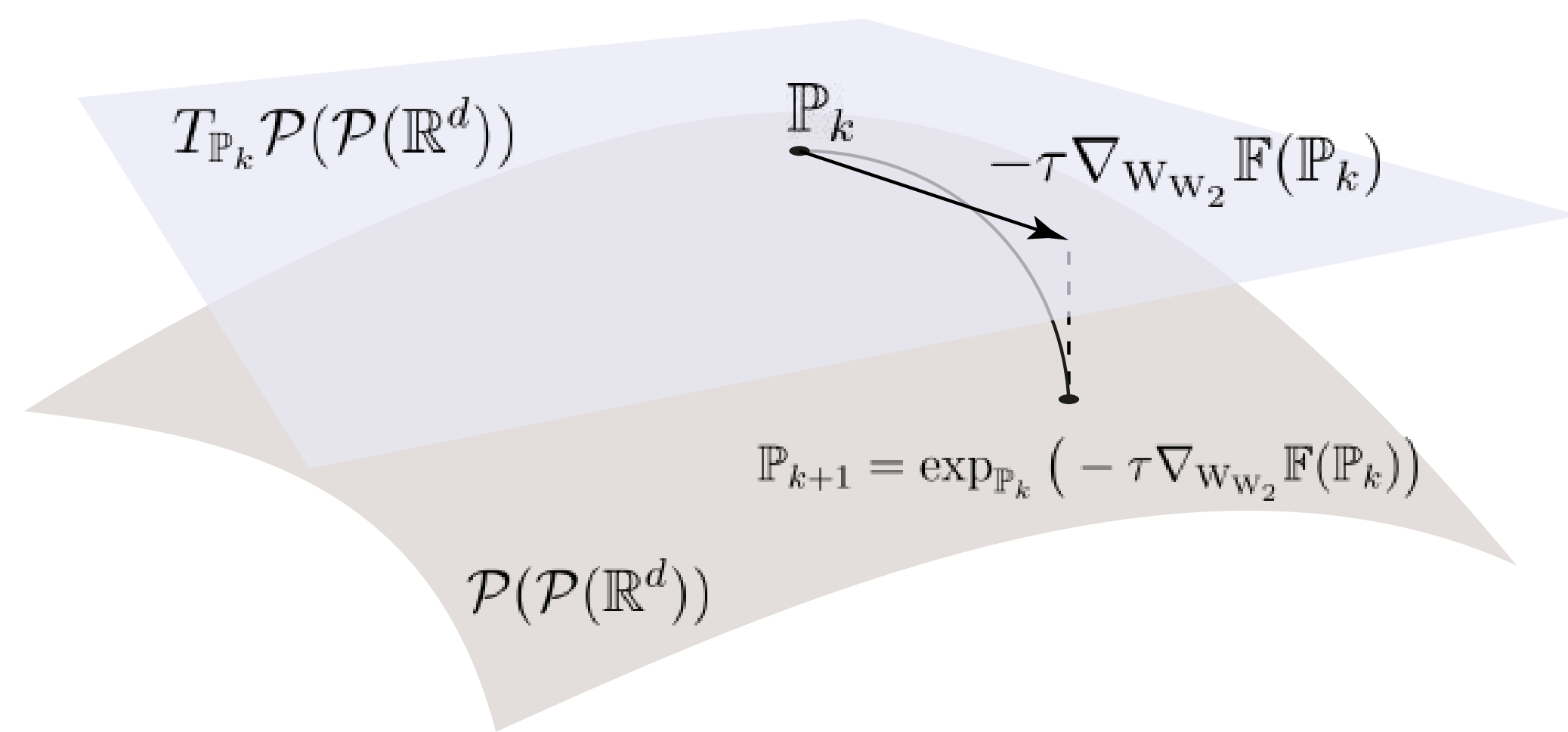
Wasserstein over Wasserstein Space

WoW distance: Let $\mathbb{P}, \mathbb{Q} \in \mathcal{P}_2(\mathcal{P}_2(\mathcal{M}))$,

$$W_{W_2}(\mathbb{P}, \mathbb{Q})^2 = \inf_{\Gamma \in \Pi(\mathbb{P}, \mathbb{Q})} \int W_2^2(\mu, \nu) d\Gamma(\mu, \nu)$$

