

# Unbalanced Gromov-Wasserstein distance

---

Thibault Séjourné

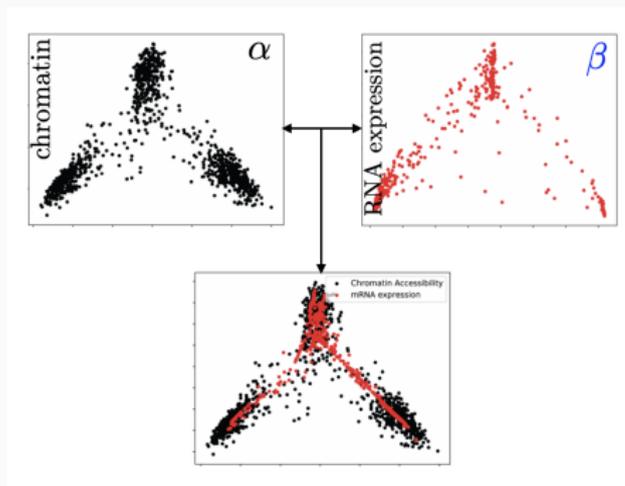
SMAI Minisymposium – 24th June, 2021

Joint work with Francois-Xavier Vialard, Gabriel Peyré, Jean Feydy and Alain Trounev

# Introduction

---

# Real world motivation: aligning genomic data<sup>1</sup>



Aligning RNA data

Data is:

- Heterogeneous ( $\neq$  dimensions)
- imbalanced

$\Rightarrow$  Need for adaptive and robust (yet meaningful) assignments !

<sup>1</sup>Demetci, Pinar, et al. Gromov-Wasserstein optimal transport to align single-cell multi-omics data.

# Optimal transport and its generalizations

Optimal transport displays three restrictions:

- Compares measures with same mass,
- Compares measures defined on the same space,
- Scales poorly in numerical solvers :  $O(n^3 \log(n))$ .

There exists extensions to overcome these issues:

- Unbalanced optimal transport,
- Gromov-Wasserstein distances,
- Entropic regularization.

# Outline of the presentation

1. Background - UOT ( ● )
2. Sinkhorn algorithm ( ● + ● )
3. Unbalanced Gromov-Wasserstein ( ● + ● )
4. Implementation of UGW ( ● + ● + ● )

## Unbalanced OT

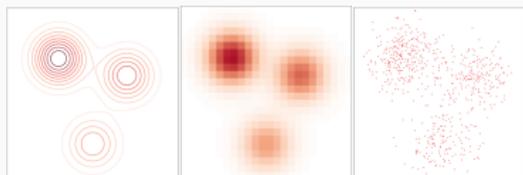
---

# Representing probabilities and measures

Several models for measures, most commonly pointclouds.

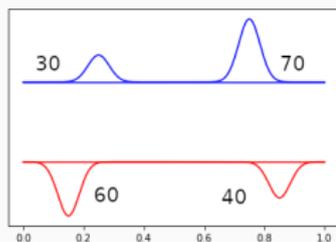
$$\text{Measure } \alpha = \sum_i \alpha_i \delta_{x_i},$$

$$\text{Mass } m(\alpha) = \sum_i \alpha_i.$$



Important example: Gaussian densities in  $\mathbb{R}$  with  $\alpha \propto \sum_i p(x_i) \delta_{x_i}$

- $\beta$  = mixture 70/30
- $\alpha$  = mixture 40/60



# Optimal Transport (OT)

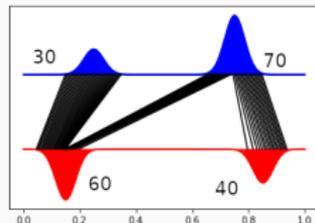
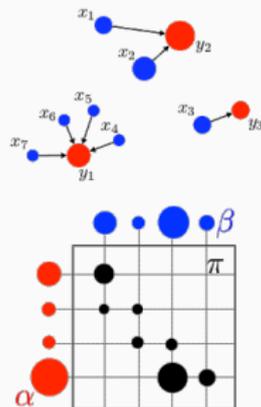
## Balanced Optimal Transport Distance<sup>2</sup>

$$\text{OT}(\alpha, \beta) \stackrel{\text{def.}}{=} \min_{\pi \geq 0} \left\{ \sum_{i,j} C_{ij} \pi_{ij} : \begin{array}{l} \pi \mathbf{1} = \alpha \\ \pi^\top \mathbf{1} = \beta \end{array} \right\}.$$

Called p-Wasserstein distance for  $C = d^p$ .

**Intuition:** Moving  $\pi_{ij}$  grams from  $x_i$  to  $y_j$  costs  $\pi_{ij} \times C_{ij}$ .

**Choice of C**  $\rightarrow$  Choice of geometric prior.



<sup>2</sup>Kantorovich, L. (1942). On the transfer of masses (in Russian).

