# Mirror and Preconditioned Gradient Descent in Wasserstein Space

Clément Bonet[1], Théo Uscidda[1], Adam David[2], Pierre-Cyril Aubin-Frankowski[3], Anna Korba[1]

[1]ENSAE, CREST, Institut Polytechnique de Paris; [2]TU Berlin; [3]TU Wien

## Contributions

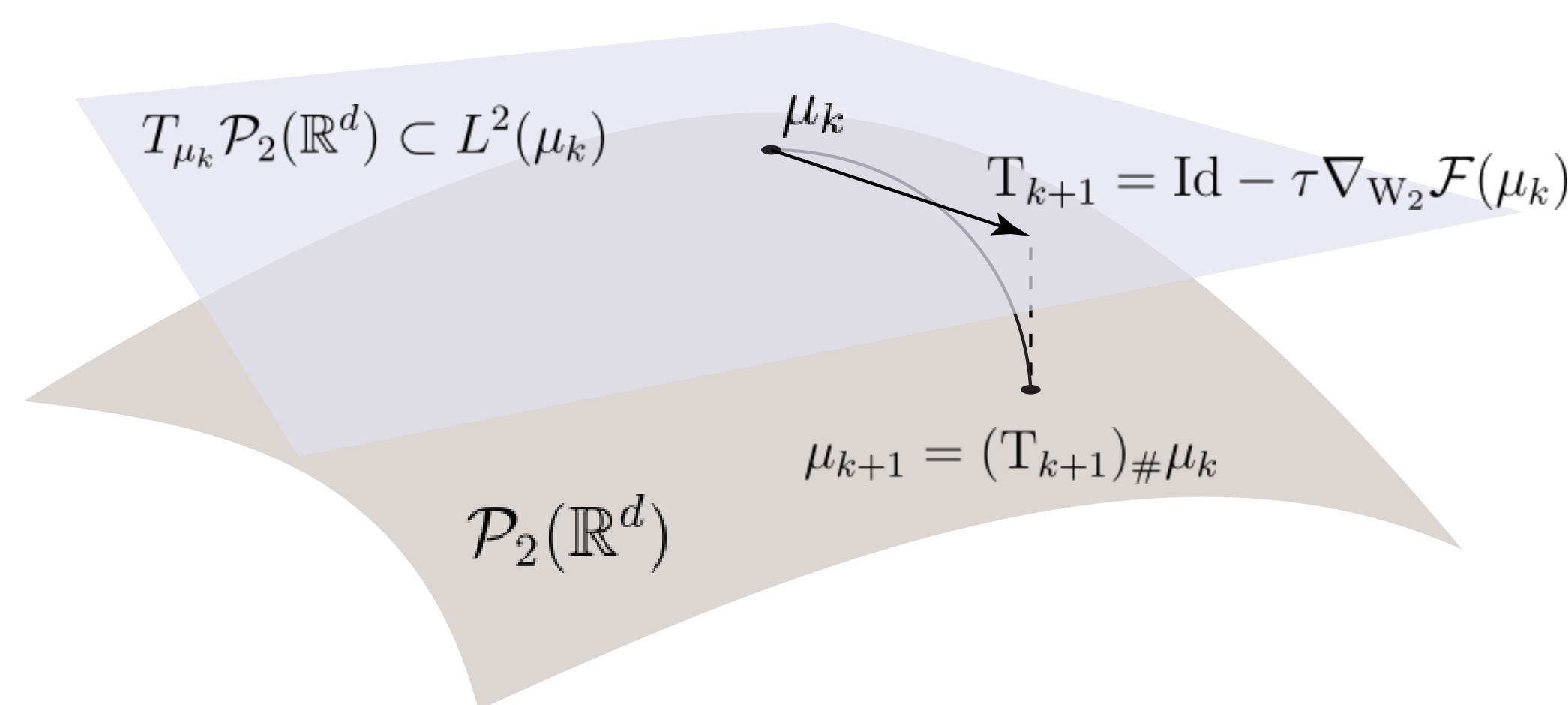**Goal:** $\min_{\mu \in \mathcal{P}_2(\mathbb{R}^d)} \mathcal{F}(\mu)$ for $\mathcal{F} : \mathcal{P}_2(\mathbb{R}^d) \to \mathbb{R}$

