# Dataset Dynamics via Gradient Flows in Probability Space

David Alvarez-Melis
Nicolò Fusi

## Gradient Flows…

- GF: curve $x(t)$ of steepest descent on functional $F : \mathcal{X} \to \mathbb{R}$
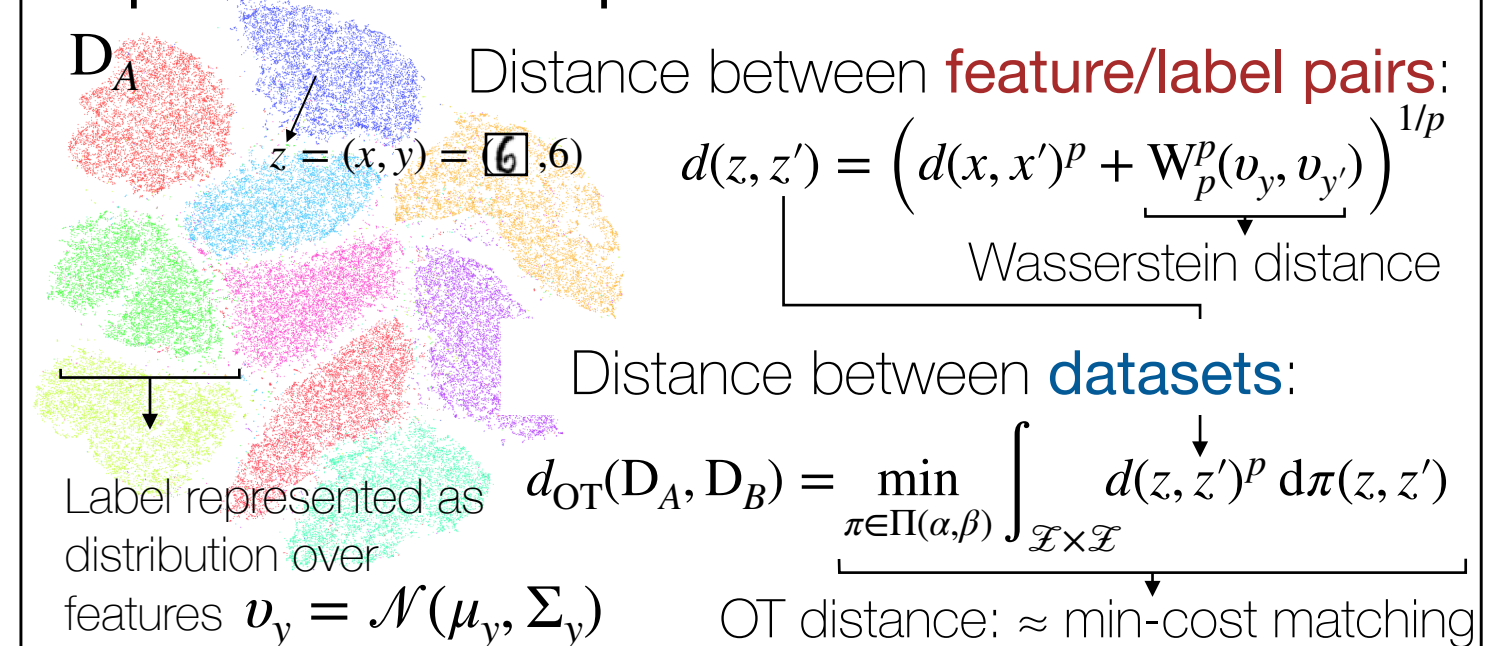
- In Euclidean space $\mathcal{X}$:
$$x'(t) = -\nabla F(x(t))$$

