

Reconstruction de trajectoires cellulaires par Transport Optimal

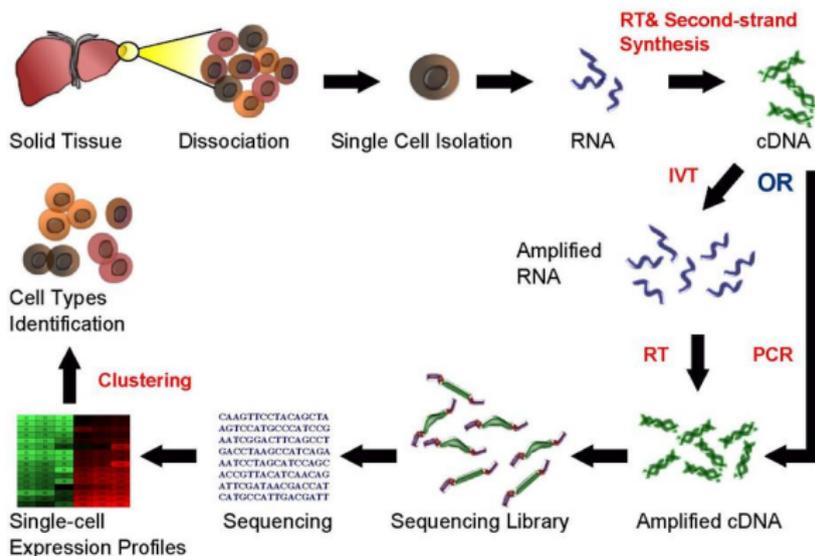
David ROBIN

Stage de Licence

Juin-Août 2018

Séquençage ARN sur cellule unique

Single Cell RNA Sequencing Workflow

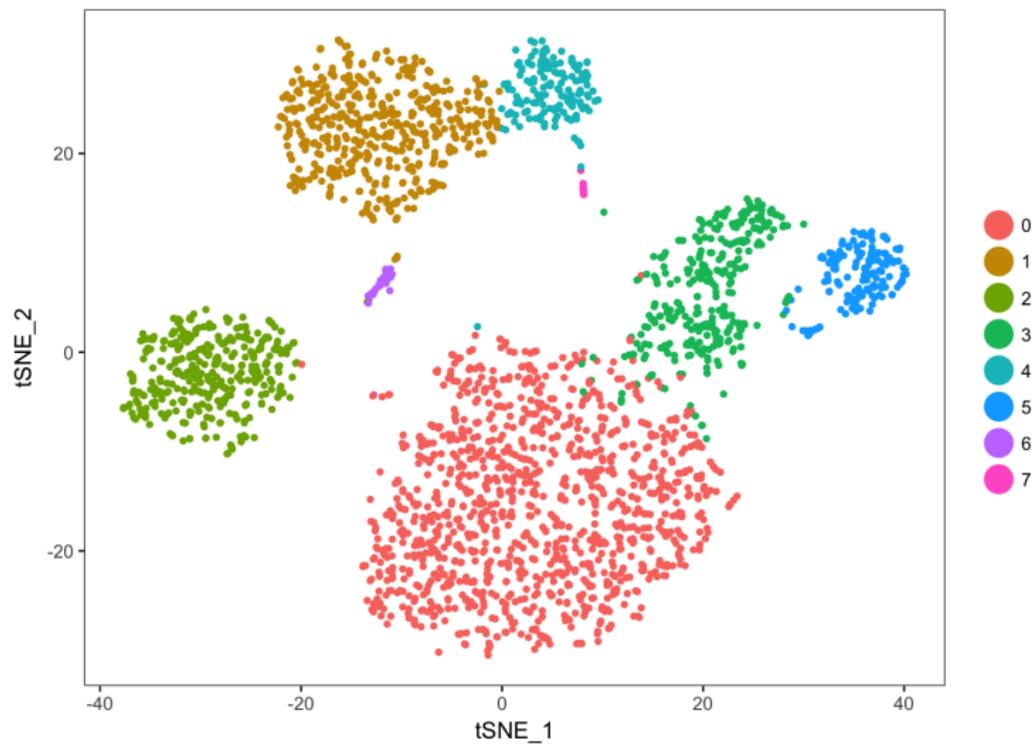


Processus scRNA-Seq destructif → pas d'évolution temporelle

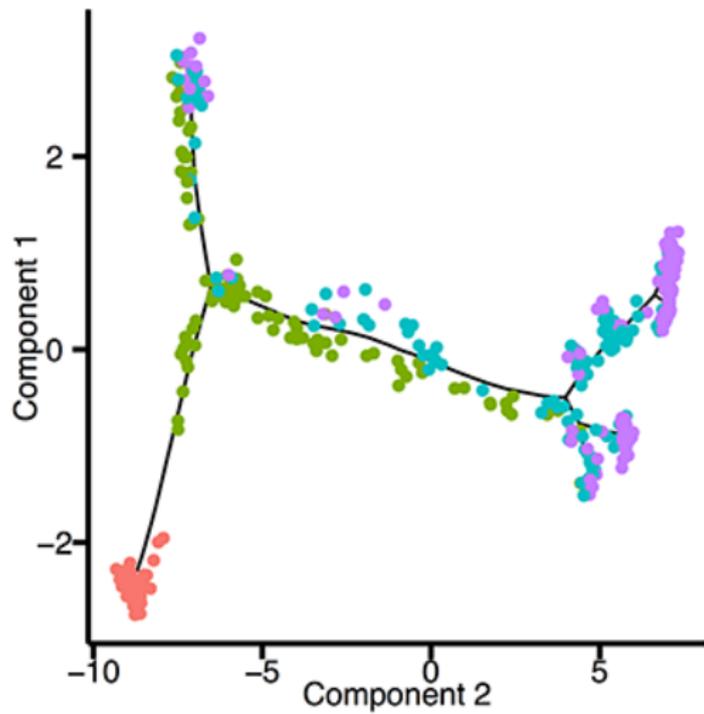
Espace d'expression génétique

	Cellule 1	Cellule 2	Cellule 3	Cellule 4	...
gene 1	32	33	124	10	...
gene 2	12	14	174	14	...
gene 3	14	12	12	12	...
gene 4	121	131	12	14	...
⋮	⋮	⋮	⋮	⋮	

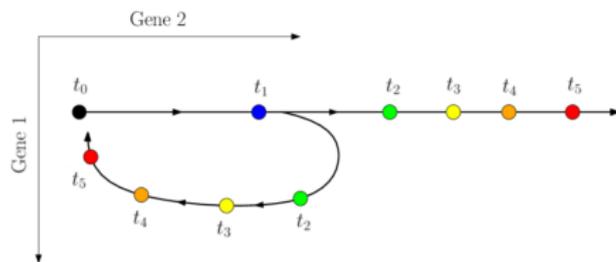
Espace d'expression génétique



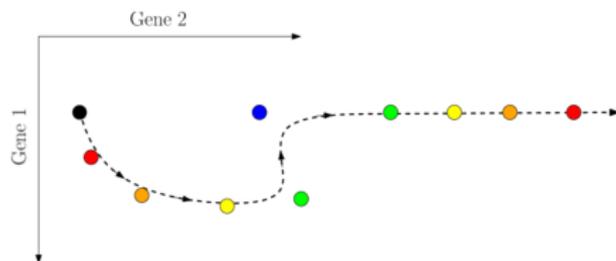
Pseudotemps



Échec du pseudotemps



Trajectoire réelle



Trajectoire inférée