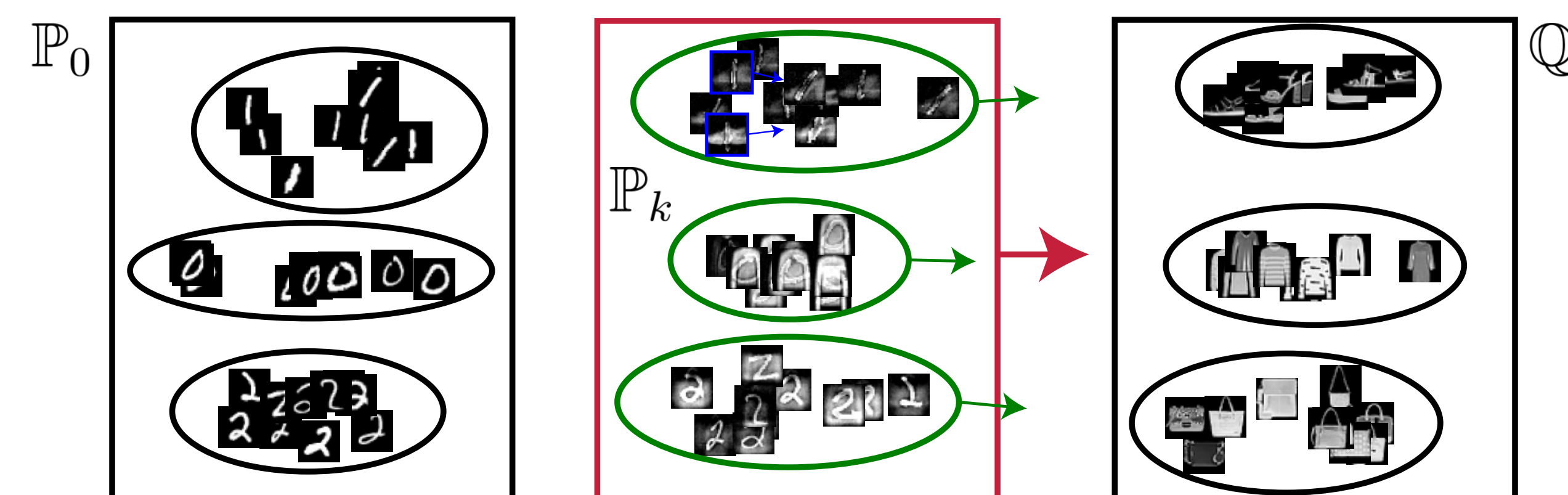
→ Riemannian structure

WoW Gradient Descent: $\mathbb{P}_{k+1} = \exp_{\mathbb{P}_k}(-\tau \nabla_{W_{W_2}} \mathbb{F}(\mathbb{P}_k))$



In practice: For $\mathcal{M} = \mathbb{R}^d$, $\mathbb{P}_k = \frac{1}{C} \sum_{c=1}^C \delta_{\mu_k^{c,n}}$

$$\forall k \geq 0, x_{i,k+1}^c = x_{i,k}^c - \tau \nabla_{W_{W_2}} \mathbb{F}(\mathbb{P}_k)(\mu_k^{c,n})(x_{i,k}^c)$$