- In Probability space $\mathcal{P}(\mathcal{X})$:
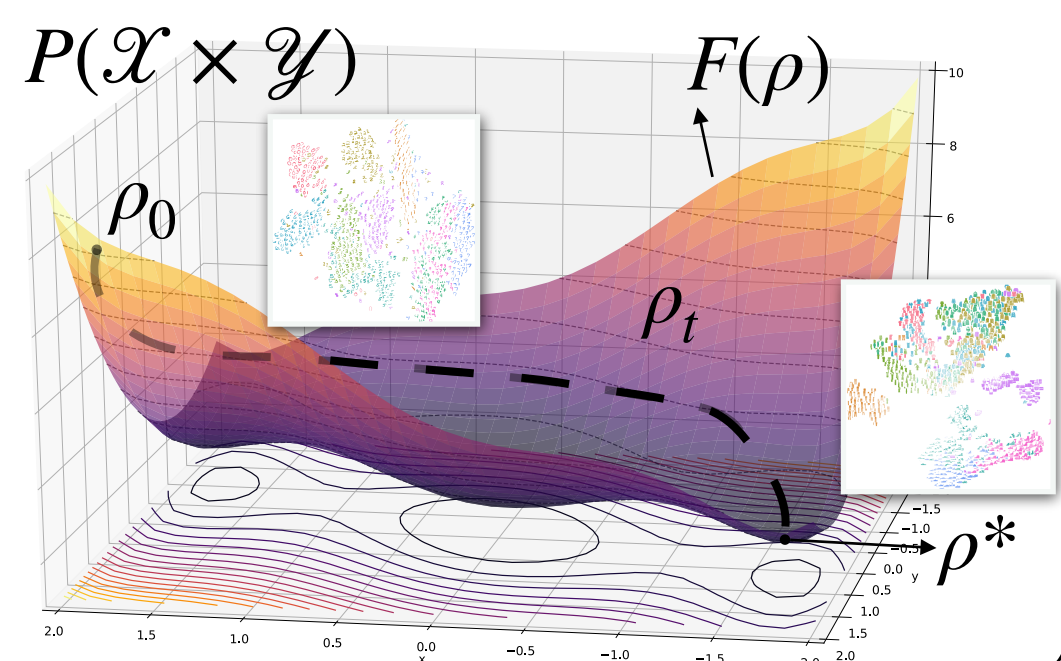$$\partial_t \rho_t = \nabla \cdot \left( \rho_t \nabla \frac{\delta F}{\delta \rho} \right)$$

## Motivation & Summary

- Dataset **transformation**: ubiquitous in ML, from augmentation to generation

- Need often arises because available (generic) data $\neq$ needed (task-specific) data

- Here: general, principled, efficient **labeled** dataset **transformation by optimization**

- Solved using **gradient flows**: guarantees, efficient computation, yields full path

- Vision: **data-centric learning** paradigm, complementary to model-centric one

## Optimal Transport Dataset Distance

$\mathbf{D}_A$

$z = (x, y) = (\boxed{6}, 6)$

Distance between **feature/label pairs**:
$$d(z, z') = \left( d(x, x')^p + W_p^p(v_y, v_{y'}) \right)^{1/p}$$
Wasserstein distance

Label represented as distribution over features $v_y = \mathcal{N}(\mu_y, \Sigma_y)$

Distance between **datasets**:
$$d_{OT}(\mathbf{D}_A, \mathbf{D}_B) = \min_{\pi \in \Pi(\alpha, \beta)} \int_{\mathcal{Z} \times \mathcal{Z}} d(z, z')^p \, d\pi(z, z')$$
OT distance: $\approx$ min-cost matching