# Unbalanced OT

**Idea:** Soften the constraint  $\pi \mathbf{1} = \alpha$

$$\rightarrow \text{KL}(\pi \mathbf{1} | \alpha) = \sum_i \log\left(\frac{\pi_{1,i}}{\alpha_i}\right) \alpha_i - m(\pi \mathbf{1}) + m(\alpha)$$

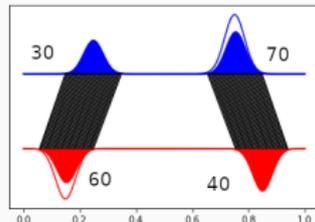
## Definition - Unbalanced OT<sup>3</sup>

For any **positive** measures  $(\alpha, \beta)$  one defines

$\text{UOT}_\rho(\alpha, \beta) = \inf_{\pi \geq 0} \mathcal{L}_{\text{UOT}}(\pi)$  where

$$\mathcal{L}_{\text{UOT}}(\pi) \stackrel{\text{def.}}{=} \sum_{i,j} C_{ij} \pi_{ij} + \rho \text{KL}(\pi \mathbf{1} | \alpha) + \rho \text{KL}(\pi^\top \mathbf{1} | \beta).$$

- **2 choices:** transport vs create/destroy
- **Other penalties:** TV, or Csiszàr div  $D_\varphi$ .
- **Balanced OT** =  $\rho \rightarrow \infty$  or  $D_\varphi = \iota(=)$ .



<sup>3</sup>Liero, M., Mielke, A., & Savaré, G. (2018). Optimal entropy-transport problems and a new Hellinger–Kantorovich distance between positive measures.

# Entropic Optimal Transport

---

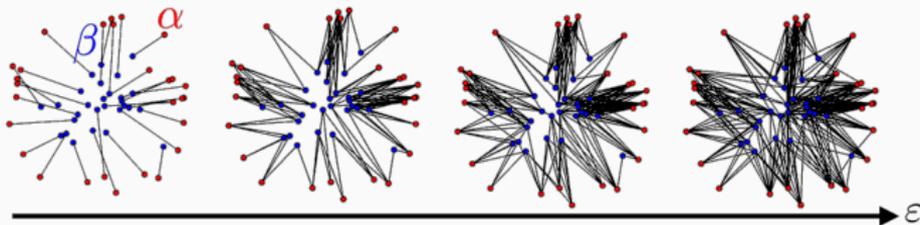
# Regularization of OT

**Reminder:** OT is computationally expensive.

**Idea:** Add an entropic penalty  $\varepsilon \text{KL}(\pi | \alpha \otimes \beta)$ .

## Entropic Unbalanced OT<sup>4 5</sup>

$$\text{UOT}_{\varepsilon, \rho}(\alpha, \beta) \stackrel{\text{def.}}{=} \inf_{\pi \geq 0} \sum_{i,j} C_{ij} \pi_{ij} + \rho \text{KL}(\pi \mathbf{1} | \alpha) + \rho \text{KL}(\pi^\top \mathbf{1}, \beta) + \varepsilon \text{KL}(\pi | \alpha \otimes \beta)$$



<sup>4</sup> Cuturi, M. (2013). Sinkhorn distances: Lightspeed computation of optimal transport.

<sup>5</sup> Chizat, L., Peyré, G., Schmitzer, B., & Vialard, F. X. (2018). Scaling algorithms for unbalanced optimal transport problems.

# Duality of regularized Balanced OT

The dual for  $\text{UOT}_{\varepsilon,\rho}$  reads

$$\begin{aligned} \text{UOT}_{\varepsilon,\rho}(\alpha, \beta) = \sup_{f,g} & \sum_i \rho(1 - e^{-f_i/\rho})\alpha_i + \sum_j \rho(1 - e^{-g_j/\rho})\beta_j \\ & - \varepsilon \sum_{i,j} \left( e^{\frac{f_i+g_j-C_{ij}}{\varepsilon}} - 1 \right) \alpha_i \beta_j. \end{aligned}$$

We consider **alternate dual ascent** to compute  $\text{UOT}_{\varepsilon,\rho}$ :

## Alternate dual ascent

Given any initialization  $f_0$ . At time  $t$  one has  $(f_t, g_t)$ . Then iterate until convergence:

1. Fix  $f_t$  and find optimal  $g$  in the dual  $\rightarrow g_{t+1}$ ,
2. Fix  $g_{t+1}$  and find optimal  $f$  in the dual  $\rightarrow f_{t+1}$ .

# Duality of regularized Balanced OT

The dual for  $\text{UOT}_{\varepsilon, \rho}$  reads

$$\begin{aligned} \text{UOT}_{\varepsilon, \rho}(\alpha, \beta) = \sup_{f, g} & \sum_i \rho(1 - e^{-f_i/\rho}) \alpha_i + \sum_j \rho(1 - e^{-g_j/\rho}) \beta_j \\ & - \varepsilon \sum_{i, j} \left( e^{\frac{f_i + g_j - C_{ij}}{\varepsilon}} - 1 \right) \alpha_i \beta_j. \end{aligned}$$

**Unbalanced Sinkhorn algorithm = Alternate dual ascent**

$$\begin{aligned} f_i &\leftarrow \frac{\rho}{\varepsilon + \rho} \left[ -\varepsilon \log \sum_j e^{(g_j - C_{ij})/\varepsilon} \beta_j \right], \\ g_j &\leftarrow \frac{\rho}{\varepsilon + \rho} \left[ -\varepsilon \log \sum_i e^{(f_i - C_{ij})/\varepsilon} \alpha_i \right]. \end{aligned}$$

Rmk: Solve dual  $\Rightarrow$  Solve primal:  $\pi_{ij}^* = \exp((f_i^* + g_j^* - C_{ij})/\varepsilon) \alpha_i \beta_j$ .

