

POLINA ARSENTEVA^{1,2} & VINCENT PAGET² & OLIVIER GUIPAUD² & FABIEN MILLIAT² & HERVÉ CARDOT¹ & MOHAMED AMINE BENADJAOU³

¹IMB, Université de Bourgogne

²IRSN, PSE-SANTE/SERAMED/LRMED

³IRSN, PSE-SANTE/SERAMED

I. Introduction

Context: treatment of cancer by radiotherapy.

- Problem:** undesirable effects for the healthy tissues situated in close proximity to the irradiated tumors.
- Ultimate goal:** compare different radiation treatment configurations in order to suggest a treatment associated with minimal risk.

Focus of this work: the dynamic of endothelial cells' response to irradiation.

- Why endothelium?** Key cell compartment for the healthy tissue radiation response and the occurrence of side effects.

Goal:

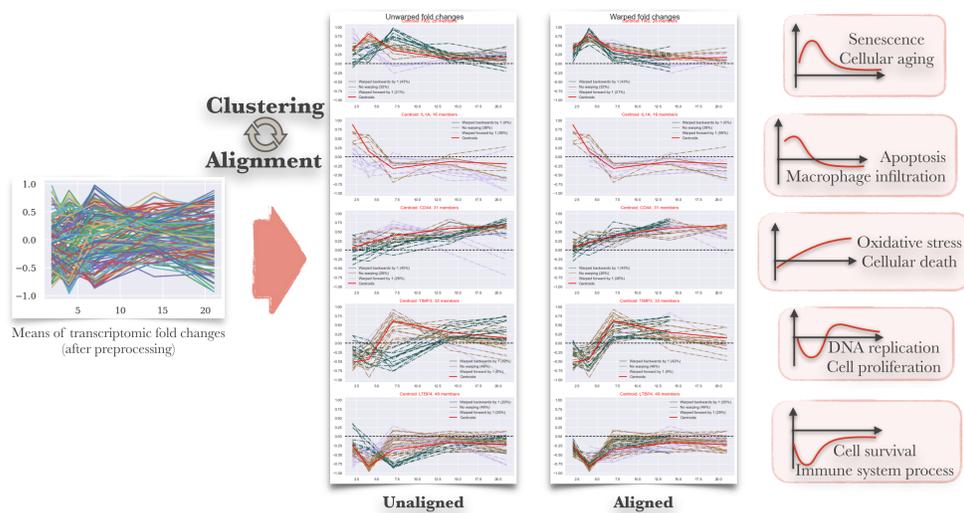
- Extract key features of data describing response to different types of irradiation
- Propose powerful visualisation tools for biological interpretation of the results

Tool:

Clustering & Alignment
Network Inference

IV. Clustering & Alignment of real fold changes

Applied to a real transcriptomic dataset, the procedure allowed to identify **5 distinct gene response types**. Enrichment analysis with Pathway Studio shows the existence of associated distinct **cellular functions**:



V. Network Inference

- Similarity** measure based on the post-warping distances:

$$Sim(\hat{\Gamma}_i, \hat{\Gamma}_{i'}) = \frac{\max_{(a,b) \in \{1, \dots, n_e\}^2} \mathcal{O}W_{ab} - \mathcal{O}W_{ii'}}{\max_{(a,b) \in \{1, \dots, n_e\}^2} \mathcal{O}W_{ab}}$$