- Study two optimization schemes of the form

$$\begin{cases} T_{k+1} = \operatorname{argmin}_{T \in L^2(\mu_k)} \ d(T, \mathrm{Id}) + \langle \nabla_{W_2} \mathcal{F}(\mu_k), T - \mathrm{Id} \rangle_{L^2(\mu_k)} \\ \mu_{k+1} = (T_{k+1})_{\#} \mu_k \end{cases}$$

- Provide descent and convergence conditions
- Verification of the benefit on experiments



## Wasserstein Space

**Wasserstein gradient:** For $\mathcal{F} : \mathcal{P}_2(\mathbb{R}^d) \to \mathbb{R}$, $\gamma \in \Pi_o(\mu, \nu)$,

$$\mathcal{F}(\nu) = \mathcal{F}(\mu) + \int \langle \nabla_{W_2} \mathcal{F}(\mu)(x), y - x \rangle \, d\gamma(x,y) + o(W_2(\mu,\nu))$$

For $\mathcal{F} : \mathcal{P}_2(\mathbb{R}^d) \to \mathbb{R}$, define $\tilde{\mathcal{F}}_\mu(T) := \mathcal{F}(T_{\#}\mu)$.
If $\mathcal{F}$ $W_2$-differentiable, $\nabla \tilde{\mathcal{F}}_\mu(T) = \nabla_{W_2} \mathcal{F}(T_{\#}\mu) \circ T$.

**Examples:** potentials $\mathcal{V}_V(\mu) = \int V d\mu$, interactions $\mathcal{W}_W(\mu) = \iint W(x-y) d\mu(x) d\mu(y)$, entropy $\mathcal{H}(\mu) = \int \log(\mu(x)) d\mu(x)$. $\nabla_{W_2} \mathcal{V}_V(\mu) = \nabla V$, $\nabla_{W_2} \mathcal{W}_W(\mu) = \nabla W \star \mu$, $\nabla_{W_2} \mathcal{H}(\mu) = \nabla \log \mu$

## Bregman Divergence and Convexity

**Bregman divergence:** Let $\phi_\mu : L^2(\mu) \to \mathbb{R}$, $T, S \in L^2(\mu)$,

$$d_{\phi_\mu}(T,S) = \phi_\mu(T) - \phi_\mu(S) - \langle \nabla \phi_\mu(S), T - S \rangle_{L^2(\mu)}$$

**Relative smoothness/convexity** along $t \mapsto \mu_t$ with $\mu_t = (T_t)_{\#}\mu$, $T_t = (1-t)S + tT$ for $S, T \in L^2(\mu)$. $\mathcal{F}$ is $\beta$-smooth (resp. $\alpha$-convex) relative to $\mathcal{G}$ along $t \mapsto \mu_t$ if for all $s, t \in [0,1]$, $d_{\tilde{\mathcal{F}}_\mu}(T_s, T_t) \leq \beta d_{\tilde{\mathcal{G}}_\mu}(T_s, T_t)$ (resp. $d_{\tilde{\mathcal{F}}_\mu}(T_s, T_t) \geq \alpha d_{\tilde{\mathcal{G}}_\mu}(T_s, T_t)$).

- For $\mathcal{F} = \mathcal{V}_V$, $\mathcal{G} = \mathcal{V}_U$: holds provided $V$ $\beta$-smooth (resp. $\alpha$-convex) relative to $U$
- For $\mathcal{F} = \mathcal{W}_W$, $\mathcal{G} = \mathcal{W}_K$: holds provided $W$ $\beta$-smooth (resp. $\alpha$-convex) relative to $K$
- $\mathcal{F} = \mathcal{V}_V + \mathcal{H}$ 1-convex relative to $\mathcal{V}_V$ and $\mathcal{H}$

## Implementation of the Schemes

**Mirror descent:** $d = \frac{1}{\tau} d_{\phi_\mu}$, by FOC: $\nabla \phi_\mu(T_{k+1}) = \nabla \phi_\mu(\mathrm{Id}) - \tau \nabla_{W_2} \mathcal{F}(\mu_k)$
For $\phi_\mu(T) = \int V \circ T \, d\mu = \mathcal{V}_V(T_{\#}\mu)$, $T_{k+1} = \nabla V^* \circ (\nabla V - \tau \nabla_{W_2} \mathcal{F}(\mu_k))$
In general: Newton method