## … in dataset space



$P(\mathcal{X} \times \mathcal{Y})$    $F(\rho)$

$\rho_0$    $\rho_t$    $\rho^*$

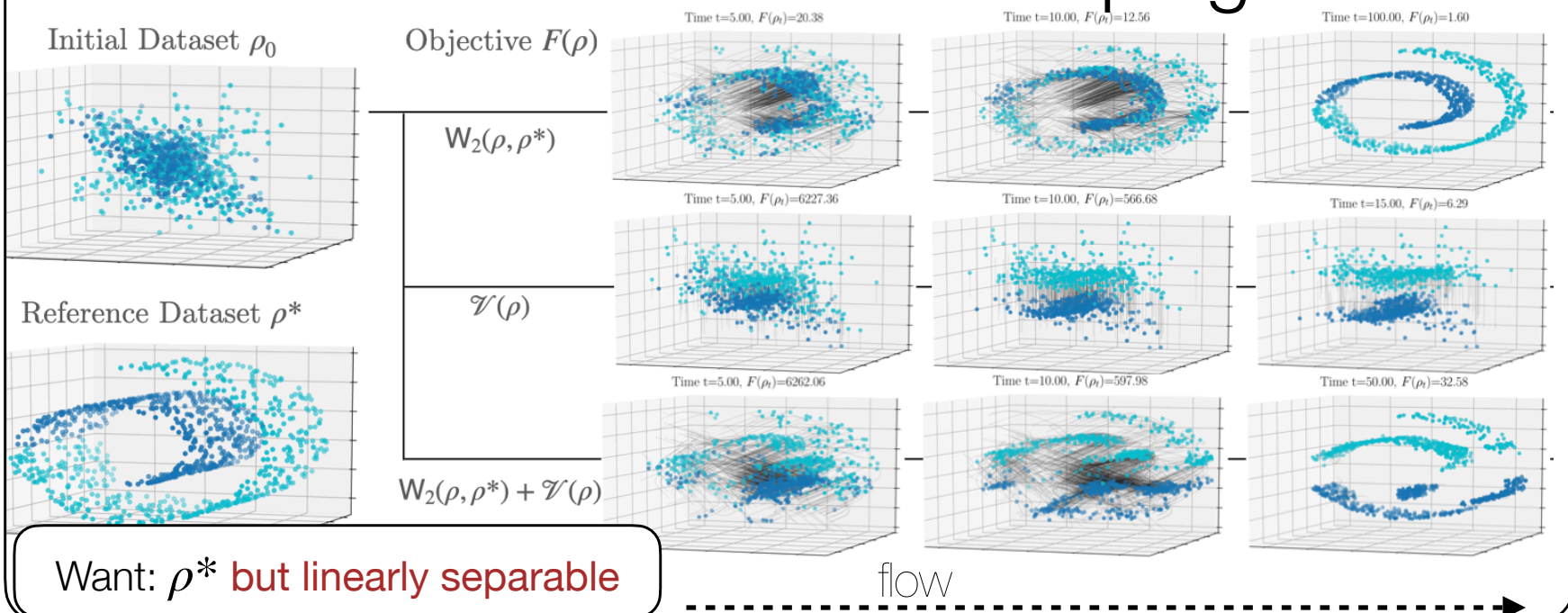## Dataset Transformation using Gradient Flows

- Model $F(z)$ via well-behaved functionals:

$$\mathcal{F}(\rho) = \int f(\rho(z))dz \quad \text{internal energy}$$

$$\mathcal{V}(\rho) = \int V(z)d\rho \quad \text{potential}$$

$$\mathcal{W}(\rho) = \frac{1}{2} \iint W(z - z')d\rho(z)d\rho(z') \quad \text{interaction}$$

$$\mathcal{T}(\rho) = \text{OTDD}(\rho, \beta) \quad \text{distance}$$

Objective: $\min_{\rho \in \mathcal{P}(\mathcal{X} \times \mathcal{Y})} F(\rho)$

Flow: $\partial_t \rho_t(z) = \nabla \cdot \left( \rho_t(z) \nabla \frac{\delta F}{\delta \rho}(z) \right)$