Modélisation

- ▶ Ensemble des gènes: G
- ▶ Cellule: $x \in \mathbb{R}^G$
- ▶ Population: $\mathbb{P}_t = \sum_{i < n} \delta_{x_i}, \quad x_i \in \mathbb{R}^G$

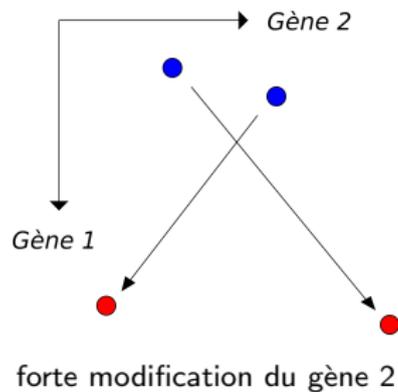
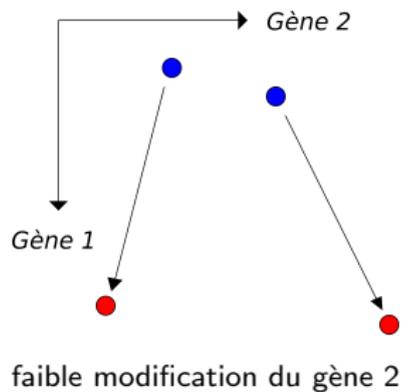
Modélisation

- ▶ Ensemble des gènes: G
- ▶ Cellule: $x \in \mathbb{R}^G$
- ▶ Population: $\mathbb{P}_t = \sum_{i < n} \delta_{x_i}, \quad x_i \in \mathbb{R}^G$
- ▶ Couplages entre \mathbb{P} et \mathbb{Q} : $\pi \in \Gamma(\mathbb{P}, \mathbb{Q})$

$$\forall A \subset \mathbb{R}^G, \quad \int_{x \in A} \int_{y \in \mathbb{R}^G} \pi(x, y) \cdot dx dy = \mathbb{P}(A)$$

$$\forall B \subset \mathbb{R}^G, \quad \int_{y \in B} \int_{x \in \mathbb{R}^G} \pi(x, y) \cdot dx dy = \mathbb{Q}(B)$$