**Preconditioned gradient descent:**

$$d(T, S) = \phi_\mu^h((S-T)/\tau)\tau = \int h((S(x) - T(x))/\tau)\tau \, d\mu(x)$$

FOC: $T_{k+1} = \mathrm{Id} - \tau \nabla h^* \circ \nabla_{W_2} \mathcal{F}(\mu_k)$

For $\mu_k = \frac{1}{n} \sum_{i=1}^n \delta_{x_i^k}$, for all $k \geq 0$, $i \in \{1, \ldots, n\}$, $x_i^{k+1} = T_{k+1}(x_i^k)$.

## Theory of Mirror Descent in Wasserstein Space

Let $\beta > 0$, $\tau \leq \frac{1}{\beta}$. For any $\mu \in \mathcal{P}_2(\mathbb{R}^d)$, let $\phi_\mu : L^2(\mu) \to \mathbb{R}$ be strictly convex, proper and differentiable. Assume $\phi_\mu(T) = \phi(T_{\#}\mu)$ for $\phi : \mathcal{P}_2(\mathbb{R}^d) \to \mathbb{R}$.
Define $W_\phi(\mu, \nu) = \inf_{\gamma \in \Pi(\mu,\nu)} \phi(\nu) - \phi(\mu) - \int \langle \nabla_{W_2} \phi(\mu)(y), x - y \rangle \, d\gamma(x,y)$.

**Assumptions:** Let $T_{\phi_{\mu_k}}^{\mu_k, \mu^*} = \operatorname{argmin}_{T, T_{\#}\mu_k = \mu^*} d_{\phi_{\mu_k}}(T, \mathrm{Id})$. For all $k \geq 0$,

1. $\mathcal{F}$ is $\beta$-smooth relative to $\phi$ along $t \mapsto ((1-t)\mathrm{Id} + tT_{k+1})_{\#}\mu_k$
2. $\mathcal{F}$ is $\alpha$-convex relative to $\phi$ along the curves $t \mapsto ((1-t)\mathrm{Id} + tT_{\phi_{\mu_k}}^{\mu_k, \mu^*})_{\#}\mu_k$
3. $d_{\phi_{\mu_k}}(T_{\phi_{\mu_k}}^{\mu_k, \mu^*}, \mathrm{Id}) = W_\phi(\mu^*, \mu_k)$ and $d_{\phi_{\mu_k}}(T_{\phi_{\mu_k}}^{\mu_k, \mu^*}, T_{k+1}) \geq W_\phi(\mu^*, \mu_{k+1})$ (True $e.g.$ if $\mu_k, \mu_{k+1} \in \mathcal{P}_{2,ac}(\mathbb{R}^d)$ and $\nabla_{W_2} \phi(\mu_k), \nabla_{W_2} \phi(\mu_{k+1})$ invertibles)

### Convergence Results

- Under 1), for all $k \geq 0$ $\mathcal{F}(\mu_{k+1}) \leq \mathcal{F}(\mu_k) - \frac{1}{\tau} d_{\phi_{\mu_k}}(\mathrm{Id}, T_{k+1})$
- Under 1), 2), 3), for all $k \geq 1$, $\mathcal{F}(\mu_k) - \mathcal{F}(\mu^*) \leq \frac{1-\alpha\tau}{k\tau} W_\phi(\mu^*, \mu_0)$