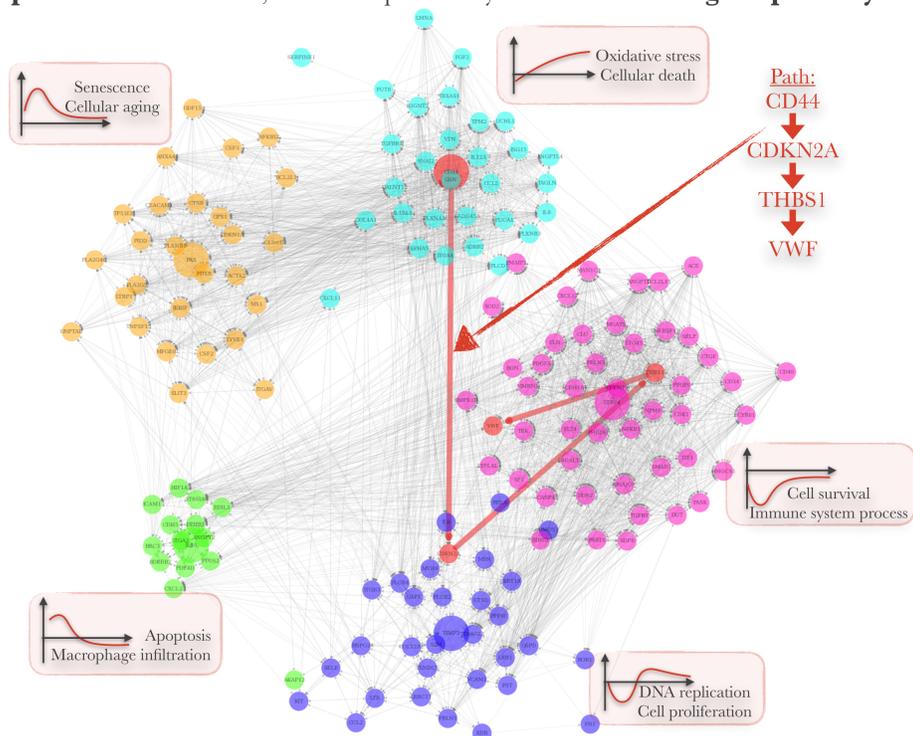
- Binary **adjacency matrix** $X = (X_{ii'})_{(i,i') \in \{1, \dots, n_e\}^2}$ describing directed graph:

$$X_{ii'} = \mathbb{1}_{Sim(\hat{\Gamma}_i, \hat{\Gamma}_{i'}) \geq q} \times \mathbb{1}_{\mathcal{O}W_{ii'} \geq 0}$$

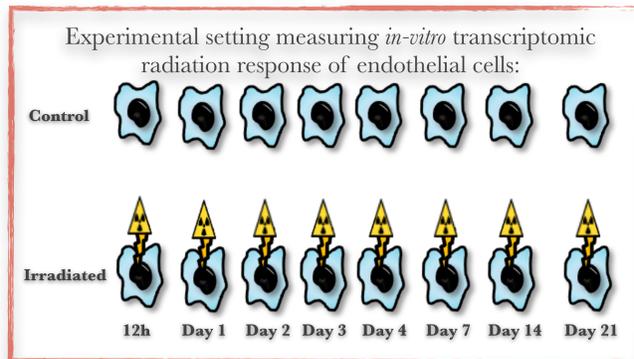
Edge $i \rightarrow i'$ exists if i and i' are at least as similar as the thresholding parameter q

Edge $i \rightarrow i'$ exists if the optimal warp is non-negative

Network is visualised in a **block form**, where blocks are the inferred clusters. Gene **paths** can be extracted, which are potentially indicative of **biological pathways**:



II. From data to radio-induced fold changes



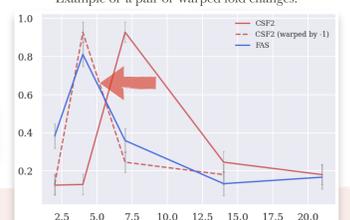
Challenges:

- Independent (destructive sampling) & unequally spaced time points
- Not time series data
- Small number of time points
- Functional data approach not applicable
- ~100-500 observed genes
- Correlations & computational cost to consider
- Multiple replicates
- Uncertainties to consider