L'hypothèse du Transport Optimal



Couplage et carte de transport

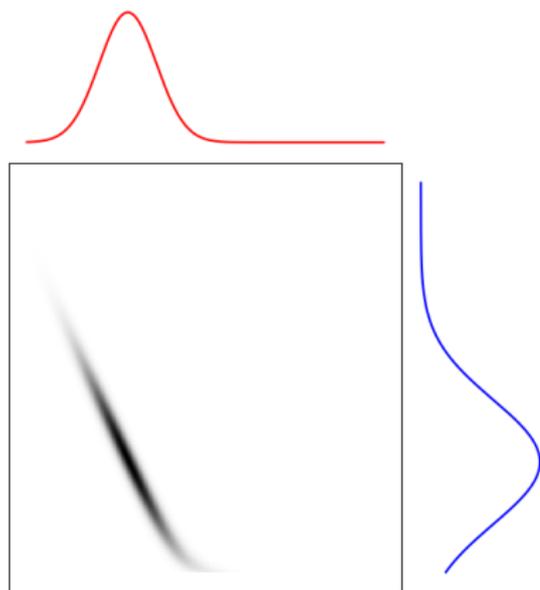


Image par une carte π :

$$\pi_{\#} : \mu \mapsto \int_{x \in \mathbb{R}^G} \frac{\pi(x, \cdot)}{\mathbb{P}(x)} \cdot d\mu(x)$$

$$\pi_{\#} \mathbb{P} = \mathbb{Q}$$

Pré-image par une carte π :

$$\pi^{\#} : \nu \mapsto \int_{y \in \mathbb{R}^G} \frac{\pi(\cdot, y)}{\mathbb{Q}(y)} \cdot d\nu(y)$$

$$\pi^{\#} \mathbb{Q} = \mathbb{P}$$

Recherche de couplage

Coût total du transport:

$$\int_x \int_y c(x,y) \cdot \pi(x,y) \cdot dx dy$$

En version discrète:

$$\langle c, \pi \rangle = \sum_x \sum_y c(x,y) \cdot \pi(x,y)$$

Problème de minimisation:

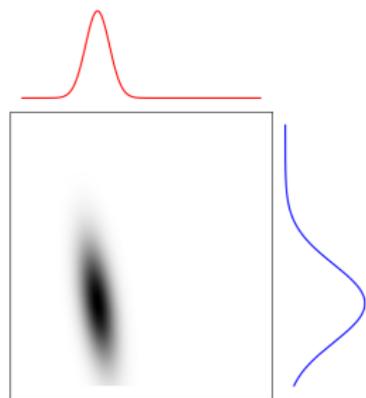
$$\min_{\pi \in \Gamma(\mathbb{P}, \mathbb{Q})} \langle c, \pi \rangle$$

$$\text{où } \Gamma(\mathbb{P}, \mathbb{Q}) = \{ \pi : \pi \cdot \mathbf{1} = \mathbb{P}, \pi^T \cdot \mathbf{1} = \mathbb{Q} \}$$

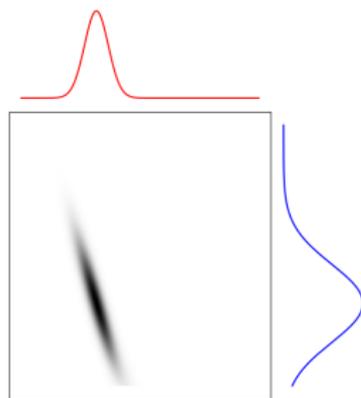
Régularisation

Problème régularisé:

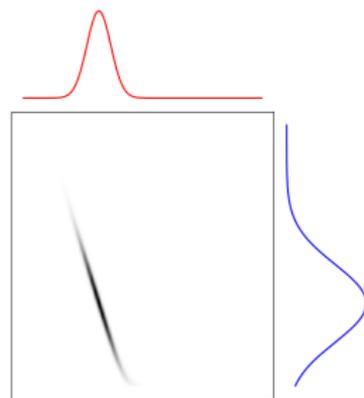
$$\min_{\pi \in \Gamma(\mathbb{P}, \mathbb{Q})} \langle c, \pi \rangle - \varepsilon \mathcal{H}(\pi)$$