## Theory of Preconditioned GD in Wasserstein Space

Let $\beta > 0$, $\tau \leq \frac{1}{\beta}$ and $\bar{T} = \operatorname{argmin}_{T, T_{\#}\mu_k = \mu^*} d_{\tilde{\mathcal{F}}_{\mu_k}}(\mathrm{Id}, \bar{T})$.
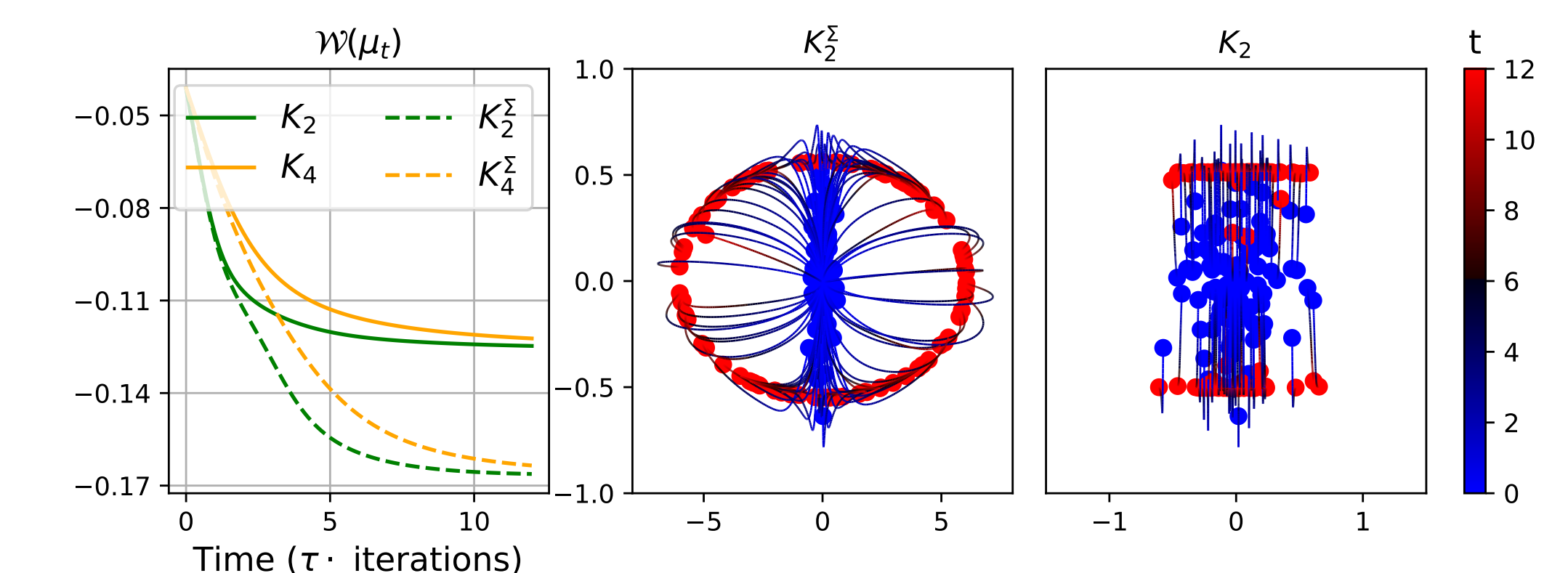
**Assumptions:** For all $k \geq 0$,

1. $\mathcal{F}$ convex along $t \mapsto ((1-t)\mathrm{Id} + tT_{k+1})_{\#}\mu_k$
2. $d_{\phi_{\mu_k}^{h^*}}(\nabla_{W_2} \mathcal{F}(\mu_{k+1}) \circ T_{k+1}, \nabla_{W_2} \mathcal{F}(\mu_k)) \leq \beta d_{\tilde{\mathcal{F}}_{\mu_k}}(\mathrm{Id}, T_{k+1})$
3. $\alpha d_{\tilde{\mathcal{F}}_{\mu_k}}(\mathrm{Id}, \bar{T}) \leq d_{\phi_{\mu_k}^{h^*}}(\nabla_{W_2} \mathcal{F}(\bar{T}_{\#}\mu_k) \circ \bar{T}, \nabla_{W_2} \mathcal{F}(\mu_k))$

### Convergence Results

- Under 1), 2), $\phi_{\mu_{k+1}}^{h^*}(\nabla_{W_2} \mathcal{F}(\mu_{k+1})) \leq \phi_{\mu_k}^{h^*}(\nabla_{W_2} \mathcal{F}(\mu_k)) - \frac{1}{\tau} d_{\tilde{\mathcal{F}}_{\mu_k}}(T_{k+1}, \mathrm{Id})$
- Under 1), 2), 3), $\phi_{\mu_k}^{h^*}(\nabla_{W_2} \mathcal{F}(\mu_k)) - h^*(0) \leq \frac{1-\tau\alpha}{\tau k}(\mathcal{F}(\mu_0) - \mathcal{F}(\mu^*))$