Response: Time point $t \in \{t_1, t_2, \dots, t_p\}$
Gene $i \in \{1, 2, \dots, n_e\}$
Experimental condition: $k = \begin{cases} 0 & \text{if control} \\ 1 & \text{if irradiated} \end{cases}$
Replicate $j \in \{1, 2, \dots, n_r\}$

$$\text{Fold Change}(\text{time}) = \text{Effect}_{\text{irradiated}}(\text{time}) - \text{Effect}_{\text{control}}(\text{time})$$

How does alignment (warping) work?
Example of a pair of warped fold changes:



III. Clustering & Alignment

- Fold changes' estimators as **random variables**: for a pair of genes i and i' ,

$$\Gamma_i = \begin{pmatrix} \overline{Y}_{i1}^t - \overline{Y}_{i0}^t & \dots & \overline{Y}_{i1}^{t_p} - \overline{Y}_{i0}^{t_p} \end{pmatrix} \quad \Sigma_{\Gamma_i} = \begin{pmatrix} \sigma_{\Gamma_i}^2 & & 0 \\ & \ddots & \\ 0 & & \sigma_{\Gamma_i}^2 \end{pmatrix} \quad \text{where } \sigma_{\Gamma_i}^2 = \frac{\sum_{j=1}^{n_r} [(Y_{ij}^t - \overline{Y}_{i1}^t)^2 + (Y_{ij}^{t_p} - \overline{Y}_{i0}^{t_p})^2]}{n_r - 1}$$

Joint distribution:

$$\hat{\Gamma}_{ii'} \sim \mathcal{N} \left(\begin{bmatrix} \Gamma_i \\ \Gamma_{i'} \end{bmatrix}, \begin{bmatrix} \Sigma_{\Gamma_i} & K \\ K^T & \Sigma_{\Gamma_{i'}} \end{bmatrix} \right)$$

$$K = \begin{pmatrix} \rho_{\Gamma_i \Gamma_{i'}} & & 0 \\ & \ddots & \\ 0 & & \rho_{\Gamma_i \Gamma_{i'}} \end{pmatrix} \quad \text{where } \rho_{\Gamma_i \Gamma_{i'}} = \frac{\sum_{j=1}^{n_r} [(Y_{ij}^t - \overline{Y}_{i1}^t)(Y_{i'j}^t - \overline{Y}_{i'1}^t) + (Y_{ij}^{t_p} - \overline{Y}_{i0}^{t_p})(Y_{i'j}^{t_p} - \overline{Y}_{i'0}^{t_p})]}{n_r - 1}$$

- Distance \hat{d}_2^2 :** L^2 -distance between normal r. v. with estimated joint distribution:

$$\hat{d}_2^2(\hat{\Gamma}_i, \hat{\Gamma}_{i'}) = \sum_{l=1}^p (\Gamma_i^l - \Gamma_{i'}^l)^2 + \sum_{l=1}^p \sigma_{\Gamma_i}^2 + \sum_{l=1}^p \sigma_{\Gamma_{i'}}^2 - 2 \sum_{l=1}^p \rho_{\Gamma_i \Gamma_{i'}}$$

- Base algorithm:** k-medoids initiated with k-means++.

- Integrating alignment:** clustering of warped fold changes pairs

$$\begin{bmatrix} \hat{\Gamma}_i \circ \mathcal{W}_s \\ \hat{\Gamma}_{i'} \circ \mathcal{W}_s \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \Gamma_i \circ \mathcal{W}_s \\ \Gamma_{i'} \circ \mathcal{W}_s \end{bmatrix}, \begin{bmatrix} \Sigma_{\Gamma_i} \circ \mathcal{W}_s & K \circ \mathcal{W}_s \\ (K \circ \mathcal{W}_s)^T & \Sigma_{\Gamma_{i'}} \circ \mathcal{W}_s \end{bmatrix} \right)$$

Time warp aligning a pair of time vectors by step s : $\mathcal{W}_s: \mathcal{T}^2 \rightarrow \mathcal{T}^2$