$\varepsilon = 10^4$



$\varepsilon = 10^3$



$\varepsilon = 10^2$

Croissance cellulaire

Hypothèse:

la position d'une cellule détermine son taux de croissance

$$r(x,y) \cdot g(x)^{\Delta t}$$

Nouvelles marginales:

$$\mathbb{P} \odot g^{\Delta t} \quad \text{et} \quad \bar{g} \cdot \mathbb{Q}$$

Transport non-équilibré

$$\min_{\pi} \langle c, \pi \rangle - \varepsilon \mathcal{H}(\pi) + \lambda_1 \cdot \text{KL}(\pi \cdot \mathbf{1} \parallel \mathbb{P} \odot g^{\Delta t}) + \lambda_2 \cdot \text{KL}(\pi^T \cdot \mathbf{1} \parallel \bar{g} \mathbb{Q})$$

avec KL la divergence de Kullback-Leibler:

$$\text{KL}(P \parallel Q) = \sum_x P(x) \cdot \log \frac{P(x)}{Q(x)}$$

Algorithme de Sinkhorn

$$\min_{\pi \in \Gamma(\mathbb{P}, \mathbb{Q})} \langle c, \pi \rangle - \varepsilon \mathcal{H}(\pi) = \min_{\pi \in \Gamma(\mathbb{P}, \mathbb{Q})} \varepsilon \text{KL}(\pi \parallel K)$$

$$\text{où } K = \exp(-c/\varepsilon)$$

On projette alternativement sur \mathcal{C}_1 et \mathcal{C}_2 :

$$\mathcal{C}_1 = \{ \pi : \pi \cdot \mathbf{1} = \mathbb{P} \} \quad \text{et} \quad \mathcal{C}_2 = \{ \pi : \pi^T \cdot \mathbf{1} = \mathbb{Q} \}$$

$$\mathcal{P}_{\mathcal{C}_1}(\pi) = \text{diag} \left(\frac{\mathbb{P}}{\pi \cdot \mathbf{1}} \right) \cdot \pi \quad \text{et} \quad \mathcal{P}_{\mathcal{C}_2}(\pi) = \pi \cdot \text{diag} \left(\frac{\mathbb{Q}}{\pi^T \cdot \mathbf{1}} \right)$$

Benamou et al., 2015,

Iterative Bregman projections for Regularized Transportation Problems

Algorithme de Sinkhorn

$$\pi^{(n)} = \text{diag}(a^{(n)}) \cdot K \cdot \text{diag}(b^{(n)})$$

$$a^{(n)} = \frac{\mathbb{P}}{K \cdot b^{(n)}} \quad \text{et} \quad b^{(n+1)} = \frac{\mathbb{Q}}{K^T \cdot a^{(n)}}$$

Extension au transport non-équilibré

$$\min_{\pi} \langle c, \pi \rangle - \varepsilon \mathcal{H}(\pi) + F_1(\pi \cdot \mathbf{1}, \mathbb{P}) + F_2(\pi^T \cdot \mathbf{1}, \mathbb{Q})$$

$$a^{(n)} = \frac{\text{prox}_{F_1/\varepsilon}(K \cdot b^{(n)})}{K \cdot b^{(n)}} \quad \text{et} \quad b^{(n+1)} = \frac{\text{prox}_{F_2/\varepsilon}(K^T \cdot a^{(n)})}{K^T \cdot a^{(n)}}$$

$$\text{où } \text{prox}_{F/\varepsilon}(z) = \underset{s}{\text{argmin}} F(s) + \varepsilon \text{KL}(s \parallel z)$$

Chizat et al., *Scaling algorithms for Unbalanced Transport Problems*, 2016 (ArXiv)

Extension au transport non-équilibré

$$F_1 = \lambda_1 \text{KL}(\cdot \| \mathbb{P}) \quad \text{et} \quad F_2 = \lambda_2 \text{KL}(\cdot \| \mathbb{Q})$$

donnent

$$a^{(n)} = \left(\frac{\mathbb{P}}{K \cdot b^{(n)}} \right)^{\frac{\lambda_1}{\lambda_1 + \varepsilon}} \quad \text{et} \quad b^{(n+1)} = \left(\frac{\mathbb{Q}}{K^T \cdot a^{(n)}} \right)^{\frac{\lambda_2}{\lambda_2 + \varepsilon}}$$

Réduction de dimension

- ▶ Calcul de la matrice de coût: $\mathcal{O}(n^2 G)$
- ▶ *Analyse en composante principale* pour diminuer G
- ▶ Annule les dimensions ne contenant que du bruit

Stabilisation logarithmique

Stabilisation des variables:

$$a = \tilde{a} \cdot \exp(u/\varepsilon) \quad \text{et} \quad b = \tilde{b} \cdot \exp(v/\varepsilon)$$

Noyau stabilisé:

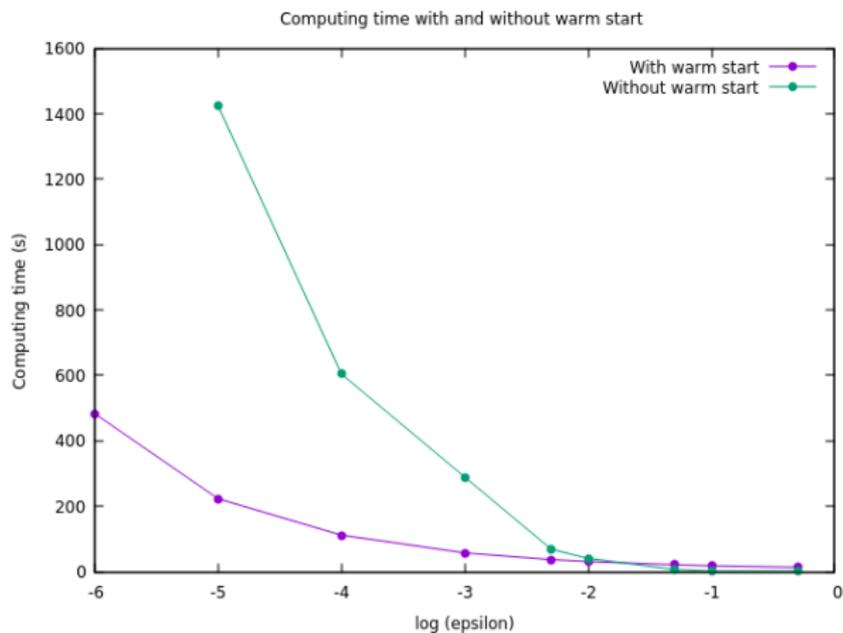
$$\tilde{K}_{i,j} = \exp((u_i + v_j - c_{i,j})/\varepsilon)$$

Nouvelles itérations:

$$\tilde{a}^{(n)} = \left(\frac{\mathbb{P}}{K \cdot \tilde{b}^{(n)}} \right)^{\frac{\lambda_1}{\lambda_1 + \varepsilon}} \odot \exp\left(-\frac{u}{\lambda_1 + \varepsilon}\right)$$

$$\tilde{b}^{(n+1)} = \left(\frac{\mathbb{Q}}{K^T \cdot \tilde{a}^{(n)}} \right)^{\frac{\lambda_2}{\lambda_2 + \varepsilon}} \odot \exp\left(-\frac{v}{\lambda_2 + \varepsilon}\right)$$

Départ à chaud



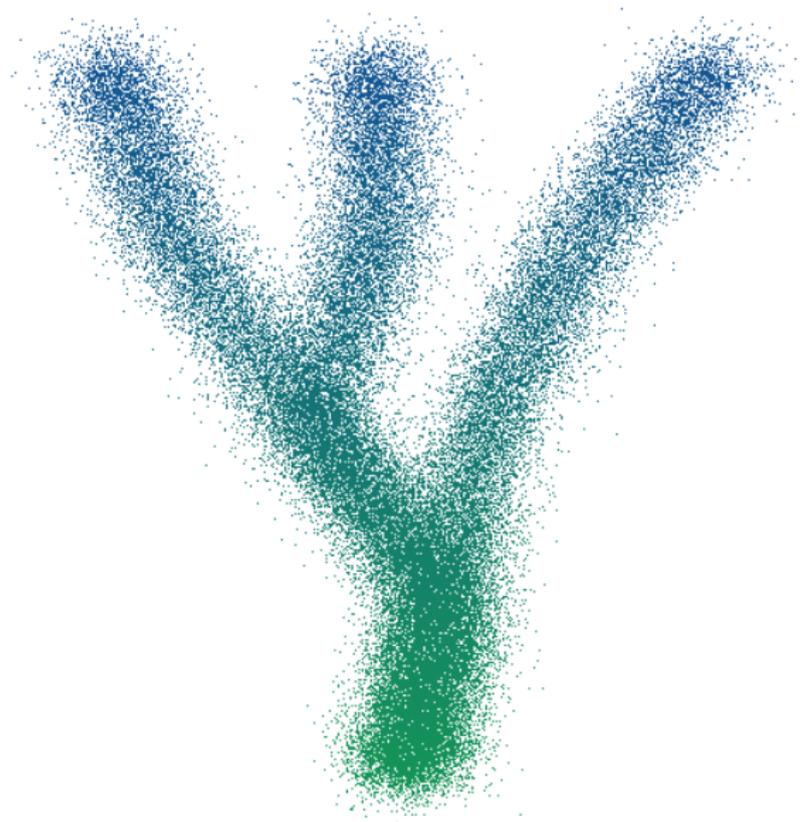
Interpolation

Population interpolée entre P et Q à la fraction α :