## Mirror Descent Experiments

**Minimization** of an interaction energy $\mathcal{F}(\mu) = \mathcal{W}_W(\mu)$ with $W(z) = \frac{1}{4}\|z\|_{\Sigma^{-1}}^4 - \frac{1}{2}\|z\|_{\Sigma^{-1}}^2$ and $\phi(\mu) = \mathcal{W}_K(\mu)$ with $K_2^\Sigma(z) = \frac{1}{2}\|z\|_{\Sigma^{-1}}^2$, $K_2 = K_2^{I_2}$, $K_4^\Sigma(z) = \frac{1}{4}\|z\|_{\Sigma^{-1}}^4 + \frac{1}{2}\|z\|_{\Sigma^{-1}}^2$, $K_4 = K_4^{I_2}$.
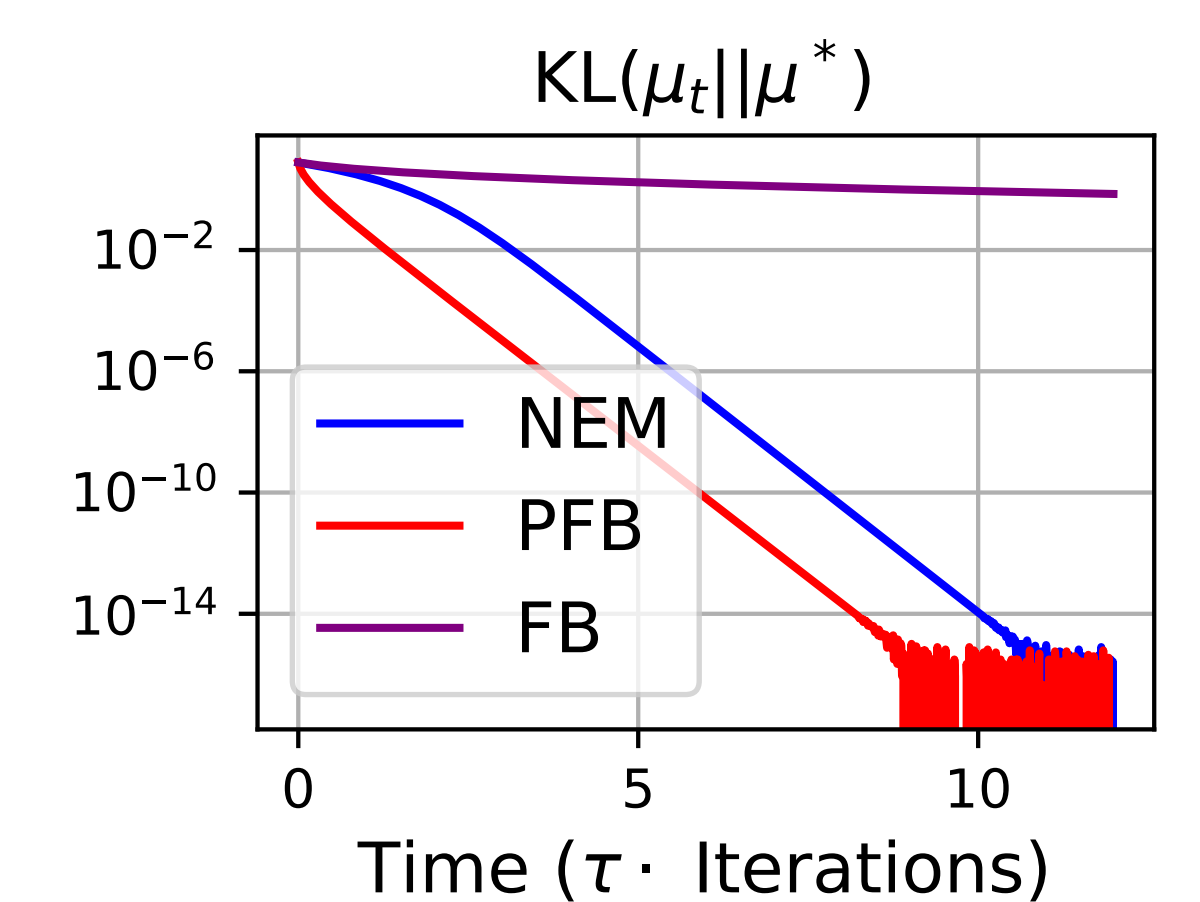