Minimization of the MMD on $\mathcal{P}_2(\mathcal{P}_2(\mathbb{R}^d))$

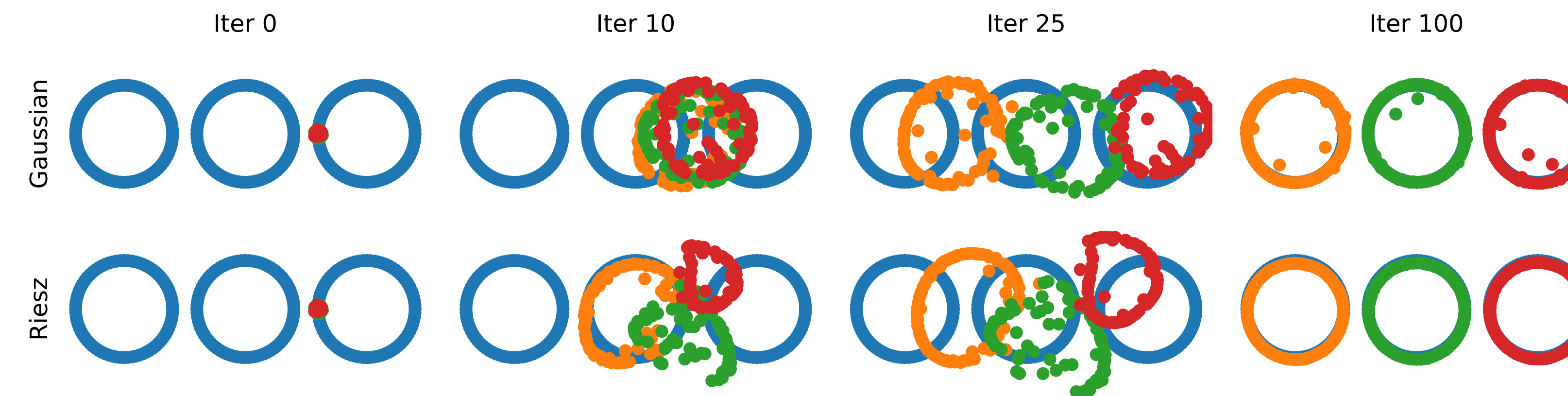
Goal: minimize $\mathbb{F}(\mathbb{P}) = \frac{1}{2} \text{MMD}_K^2(\mathbb{P}, \mathbb{Q}) = \mathbb{V}(\mathbb{P}) + \mathbb{W}(\mathbb{P}) + \text{cst}$, with $\mathbb{V}(\mathbb{P}) = \int \mathcal{V}(\mu) d\mathbb{P}(\mu)$, $\mathcal{V}(\mu) = -\int K(\mu, \nu) d\mathbb{Q}(\nu)$, $\mathbb{W}(\mathbb{P}) = \frac{1}{2} \iint K(\mu, \nu) d\mathbb{P}(\mu) d\mathbb{P}(\nu)$

SW distance: $\text{SW}_2^2(\mu, \nu) = \int_{S^{d-1}} W_2^2(P_\#^\theta \mu, P_\#^\theta \nu) d\sigma(\theta)$, $P^\theta(x) = \langle x, \theta \rangle$

Kernel: $K(\mu, \nu) = e^{-\frac{1}{2h} \text{SW}_2^2(\mu, \nu)}$ (Gaussian) or $K(\mu, \nu) = -\text{SW}_2(\mu, \nu)$ (Riesz)

Computational complexity: $O(C^2 L n \log n)$

$\nabla_{W_{W_2}} \mathbb{F}(\mathbb{P})(\mu^{c,n})(x_i^c) = n C \nabla_{i,c} F(\mathbf{x})$ for $\mathbf{x} = (x_i^c)_{i,c}$: obtained in closed-form or by auto-differentiation of $F(\mathbf{x}) := \mathbb{F}(\mathbb{P})$



Wasserstein over Wasserstein Gradients

For $(x, v) \in T\mathcal{M}$, define $\pi^{\mathcal{M}}((x, v)) = x$.

Couplings. For any $\gamma \in \mathcal{P}_2(T\mathcal{M})$, let $\phi^{\mathcal{M}}(\gamma) = \pi_{\#}^{\mathcal{M}} \gamma$, $\phi^{\text{exp}}(\gamma) = \exp_{\#} \gamma$.

$$\exp_{\mathbb{P}}^{-1}(\mathbb{Q}) := \{\tilde{\Gamma} \in \mathcal{P}_2(\mathcal{P}_2(T\mathcal{M})) \mid \phi_{\#}^{\mathcal{M}} \tilde{\Gamma} = \mathbb{P}, \phi^{\text{exp}} \tilde{\Gamma} = \mathbb{Q}, \iint \|v\|_x^2 d\gamma(x, v) d\tilde{\Gamma}(\gamma) = W_{W_2}^2(\mathbb{P}, \mathbb{Q})\}.$$