have pop. convergence guarantees, simple 'derivative', tractable

can model various useful objectives on datasets

## Practical Implementation

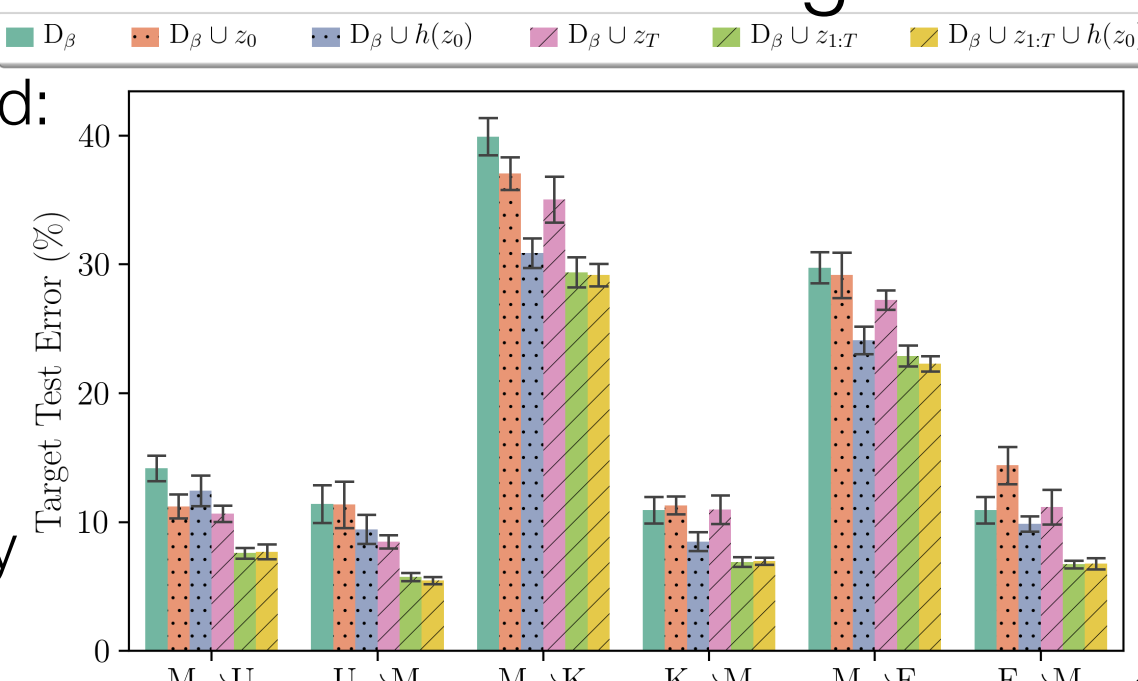- Flow discretized in time (Euler) & space (particles):
$$z_{t+1}^{(i)} = z_t^{(i)} - \gamma \nabla_z F(z_t^{(i)}), \quad i \in \{i, \dots, n\}$$

- Distribution is approximated as $\hat{\rho}_t = \sum p_i \delta_{z_t^{(i)}}$

- Implemented using automatic differentiation

- Challenge for $F = \text{OTDD}$: how to update labels, we propose three types of update schemes

- Unlabeled data? Semi- and un-supervised flows

## Flows for Dataset Shaping



Initial Dataset $\rho_0$

Objective $F(\rho)$

$W_2(\rho, \rho^*)$

$\mathcal{V}(\rho)$

$W_2(\rho, \rho^*) + \mathcal{V}(\rho)$

Reference Dataset $\rho^*$

flow

Want: $\rho^*$ but linearly separable

## Flows for Transfer Learning

- Semi-supervised:
1. run flow with labeled data
2. fit parametrized model of flow

- Flowed data helps, especially full trajectories!



Legend: $D_\beta$   $D_\beta \cup z_0$   $D_\beta \cup h(z_0)$   $D_\beta \cup z_T$   $D_\beta \cup z_{1:T}$   $D_\beta \cup z_{1:T} \cup h(z_0)$

Target Test Error (%)

M→U   U→M   M→K   K→M   M→F   F→M

## Flows for Model Re-Purposing

- ResNet trained on CIFAR10, **frozen**

- Target dataset: Camelyon10

- Flow: CAM→CIFAR

- High acc. on flowed data!



Accuracy

— Acc. on Flowed Data
-- Flow Objective
···· Baseline Accuracy