**Minimization** of

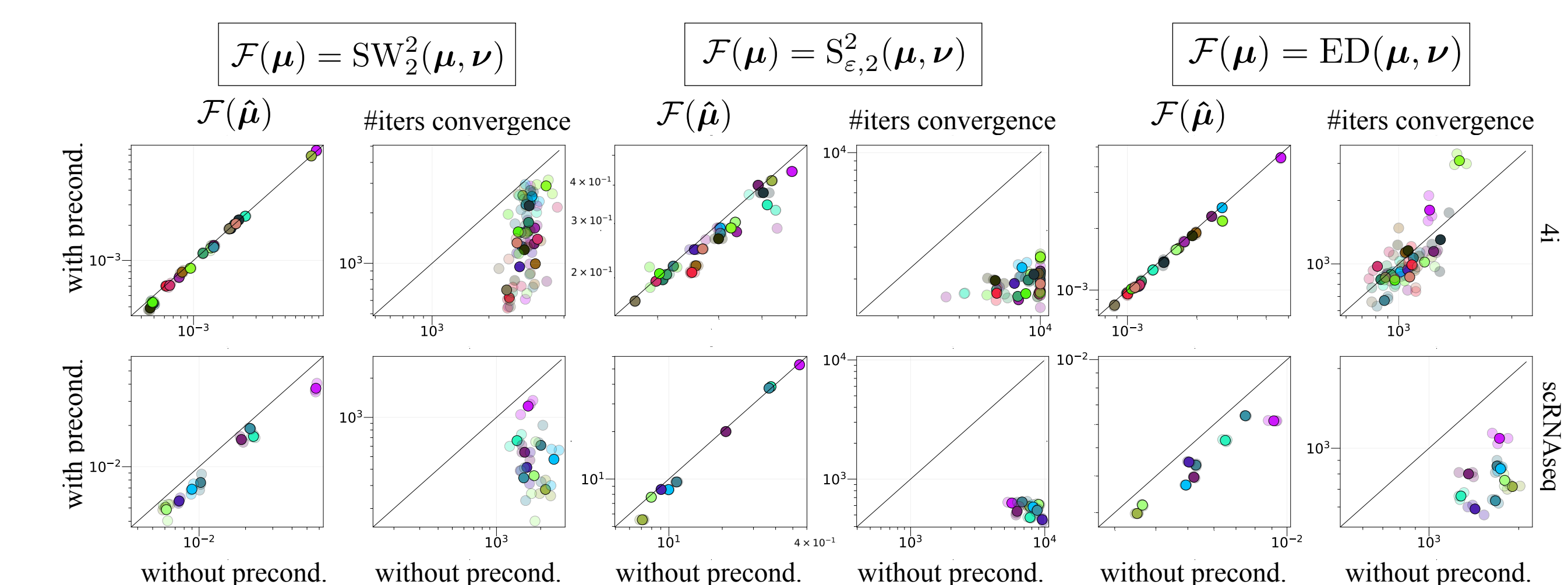$$\mathcal{F}(\mu) = \mathcal{V}_V(\mu) + \mathcal{H}(\mu),$$

for $V(x) = \frac{1}{2} x^T \Sigma^{-1} x$ with

$\phi(\mu) = \int \frac{1}{2}\|x\|_2^2 \, d\mu(x)$ (FB),
$\phi(\mu) = \mathcal{V}_V(\mu)$ (PFB),
$\phi(\mu) = \mathcal{H}(\mu)$ (NEM).



## Preconditioned GD for Single Cells

**Minimize** $\mathcal{F}(\mu) = D(\mu, \nu)$ with $\mu_0$ untreated cells and $\nu$ perturbed cells. Use $h^*(x) = (\|x\|_2^a + 1)^{1/a} - 1$ with $a \in \{1.25, 1.5, 1.75\}$ which is well suited to minimize functions growing in $\|x - x^*\|^{a/(a-1)}$.



- Rows: 2 profiling technologies
- Points: For treatment $i$, $z_i = (x_i, y_i)$ with $x_i$ value of $\mathcal{F}(\hat{\mu}) = D(\hat{\mu}, \nu)$ (1st subcolumn) or number of iterations to converge (2nd subcolumn) without preconditioning and $y_i$ with preconditioning

→ Points below the diagonal: **Preconditioned GD provides a better minimum or converges faster**