$$\mu : I_{i,j} \mapsto \tilde{\pi}_{i,j} = r_{i,j} \cdot g_i^{\alpha \cdot \Delta t}$$

$$\text{où } I_{i,j} = (1 - \alpha) \cdot P_i + \alpha \cdot Q_j$$

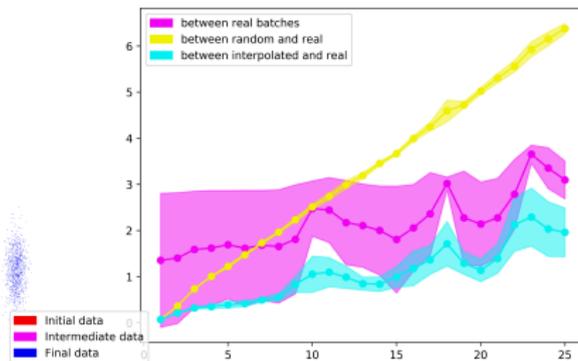
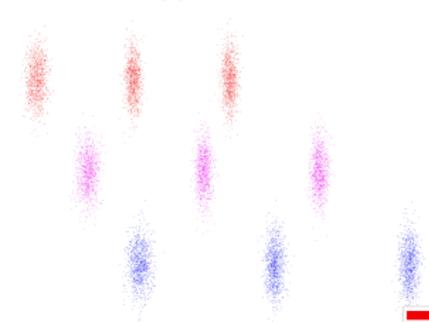
Simulation



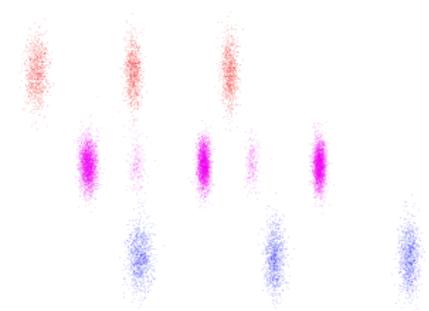
Résultats sur simulation

Optimal transport validation from time points 5.0 to 7.0 ($\epsilon = 0.05$)

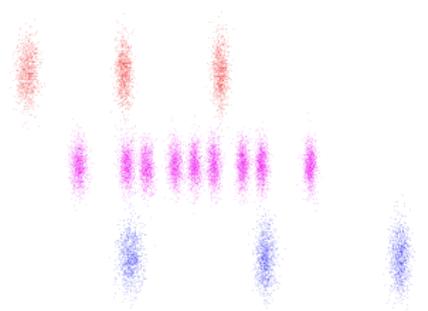
Real population (3237 cells)



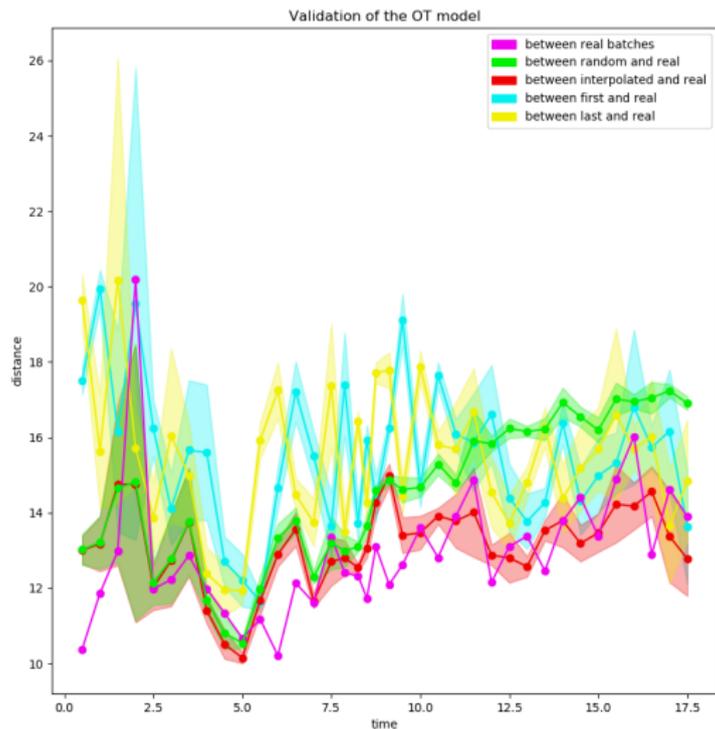
Interpolated population (9984 cells)



Randomized population (9995 cells)