# Unbalanced Gromov-Wasserstein

---

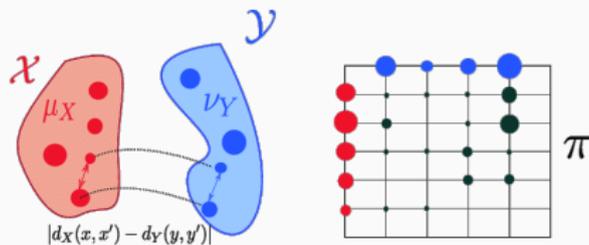
# Balanced Gromov-Wasserstein distance

**mm-space:**  $\mathcal{X} = (X, d^{(X)}, \alpha)$  with  $(X, d^{(X)})$  complete separable,  $\alpha$  positive measure

## Definition - GW distance<sup>6</sup>

Take  $\mathcal{X} = (X, d^{(X)}, \alpha)$  and  $\mathcal{Y} = (Y, d^{(Y)}, \beta)$  equipped with **probabilities**. One defines  $GW(\mathcal{X}, \mathcal{Y}) = \inf_{\{\pi \mathbf{1} = \alpha, \pi^\top \mathbf{1} = \beta\}} \mathcal{G}(\pi)$  where

$$\mathcal{G}(\pi) \stackrel{\text{def.}}{=} \sum_{i,j,k,l} \left( d_{ik}^{(X)} - d_{jl}^{(Y)} \right)^2 \pi_{ij} \pi_{kl}.$$



<sup>6</sup> Mémoli, F. (2011). Gromov-Wasserstein distances and the metric approach to object matching.

# Isometric mm-spaces

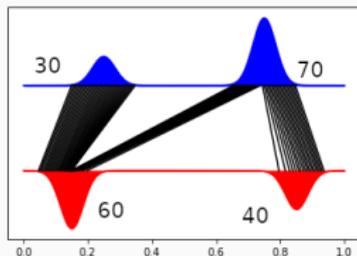
The GW distances encodes an equivalence relation of isometry.

## Isometric mm-spaces

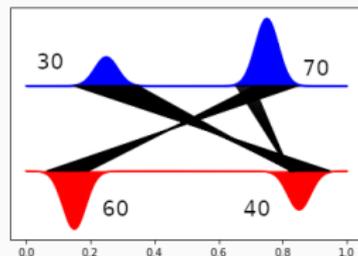
**Def:**  $\mathcal{X} \sim \mathcal{Y} \Leftrightarrow \exists \psi : X \rightarrow Y$  bijective isometry s.t.

$$d_X(x, x') = d_Y(\psi(x), \psi(x')) \quad \text{and} \quad \beta = \sum_i \alpha_i \delta_{\psi(x_i)}$$

**Prop:** With  $\lambda(t) = t^q$ ,  $GW^{\frac{1}{q}}$  distance and definite iff  $\mathcal{X} \sim \mathcal{Y}$



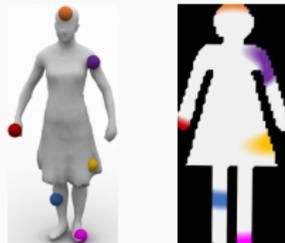
Balanced OT



GW

## Two key differences with OT

- GW is non-convex (quadratic assignment program)
- $(\mathcal{X}, \mathcal{Y})$  can differ radically in nature.<sup>7</sup>



---

<sup>7</sup>Solomon, J., Peyré, G., Kim, V. G., & Sra, S. (2016). Entropic metric alignment for correspondence problems.

# Unbalanced Gromov-Wasserstein

Define the tensor product of measures  $(\pi \otimes \pi)_{ijkl} \stackrel{\text{def.}}{=} \pi_{ij}\pi_{kl}$ .

## Definition

One defines  $UGW(\mathcal{X}, \mathcal{Y}) = \inf_{\pi \geq 0} \mathcal{L}_2(\pi)$  where

$$\mathcal{L}_{UGW}(\pi) = \sum_{i,j,k,l} \left( d_{ik}^{(X)} - d_{jl}^{(Y)} \right)^2 \pi_{ij}\pi_{kl} + \rho \text{KL}(\pi_1 \otimes \pi_1, \alpha \otimes \alpha) \\ + \rho \text{KL}(\pi_2 \otimes \pi_2, \beta \otimes \beta).$$

To be compared with

$$\mathcal{G}(\pi) = \sum_{i,j,k,l} \left( d_{ik}^{(X)} - d_{jl}^{(Y)} \right)^2 \pi_{ij}\pi_{kl},$$

$$\mathcal{L}_{UOT}(\pi) = \sum_{i,j} C_{ij}\pi_{ij} + \rho \text{KL}(\pi_1, \alpha) + \rho \text{KL}(\pi_2, \beta).$$

# Unbalanced Gromov-Wasserstein

Define the tensor product of measures  $(\pi \otimes \pi)_{ijkl} \stackrel{\text{def.}}{=} \pi_{ij}\pi_{kl}$ .

## Definition

One defines  $UGW(\mathcal{X}, \mathcal{Y}) = \inf_{\pi \geq 0} \mathcal{L}_2(\pi)$  where

$$\mathcal{L}_{UGW}(\pi) = \sum_{i,j,k,l} \left( d_{ik}^{(X)} - d_{jl}^{(Y)} \right)^2 \pi_{ij}\pi_{kl} + \rho \text{KL}(\pi_1 \otimes \pi_1 | \alpha \otimes \alpha) \\ + \rho \text{KL}(\pi_2 \otimes \pi_2 | \beta \otimes \beta).$$

To be compared with

$$\mathcal{G}(\pi) = \sum_{i,j,k,l} \left( d_{ik}^{(X)} - d_{jl}^{(Y)} \right)^2 \pi_{ij}\pi_{kl}, \\ \mathcal{L}_{UOT}(\pi) = \sum_{i,j} C_{ij}\pi_{ij} + \rho \text{KL}(\pi_1 | \alpha) + \rho \text{KL}(\pi_2 | \beta).$$

# Unbalanced Gromov-Wasserstein

Define the tensor product of measures  $(\pi \otimes \pi)_{ijkl} \stackrel{\text{def.}}{=} \pi_{ij}\pi_{kl}$ .