WoW Gradient

Let $\mathbb{F} : \mathcal{P}_2(\mathcal{P}_2(\mathcal{M})) \rightarrow \mathbb{R}$. \mathbb{F} is WoW differentiable at \mathbb{P} if there exists $\nabla_{W_{W_2}} \mathbb{F}(\mathbb{P}) : \mathcal{P}_2(\mathcal{M}) \rightarrow T\mathcal{P}_2(\mathcal{M})$ s.t. for any $\mathbb{Q} \in \mathcal{P}_2(\mathcal{P}_2(\mathcal{M}))$, $\tilde{\Gamma} \in \exp_{\mathbb{P}}^{-1}(\mathbb{Q})$, $\mathbb{F}(\mathbb{Q}) = \mathbb{F}(\mathbb{P}) + \iint \langle \nabla_{W_{W_2}} \mathbb{F}(\mathbb{P})(\pi_{\#}^{\mathcal{M}} \gamma)(x), v \rangle_x d\gamma(x, v) d\tilde{\Gamma}(\gamma) + o(W_{W_2}(\mathbb{P}, \mathbb{Q}))$.

Potentials: $\mathbb{V}(\mathbb{P}) = \int \mathcal{F}(\mu) d\mathbb{P}(\mu)$, $\nabla_{W_{W_2}} \mathbb{V}(\mathbb{P}) = \nabla_{W_2} \mathcal{F}$

Interactions: $\mathbb{W}(\mathbb{P}) = \iint \mathcal{W}(\mu, \nu) d\mathbb{P}(\mu) d\mathbb{P}(\nu)$,

$$\nabla_{W_{W_2}} \mathbb{W}(\mathbb{P})(\mu) = \int (\nabla_{W_2,1} \mathcal{W}(\mu, \nu) + \nabla_{W_2,2} \mathcal{W}(\mu, \nu)) d\mathbb{P}(\nu)$$

For $K_\nu(\mu) = K(\mu, \nu) = e^{-\frac{1}{2h} \text{SW}_2^2(\mu, \nu)}$, $\nabla_{W_{W_2}} \mathbb{F}(\mathbb{P})(\mu) = \int \nabla_{W_2} K_\nu(\mu) d(\mathbb{P} - \mathbb{Q})(\nu)$, $\nabla_{W_2} K_\nu(\mu) = -\frac{1}{h} e^{-\frac{1}{2h} \text{SW}_2^2(\mu, \nu)} \int_{S^{d-1}} \psi'_\theta(\langle x, \theta \rangle) \theta d\sigma(\theta)$, $\psi'_\theta(u) = u - F_{P_\#^\theta \nu}^{-1}(F_{P_\#^\theta \mu}(u))$.