Résultats sur données réelles

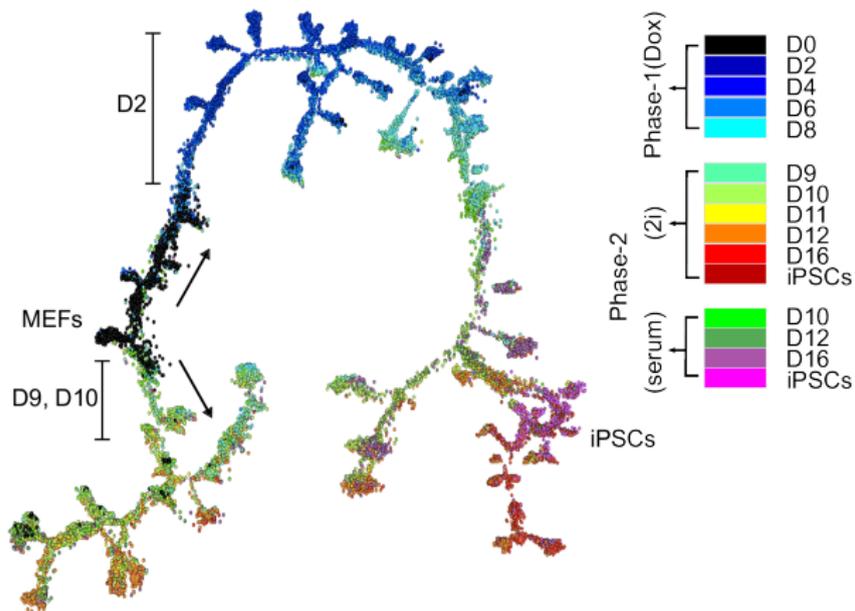


Reconstruction de trajectoires cellulaires

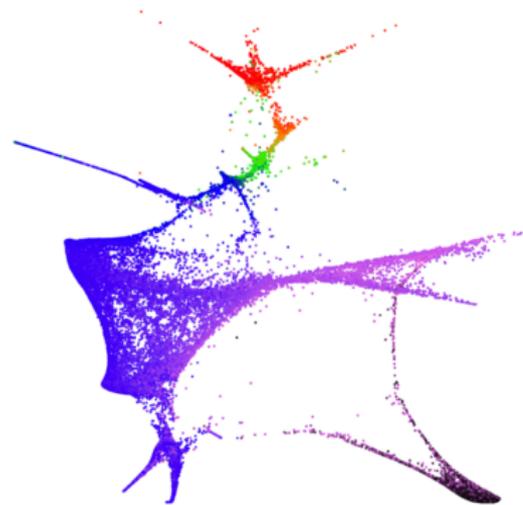
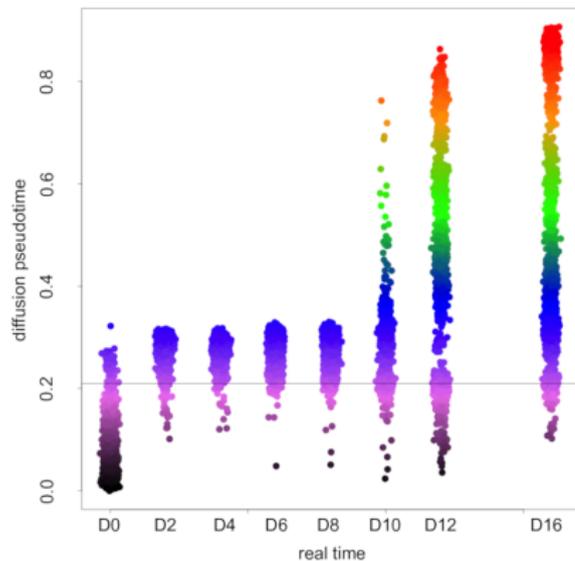
—