## Definition

One defines  $UGW(\mathcal{X}, \mathcal{Y}) = \inf_{\pi \geq 0} \mathcal{L}_2(\pi)$  where

$$\mathcal{L}_{UGW}(\pi) = \sum_{i,j,k,l} \left( d_{ik}^{(X)} - d_{jl}^{(Y)} \right)^2 \pi_{ij}\pi_{kl} + \rho \text{KL}(\pi_1 \otimes \pi_1 | \alpha \otimes \alpha) \\ + \rho \text{KL}(\pi_2 \otimes \pi_2 | \beta \otimes \beta).$$

To be compared with

$$\mathcal{G}(\pi) = \sum_{i,j,k,l} \left( d_{ik}^{(X)} - d_{jl}^{(Y)} \right)^2 \pi_{ij}\pi_{kl}, \\ \mathcal{L}_{UOT}(\pi) = \sum_{i,j} C_{ij}\pi_{ij} + \rho \text{KL}(\pi_1 | \alpha) + \rho \text{KL}(\pi_2 | \beta).$$

# Unbalanced Gromov-Wasserstein

Define the tensor product of measures  $(\pi \otimes \pi)_{ijkl} \stackrel{\text{def.}}{=} \pi_{ij}\pi_{kl}$ .

## Definition

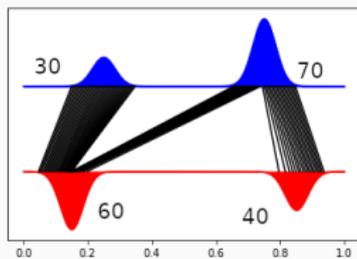
One defines  $UGW(\mathcal{X}, \mathcal{Y}) = \inf_{\pi \geq 0} \mathcal{L}_2(\pi)$  where

$$\mathcal{L}_{UGW}(\pi) = \sum_{i,j,k,l} \left( d_{ik}^{(X)} - d_{jl}^{(Y)} \right)^2 \pi_{ij}\pi_{kl} + \rho \text{KL}(\pi_1 \otimes \pi_1 | \alpha \otimes \alpha) \\ + \rho \text{KL}(\pi_2 \otimes \pi_2 | \beta \otimes \beta).$$

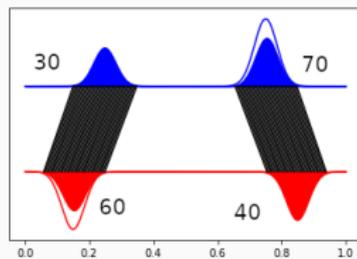
To be compared with

$$\mathcal{G}(\pi) = \sum_{i,j,k,l} \left( d_{ik}^{(X)} - d_{jl}^{(Y)} \right)^2 \pi_{ij}\pi_{kl}, \\ \mathcal{L}_{UOT}(\pi) = \sum_{i,j} C_{ij}\pi_{ij} + \rho \text{KL}(\pi_1 | \alpha) + \rho \text{KL}(\pi_2 | \beta).$$

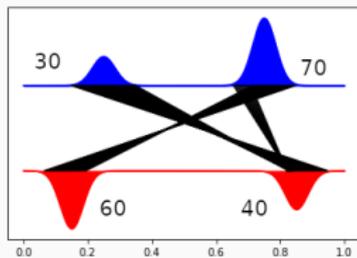
# Numeric in dimension 1



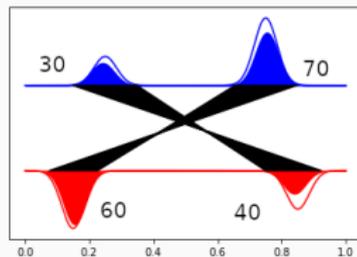
Balanced OT



UOT-KL

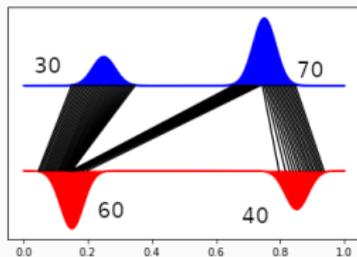


GW

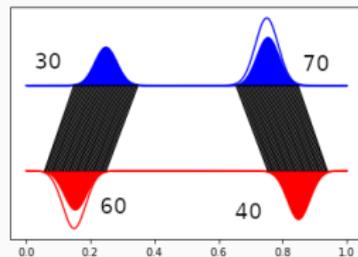


UGW

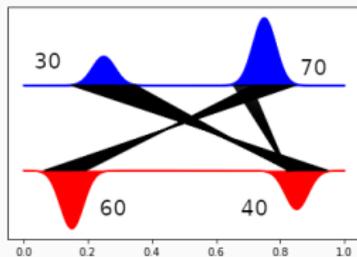
# Take home message



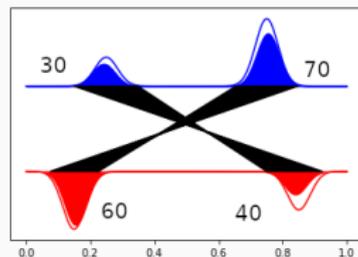
**1 to 1**



**1 to 1  $\oplus$  reweighting**



**1 to 1  $\oplus$  isometry**



**1 to 1  $\oplus$  isom.  $\oplus$  rew.**

**We can enrich assignments with a variety of priors.**

# Theoretical results and conic formulation

- UOT is not convenient to prove the triangle inequality.
  - Need to use another formulation called "conic" (COT)
- COT = OT on a lifted space  $\mathfrak{C} = X \times \mathbb{R}_+$
- **Thm 1:** UOT is definite.
  - **Thm 2:** COT is a distance between positive measures.
  - **Thm 3:** One has UOT = COT.