- Replacing \hat{d}_2^2** by a more general distance $\widehat{\text{diss}}$:

$$\widehat{\text{diss}}(\hat{\Gamma}_i \circ \mathcal{W}_s, \hat{\Gamma}_{i'} \circ \mathcal{W}_s) = \|\hat{\Gamma}_i \circ \mathcal{W}_s - \hat{\Gamma}_{i'} \circ \mathcal{W}_s\|^2 + \sum (\Sigma_{\Gamma_i} \circ \mathcal{W}_s) + \sum (\Sigma_{\Gamma_{i'}} \circ \mathcal{W}_s) - 2 \sum (K \circ \mathcal{W}_s)$$

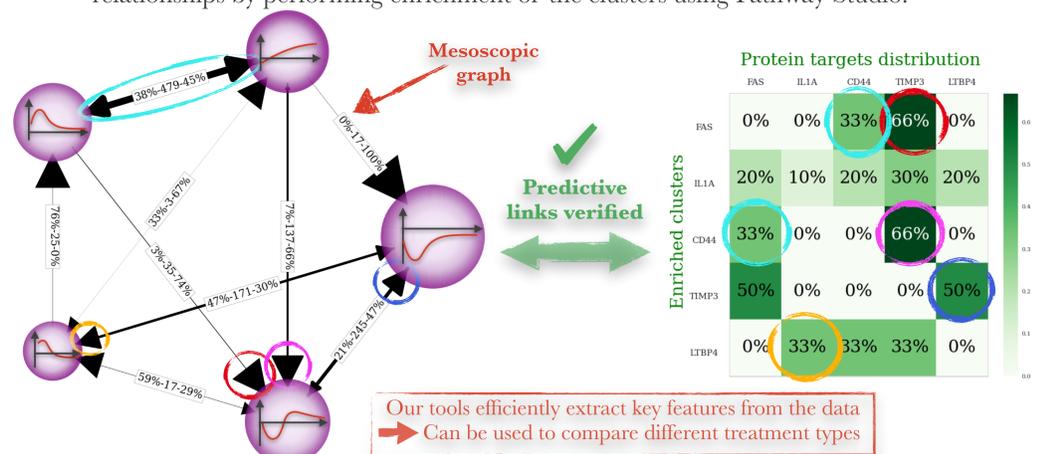
- Clustering based on the **Optimal Warping Distance** matrix:

$$\mathcal{O}W_{\mathcal{D}} = \left[\min_{s \in \mathcal{S}} \widehat{\text{diss}}(\hat{\Gamma}_i \circ \mathcal{W}_s, \hat{\Gamma}_{i'} \circ \mathcal{W}_s) \right]_{1 \leq i, i' \leq n_e}$$

Clustering & Alignment performed simultaneously with a low computational cost.

VI. Our approach vs. Biological literature

We propose a **mesoscopic** representation of the network summarising the key characteristics extracted from the data. We compare it to the known gene-protein relationships by performing enrichment of the clusters using Pathway Studio:



Our tools efficiently extract key features from the data
Can be used to compare different treatment types

References:

- Clark R. Givens, Rae Michael Shortt. (1984). A class of Wasserstein metrics for probability distributions. *Michigan Mathematical Journal*.
- David Arthur, Sergei Vassilvitskii. (2007). k-means++: the advantages of careful seeding. In Proceedings of the Eighteenth Annual ACM-SIAM Symposium on Discrete Algorithms, New Orleans, SIAM, pp. 1027-1035.
- Hae-Sang Park, Chi-Hyuck Jun. (2009). A simple and fast algorithm for k-medoids clustering. *Expert Syst. Appl.* 36. 3336-3341.
- Pathway Studio: <https://www.pathwaystudio.com>

All data were produced by the IRSN and are not public yet. Code & data soon to be available at: github.com/parsenteva.

Conflict of interest disclosure: no relevant relationships with any commercial or non-profit organisations.