Annexes

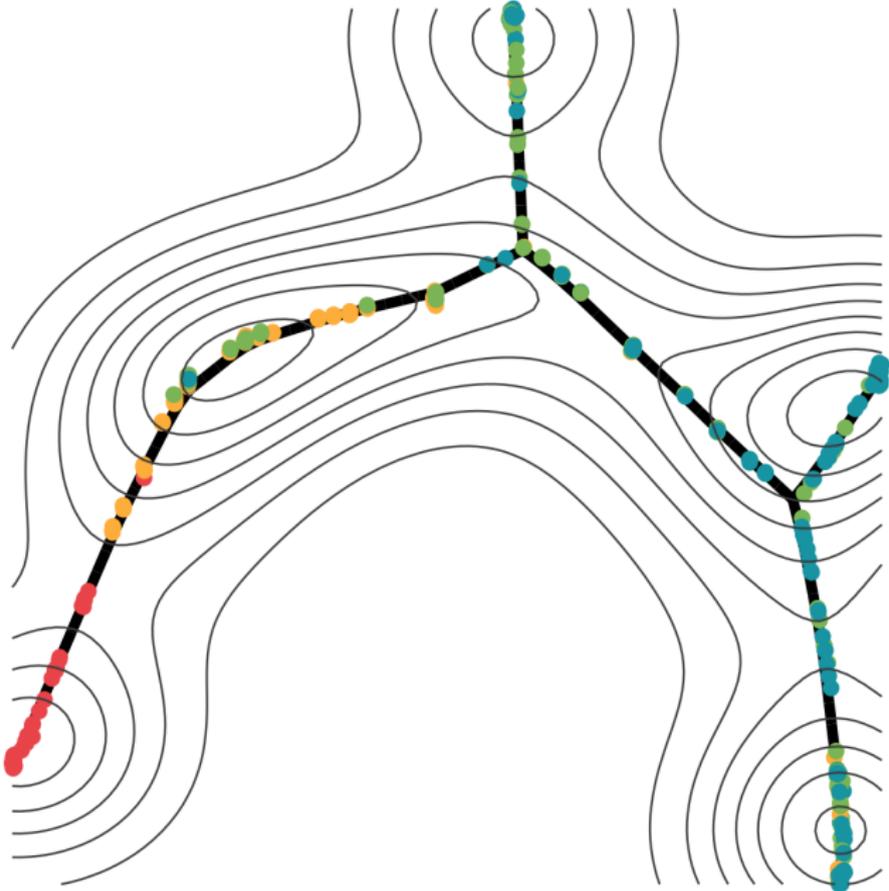
MEF: what goes wrong without time



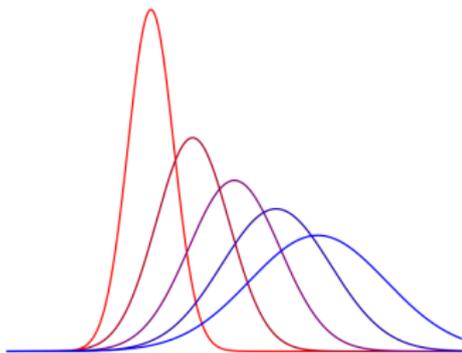
MEF: what goes wrong without time



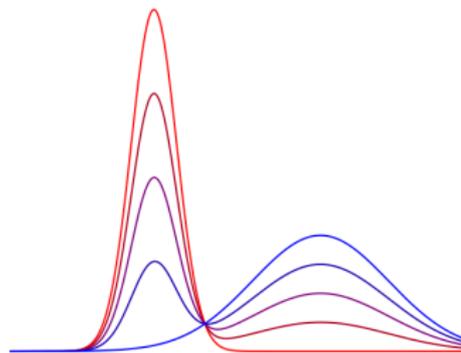
Monocle: Pseudotime



Interpolation



Interpolation par transport

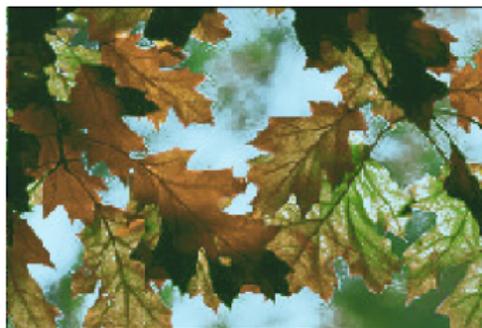


Interpolation par moyenne pondérée

Transfert de couleurs



Image originale



Résultat



Couleurs à transférer

Algorithme de Sinkhorn stabilisé

function SINKHORNSRUT ($c, p, q, \varepsilon, \lambda_1, \lambda_2$)

$$(\tilde{b}, u, v) \leftarrow (1, 0, 0)$$

$$\tilde{K}_{i,j} \leftarrow \exp(-c_{i,j}/\varepsilon) \quad \forall i,j$$

repeat

$$\tilde{a} \leftarrow \left(\frac{p}{K \cdot \tilde{b}} \right)^{\frac{\lambda_1}{\lambda_1 + \varepsilon}} \odot \exp\left(-\frac{u}{\lambda_1 + \varepsilon}\right)$$

$$\tilde{b} \leftarrow \left(\frac{q}{K^T \cdot \tilde{a}} \right)^{\frac{\lambda_2}{\lambda_2 + \varepsilon}} \odot \exp\left(-\frac{v}{\lambda_2 + \varepsilon}\right)$$

if a component of $|\tilde{a}|$ or $|\tilde{b}|$ is “too big” **then**

$$(u, v) \leftarrow (u + \varepsilon \log \tilde{a}, v + \varepsilon \log \tilde{b})$$

$$\tilde{K}_{i,j} \leftarrow \exp((u_i + v_j - c_{i,j}) / \varepsilon) \quad \forall i,j$$

$$\tilde{b} \leftarrow 1$$

until stopping criterion

return $(\tilde{a}_i \tilde{K}_{i,j} \tilde{b}_j)_{i,j}$

Force-directed layout embedding