## Theorem [S., Vialard, Peyré]

1. UGW is definite up to isometries.
2. There exists a conic formulation CGW which is a distance between mm-spaces up to isometry.
3. One has  $UGW \geq CGW$ .

# Implementation of UGW and experiments

---

**Idea:** Entropic regularization + alternate minimization

$$\begin{aligned}\text{UGW}_\varepsilon(\mathcal{X}, \mathcal{Y}) &\stackrel{\text{def.}}{=} \inf_{\pi \geq 0} \mathcal{L}_{\text{UGW}}(\pi) + \varepsilon \text{KL}(\pi \otimes \pi, (\alpha \otimes \beta)^{\otimes 2}) \\ &\geq \inf_{\pi, \gamma \geq 0} \mathcal{F}(\pi, \gamma) + \varepsilon \text{KL}(\pi \otimes \gamma, (\alpha \otimes \beta)^{\otimes 2}),\end{aligned}$$

$$\begin{aligned}\text{where } \mathcal{F}(\pi, \gamma) &\stackrel{\text{def.}}{=} \sum_{i,j,k,l} \left( d_{ik}^{(X)} - d_{jl}^{(Y)} \right)^2 \pi_{ij} \gamma_{kl} \\ &\quad + \rho \text{KL}(\pi_1 \otimes \gamma_1, \alpha \otimes \alpha) + \rho \text{KL}(\pi_2 \otimes \gamma_2, \beta \otimes \beta)\end{aligned}$$

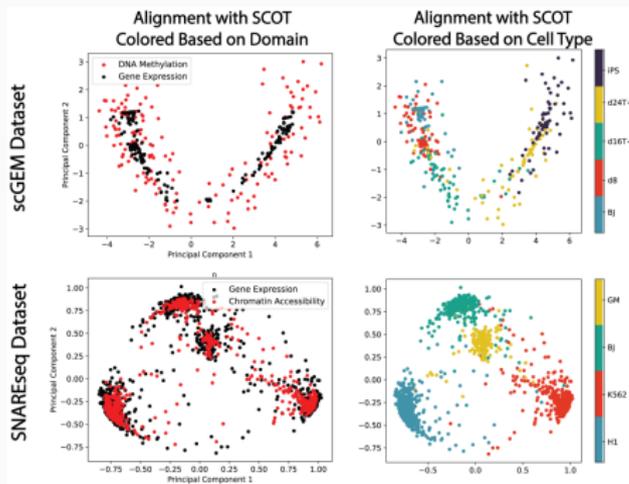
## Proposition - alternate descent $\leftrightarrow$ solve UOT

For a fixed  $\gamma$ ,  $\pi \in \arg \min_{\pi} \mathcal{F}(\pi, \gamma) + \varepsilon \text{KL}(\pi \otimes \gamma | (\alpha \otimes \beta)^{\otimes 2})$  is the solution of

$$\min_{\pi} \sum_{i,j} \tilde{c}_{ij} \pi_{ij} + \tilde{\rho} \text{KL}(\pi_1 | \alpha) + \tilde{\rho} \text{KL}(\pi_2 | \beta) \\ + \tilde{\varepsilon} \text{KL}(\pi | \alpha \otimes \beta),$$

where  $(\tilde{c}, \tilde{\rho}, \tilde{\varepsilon})$  depend on the fixed measure  $\gamma$  via a computable formula.

# Back to genomic data alignment<sup>8</sup>



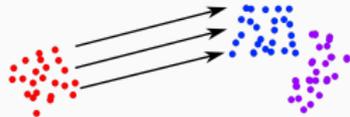
From Demetci et al.

- GW reaches state of the art performance for RNASeq data.
- UGW improves the classification performance over GW.
- See rsinghlab/SCOT on Github.

<sup>8</sup>Demetci, Pinar, et al. Gromov-Wasserstein optimal transport to align single-cell multi-omics data.

# Application - Positive Unlabeled learning

- Domain adaptation = Propagate labels on a similar dataset.
- PU learning = supervised learning but we learn from only one class.
- **Idea:** Use a transport plan to map positive samples to unlabeled positive ones.