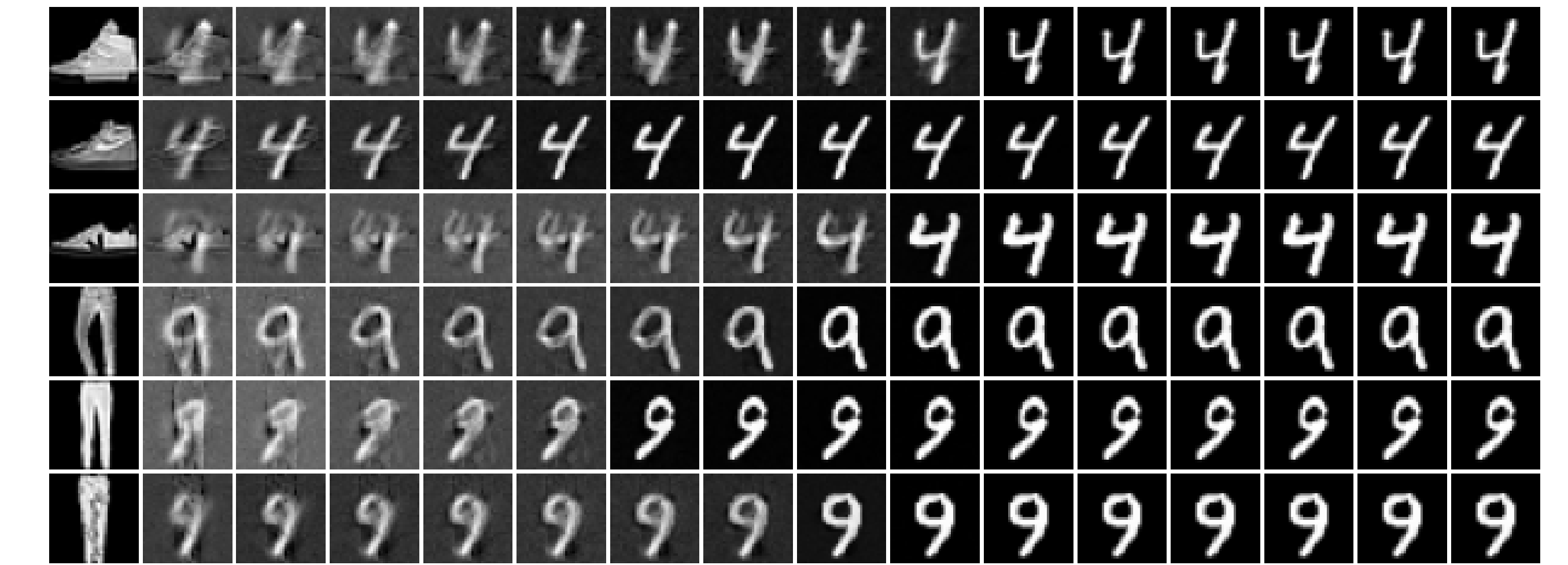
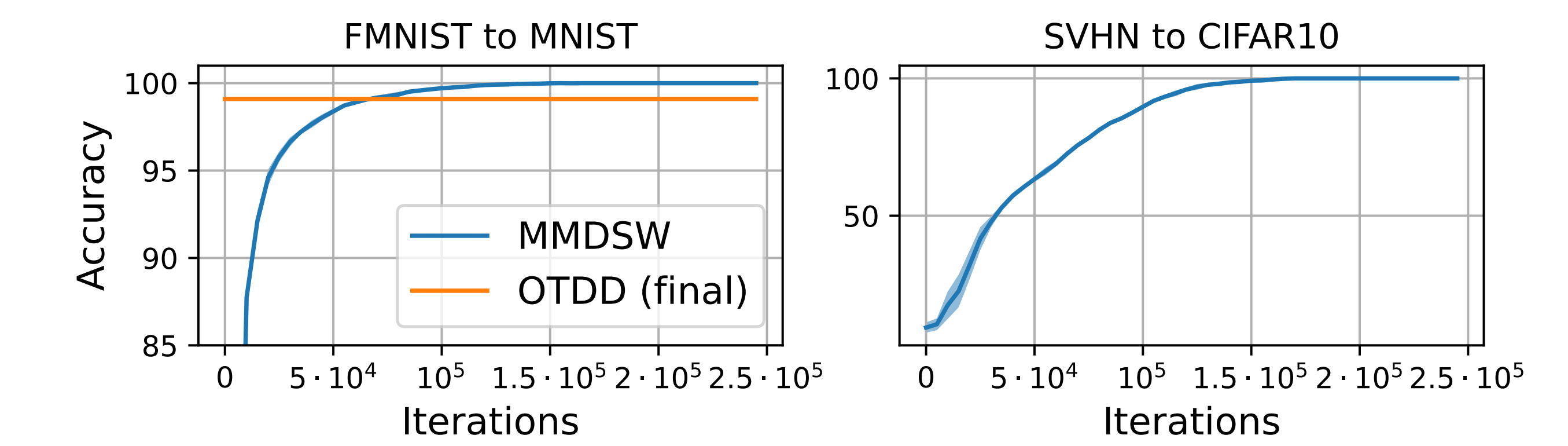
Tangent space: $T_{\mathbb{P}} \mathcal{P}_2(\mathcal{P}_2(\mathcal{M})) = \overline{\{\nabla_{W_2} \varphi, \varphi \in \text{Cyl}(\mathcal{P}_2(\mathcal{M}))\}}$