# Performance results

Dataset	prior	Init (PW)	PGW	UGW	Dataset	prior	Init (FLB)	PGW	UGW
surf-C → surf-C	0.1	<b>89.9</b>	84.9	83.9	surf-C → decaf-C	0.1	85.0	85.1	<b>85.6</b>
surf-C → surf-A	0.1	81.8	82.2	<b>83.5</b>	surf-C → decaf-A	0.1	84.2	<b>87.1</b>	83.6
surf-C → surf-W	0.1	<b>81.9</b>	81.3	80.3	surf-C → decaf-W	0.1	86.2	<b>88.6</b>	86.8
surf-C → surf-D	0.1	80.0	81.4	<b>83.2</b>	surf-C → decaf-D	0.1	84.7	<b>91.1</b>	90.7
surf-C → surf-C	0.2	<b>79.7</b>	75.7	75.4	surf-C → decaf-C	0.2	74.8	75.6	<b>75.9</b>
surf-C → surf-A	0.2	65.6	66.0	<b>76.4</b>	surf-C → decaf-A	0.2	76.2	<b>87.9</b>	82.4
surf-C → surf-W	0.2	65.1	64.3	<b>67.3</b>	surf-C → decaf-W	0.2	81.5	88.4	<b>89.9</b>
decaf-C → decaf-C	0.1	<b>93.9</b>	83.0	86.8	decaf-C → surf-C	0.1	<b>81.7</b>	81.0	81.1
decaf-C → decaf-A	0.1	80.1	81.4	<b>85.6</b>	decaf-C → surf-A	0.1	80.9	81.2	<b>82.4</b>
decaf-C → decaf-W	0.1	80.1	82.7	<b>86.1</b>	decaf-C → surf-W	0.1	82.0	81.3	<b>83.5</b>
decaf-C → decaf-D	0.1	80.6	<b>83.8</b>	83.4	decaf-C → surf-D	0.1	80.0	80.8	<b>81.5</b>
decaf-C → decaf-C	0.2	<b>90.6</b>	76.7	80.5	decaf-C → surf-C	0.2	<b>66.6</b>	63.7	65.2
decaf-C → decaf-A	0.2	62.5	68.7	<b>74.7</b>	decaf-C → surf-A	0.2	62.9	62.4	<b>69.3</b>
decaf-C → decaf-W	0.2	65.7	75.9	<b>79.2</b>	decaf-C → surf-W	0.2	65.1	61.4	<b>83.3</b>

**Table 1:** Accuracy for all tasks. The left block are domain adaptation experiments with similar features, where both PGW and UGW are initialised with PW. The right block are domain adaptation experiments with different features, and the reported init is FLB used for UGW.

# Conclusion

- Flexibility of UOT models through  $(C, \rho, \varepsilon) + \text{KL} \rightsquigarrow D_\varphi$

- Blending of UOT with GW distances
- Computations on GPUs  $\rightarrow$  UGW
- Theoretical aspects  $\rightarrow$  CGW distance

- $\text{UOT}_{\varepsilon, \rho}$  is fast to compute but lost OT properties,
- Possible to "debias"  $\text{UOT}_{\varepsilon, \rho}$  to retrieve some of them.
- **Open question:** Can we debias  $\text{UGW}_\varepsilon$ ? Which properties ?

## Implementations - github repositories

- thibsej/unbalanced-ot-functionals
- jeanfeudy/geomloss
- thibsej/unbalanced\_gromov\_wasserstein

## References

- Feydy, J., Séjourné, T., Vialard, F. X., Amari, S. I., Trounev, A., & Peyré, G. (2019). Interpolating between optimal transport and MMD using Sinkhorn divergences.
- Séjourné, T., Feydy, J., Vialard, F. X., Trounev, A., & Peyré, G. (2019). Sinkhorn Divergences for Unbalanced Optimal Transport.
- Séjourné, T., Vialard, F. X., & Peyré, G. (2020). The Unbalanced Gromov Wasserstein Distance: Conic Formulation and Relaxation.

**Thank you !**

# Correcting the entropic bias - Sinkhorn divergence

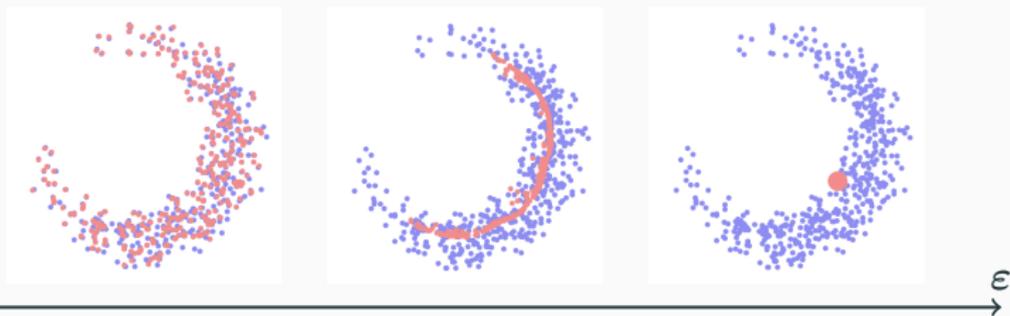
---

# Entropic bias

**Problem:**  $\mathcal{L} = \text{OT}_\varepsilon$  does not retrieve  $\beta$  for  $\varepsilon > 0$ .

Not a distance:  $\text{OT}_\varepsilon(\alpha, \alpha) > 0$ ,

$\exists \alpha \in \mathcal{M}_1^+(\mathcal{X}), \text{OT}_\varepsilon(\alpha, \beta) < \text{OT}_\varepsilon(\beta, \beta)$ .



$\Rightarrow$  **One cannot crossvalidate the parameter  $\varepsilon$ .**

# Unbalanced Sinkhorn Divergence

## Definition

Setting  $m(\mu) = \sum_i \mu_i$ , we define

$$S_{\varepsilon, \rho}(\alpha, \beta) \stackrel{\text{def.}}{=} \text{UOT}_{\varepsilon, \rho}(\alpha, \beta) - \frac{1}{2} \text{UOT}_{\varepsilon, \rho}(\alpha, \alpha) - \frac{1}{2} \text{UOT}_{\varepsilon, \rho}(\beta, \beta) + \frac{\varepsilon}{2} (m(\alpha) - m(\beta))^2.$$

It extends the balanced Sinkhorn divergence<sup>9 10</sup>.

**Remark:** When  $\alpha = \beta$ , one has  $S_{\varepsilon, \rho}(\alpha, \beta) = 0$ .

**Is it positive ? Definite ? Smooth ?**

---

<sup>9</sup>Ramdas, A., Trillos, N. G., & Cuturi, M. (2017). On wasserstein two-sample testing and related families of nonparametric tests.

<sup>10</sup>Genevay, A., Peyré, G., & Cuturi, M. (2018, March). Learning generative models with sinkhorn divergences.

## Theorem [S., Feydy, Vialard, Trounev, Peyre '19]

For any Lipschitz cost  $C$  on a compact set s.t.  $k_\varepsilon \stackrel{\text{def.}}{=} e^{-\frac{C}{\varepsilon}}$  is a positive universal kernel, for any  $\varepsilon > 0$

- $S_{\varepsilon, \rho}$  is convex, positive, definite.
- It is (weakly) differentiable.
- One has  $S_{\varepsilon, \rho}(\alpha, \beta) \rightarrow 0 \Leftrightarrow \alpha \rightarrow \beta$ .

**Corollary:** holds for  $C(x, y) = \|\psi(x) - \psi(y)\|_2^2$ , for  $\psi$  neural net.

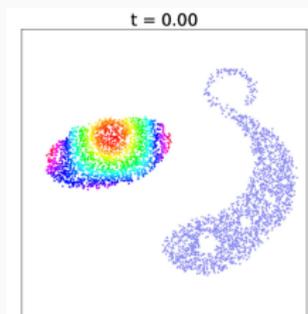
# Numerical insights on UOT and the Sinkhorn divergence

---

# Numerical experiments model

Setting adapted from [Chizat '19]<sup>11</sup>.

- Position/mass parameterization  
 $\theta = \{(x_i, r_i)_i\} \in (\mathbb{R}^d \times \mathbb{R}_+)^n$
- Model measure  $\theta \mapsto \alpha(\theta) = \sum_i^n r_i^2 \delta_{x_i}$
- Minimize  $\mathcal{L}(\alpha(\theta), \beta)$  w.r.t.  $\theta$



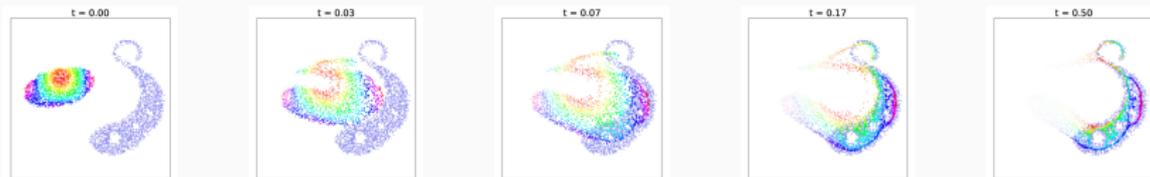
## Updates of the coordinates

$$\begin{aligned}x_i^{(t+1)} &= x_i^{(t)} - \eta_x \nabla_{x_i} \mathcal{L}(\alpha(\theta^{(t)}), \beta), \\r_i^{(t+1)} &= r_i^{(t)} \cdot \exp(-2\eta_r \nabla_{r_i} \mathcal{L}(\alpha(\theta^{(t)}), \beta))\end{aligned}$$

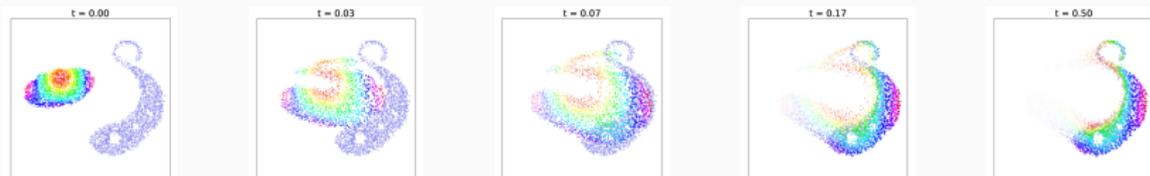
<sup>11</sup>Chizat, L. (2019). Sparse optimization on measures with over-parameterized gradient descent.

# Numerics - Gradient descent

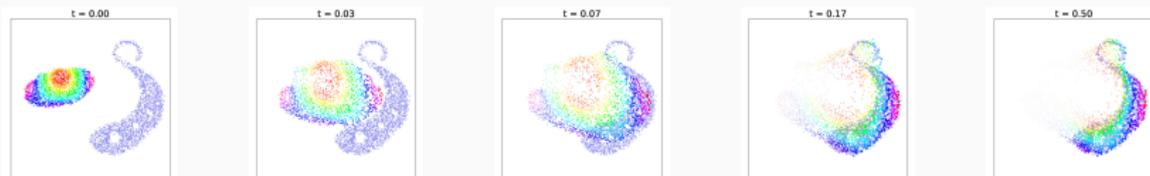
Parameters:  $C(x, y) = \|x - y\|_2^2$  on  $[0, 1]^2$ ,  $\rho = 0.3$ ,  $\eta_x = 60.0$ ,  $\eta_r = 0.3$



$$\mathcal{L} = \text{UOT}_{\epsilon, \rho}, \epsilon = 10^{-3}$$



$$\mathcal{L} = S_{\epsilon, \rho}, \epsilon = 10^{-3}$$



$$\mathcal{L} = S_{\epsilon, \rho}, \epsilon = 10^{-2}$$

## Supplementary slides

---

Define  $\varphi : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  l.s.c., convex,  $\varphi(1) = 0$ ,  $\varphi'^\infty = \lim_{x \rightarrow \infty} \frac{\varphi(x)}{x}$ .