Properties: There is at most one element in $\partial \mathbb{F}(\mathbb{P}) \cap T_{\mathbb{P}} \mathcal{P}_2(\mathcal{P}_2(\mathcal{M}))$. If $\xi \in \partial \mathbb{F}(\mathbb{P}) \cap T_{\mathbb{P}} \mathcal{P}_2(\mathcal{P}_2(\mathcal{M}))$, then ξ is a strong differential of \mathbb{F} at \mathbb{P} (i.e. the Taylor expansion holds for any $\tilde{\Gamma} \in \mathcal{P}_2(\mathcal{P}_2(T\mathcal{M}))$ s.t. $\phi_{\#}^{\mathcal{M}} \tilde{\Gamma} = \mathbb{P}$).

Domain Adaptation

Minimize $\mathbb{F}(\mathbb{P}) = \frac{1}{2} \text{MMD}_K(\mathbb{P}, \mathbb{Q})$ starting from \mathbb{P}_0 (FMNIST or SVHN) towards \mathbb{Q} (MNIST or CIFAR10).

- Pretrain a classifier on \mathbb{Q}
- Monitor accuracy of the classifier along the flow



Applications

Dataset distillation. Synthesize a dataset $\mathbb{Q} = \frac{1}{n} \sum_{c=1}^C \delta_{\nu^{c,n}}$ (n big) with a dataset $\mathbb{P} = \frac{1}{C} \sum_{c=1}^C \delta_{\mu^{c,k}}$ (k small).

→ $\min_{\mathbb{P}} \mathbb{E}_{\theta, \omega} [\text{MMD}_K^2(\phi_{\#}^{\theta, \omega} \mathbb{P}, \phi_{\#}^{\theta, \omega} \mathbb{Q})]$ with $\phi^{\theta, \omega}(\mu) = \psi_{\#}^{\theta} \mathcal{A}_{\#}^{\omega} \mu$, \mathcal{A} a random augmentation, ψ^{θ} a randomly initialized neural network.

Evaluation: train a classifier on the new synthetic dataset \mathbb{P}

Transfer learning (k-shot learning). Augment a dataset $\mathbb{Q} = \frac{1}{C} \sum_{c=1}^C \delta_{\nu^{c,k}}$ (k small) adding flowed samples starting from $\mathbb{P}_0 = \frac{1}{C} \sum_{c=1}^C \delta_{\mu^{c,n}}$ a known bigger dataset.

→ $\min_{\mathbb{P}} \text{MMD}_K^2(\mathbb{P}, \mathbb{Q})$ starting from $\mathbb{P} = \mathbb{P}_0$

Evaluation: train a classifier on the augmented dataset $\hat{\mathbb{Q}}$

Dataset distillation						Transfer learning					
Dataset	k	$\psi^{\theta} = \mathcal{A}^{\omega} = \text{Id}$	DM	MMDSW	Baselines	Dataset	k	Train on Q	MMDSW	OTDD	(Hua et al., 2023)
MNIST	1	61.1±6.5	66.5±5.5	55.8±2.0	99.4	M to F	1	26.0±5.3	40.5±4.7	30.5±4.2	36.4±3.3
	5	88.2±2.8	93.2±0.7	92.2±1.1			5	38.5±6.7	61.5±4.6	59.7±1.8	62.7±1.1
	10	95.9±0.9	97.0±0.2	97.6±0.2			10	53.9±7.9	65.4±1.5	64.0±1.4	66.2±1.0
	50	95.9±0.9	97.0±0.2	97.6±0.2			100	71.1±1.5	74.7±0.8	-	73.5±0.7
FMNIST	1	54.4±3.2	60.0±4.1	49.0±7.5	92.4	M to K	1	18.4±3.1	20.9±2.0	18.8±2.1	19.4±1.9
	5	74.6±1.0	76.7±1.0	75.3±0.7			5	25.9±4.0	37.4±2.2	31.3±1.4	39.0±1.0
	10	81.3±0.5	84.2±0.1	83.2±0.2			10	30.9±4.6	44.7±1.8	34.1±0.9	44.1±1.2
	50	81.3±0.5	84.2±0.1	83.2±0.2			100	60.1±1.1	66.8±0.8	66.3±0.9	62.4±1.2