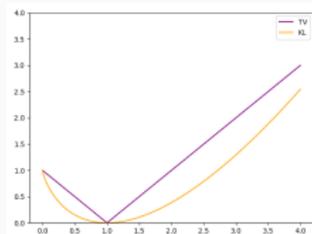
Write  $\alpha = \sum_i \alpha_i \delta_{x_i}$  and  $\beta = \sum_i \beta_i \delta_{x_i}$  (Same support  $(x_i)$ )

## Definition - $\varphi$ -divergence

$$D_\varphi(\alpha, \beta) = \sum_{\beta_i \neq 0} \varphi\left(\frac{\alpha_i}{\beta_i}\right) \beta_i + \varphi'^\infty \sum_{\beta_i = 0} \alpha_i.$$

## Examples:

- $\text{KL}(\alpha, \beta) = \sum_i (\log(\frac{\alpha_i}{\beta_i}) \alpha_i - \alpha_i + \beta_i)$ :  
 $\varphi(x) = x \log x - x + 1$ ,
- $\text{TV}(\alpha, \beta) = \sum_i |\alpha_i - \beta_i|$ :  $\varphi(x) = |x - 1|$ .



<sup>12</sup>Csiszàr, I. (1967). Information-type measures of difference of probability distributions and indirect observation.

## Alternate UGW = sequence of Sinkhorn updates

- Focus on  $\lambda(t) = t^2$  for improved time and memory complexity
- Focus on  $D_\varphi = \text{KL}$  which verifies

$$\begin{aligned} \text{KL}(\mu \otimes \nu, \alpha \otimes \beta) &= m(\nu)\text{KL}(\mu, \alpha) + m(\mu)\text{KL}(\nu, \beta) \\ &\quad + (m(\mu) - m(\alpha))(m(\nu) - m(\beta)). \end{aligned}$$

⇒ Given  $\gamma$ , minimizing w.r.t.  $\pi$  amounts to solve a regularized UOT problem.

---

## Algorithm 1 – UGW( $\mathcal{X}, \mathcal{Y}, \rho, \varepsilon$ )

---

**Input:** mm-spaces  $(\mathcal{X}, \mathcal{Y})$ , relaxation  $\rho$ , regularization  $\varepsilon$

**Output:** approximation  $(\pi, \gamma)$  minimizing  $\mathcal{F} + \varepsilon \text{KL}^\otimes$

- 1: Initialize  $(\pi, \gamma)$  and  $(f, g)$
  - 2: **while**  $(\pi, \gamma)$  has not converged **do**
  - 3:     Update  $\gamma \leftarrow \pi$  and compute the cost  $\tilde{c} \leftarrow c^{\varepsilon, \gamma}$
  - 4:     Update parameters  $(\tilde{\rho}, \tilde{\varepsilon}) \leftarrow (m(\pi)\rho, m(\pi)\varepsilon)$
  - 5:     Compute  $(f, g)$  that solves UOT( $\mu, \nu, \tilde{c}, \tilde{\rho}, \tilde{\varepsilon}$ )
  - 6:     Update  $\gamma_{ij} \leftarrow \exp \left[ (f_i + g_j - \tilde{c}_{ij}) / \tilde{\varepsilon} \right] \alpha_i \beta_j$
  - 7:     Rescale  $\gamma \leftarrow \sqrt{m(\pi) / m(\gamma)} \gamma$
  - 8: **Return**  $(\pi, \gamma)$ .
-

## Detailed algorithm

---

### Algorithm 2 – UGW( $\mathcal{X}, \mathcal{Y}, \rho, \varepsilon$ )

---

**Input:** mm-spaces  $(\mathcal{X}, \mathcal{Y})$ , relaxation  $\rho$ , regularization  $\varepsilon$

**Output:** approximation  $(\pi, \gamma)$  minimizing  $\mathcal{F} + \varepsilon \text{KL}^\otimes$

- 1: Initialize  $\pi = \gamma = \mu \otimes \nu / \sqrt{m(\mu)m(\nu)}$ ,  $g = 0$ .
  - 2: **while**  $(\pi, \gamma)$  has not converged **do**
  - 3:     Update  $\pi \leftarrow \gamma$ , then  $c \leftarrow c_\pi^\varepsilon$ ,  $\tilde{\rho} \leftarrow m(\pi)\rho$ ,  $\tilde{\varepsilon} \leftarrow m(\pi)\varepsilon$
  - 4:     **while**  $(f, g)$  has not converged **do**
  - 5:          $\forall x, f(x) \leftarrow -\frac{\tilde{\varepsilon}\tilde{\rho}}{\tilde{\varepsilon}+\tilde{\rho}} \log \left( \int e^{(g(y)-c(x,y))/\tilde{\varepsilon}} d\nu(y) \right)$
  - 6:          $\forall y, g(y) \leftarrow -\frac{\tilde{\varepsilon}\tilde{\rho}}{\tilde{\varepsilon}+\tilde{\rho}} \log \left( \int e^{(f(x)-c(x,y))/\tilde{\varepsilon}} d\mu(x) \right)$
  - 7:         Update  $\gamma(x, y) \leftarrow \exp \left[ (f(x) + g(y) - c(x, y))/\tilde{\varepsilon} \right] \mu(x)\nu(y)$
  - 8:         Rescale  $\gamma \leftarrow \sqrt{m(\pi)/m(\gamma)}\gamma$
  - 9:     Return  $(\pi, \gamma)$ .
-