# Exploration vs Exploitation of Scientific Fields
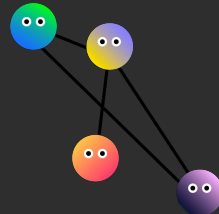
### Chakresh Kumar Singh

@chakresh_iitgn

https://chakreshiitgn.github.io/

What are the **mechanisms** underlying the knowledge discovery process?

- Rise and Fall of Research Fields
- **Exploration vs Exploitation**
- Research trajectory

### arXiv data

- First open-access **preprint repository**
- **30 years** (1986-2018)
- **175 field tags (**Physics, Maths, Computer Science, Finance / Economy, Biology)
- **1.5M articles**
- **50k authors** mapped to unique ORCID ID's .

Collected using the arXiv API

Clement, C. B., Bierbaum, M., O'Keeffe, K. P., & Alemi, A. A. (2019). **On the Use of ArXiv as a Dataset**. *arXiv preprint arXiv:1905.00075*.

High Energy Physics, Accelerator Physics, Nuclear Physics, Algebraic Geometry

Computation and Language, Vision and Pattern Recognition, Data Structures, Computation, Applications, Applied Physics

**Building the Co-Tag Network**

Article i

{ **hep-th; nlin.cd; cs.ml** }

Article j

{ **hep-th; cond-mat; acc-phy** }

**Edge Weight is defined as a function of common papers b/w two field tags**

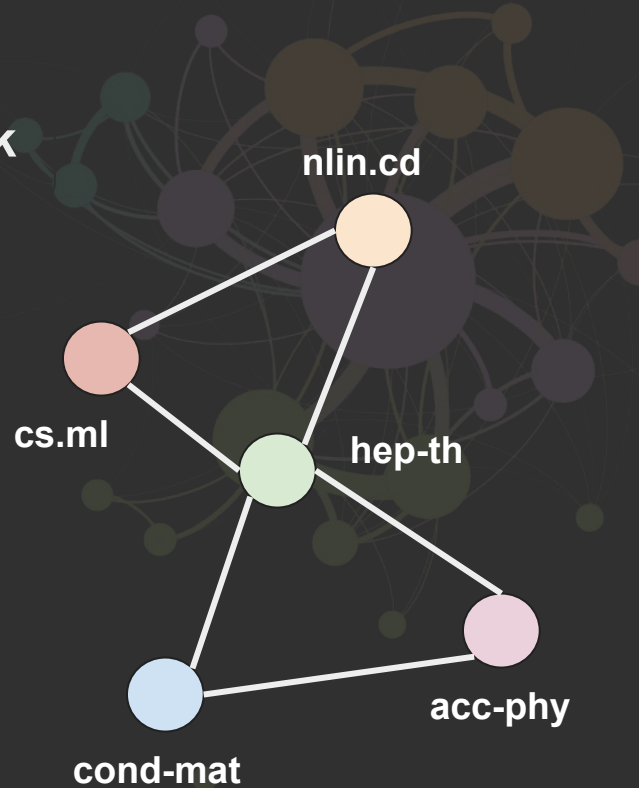**Edge Weight = $-log_{10}(p_{ij})$**

$$p_{ij} = \frac{\binom{K}{k}\binom{N-K}{n-k}}{\binom{N}{n}}$$

- N - Total Articles
- K - Articles in field i
- n - Articles in field j
- k - common articles bw i and j

**Note - Here lower p-values are more significant. We eliminate edges with p > 0.01**

- Nodes are field tags
- Size is proportional to degree
- Weighted and undirected
- Node color is the main research area
- Edge thickness proportional to weight

**Cognitive Distance**

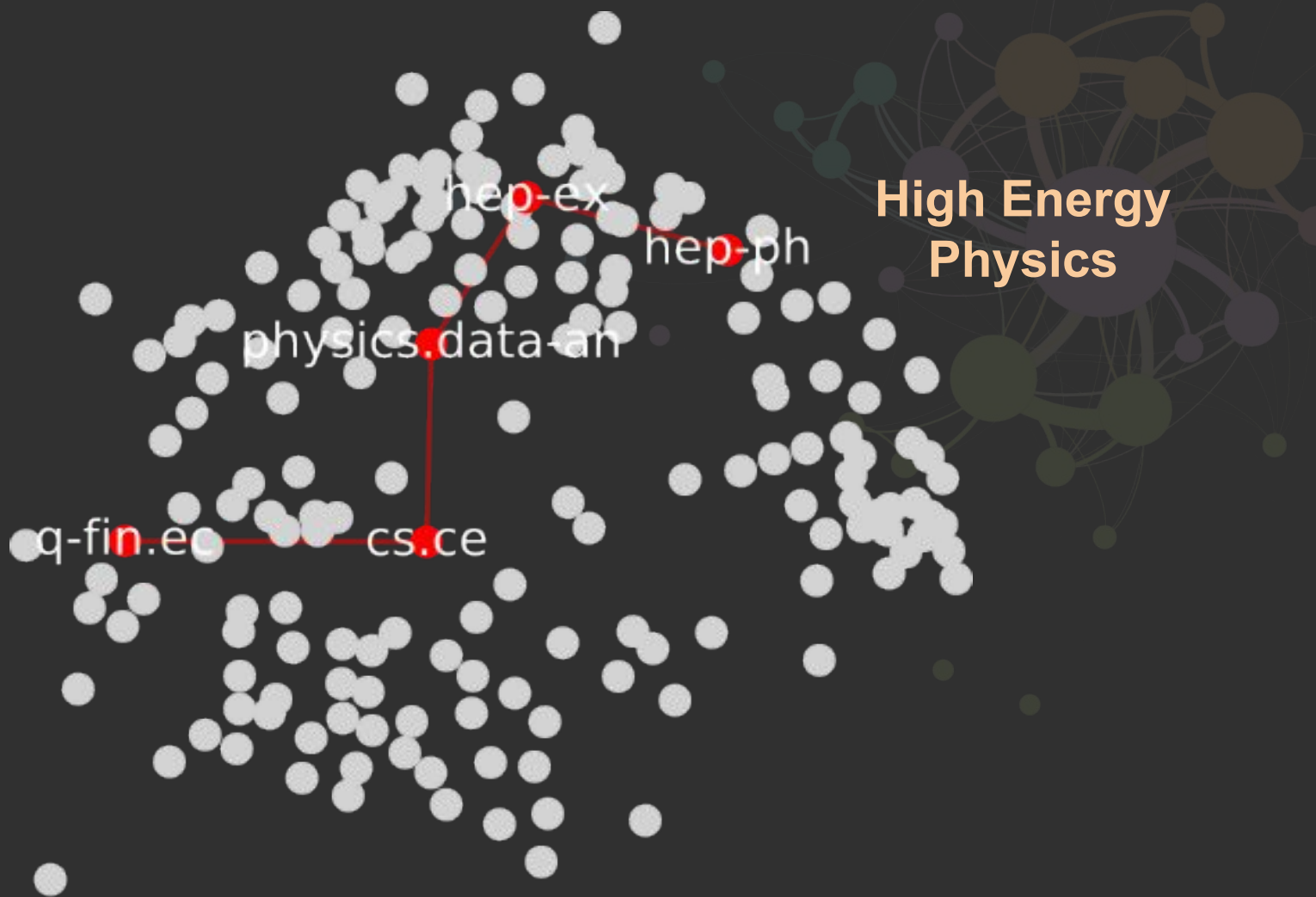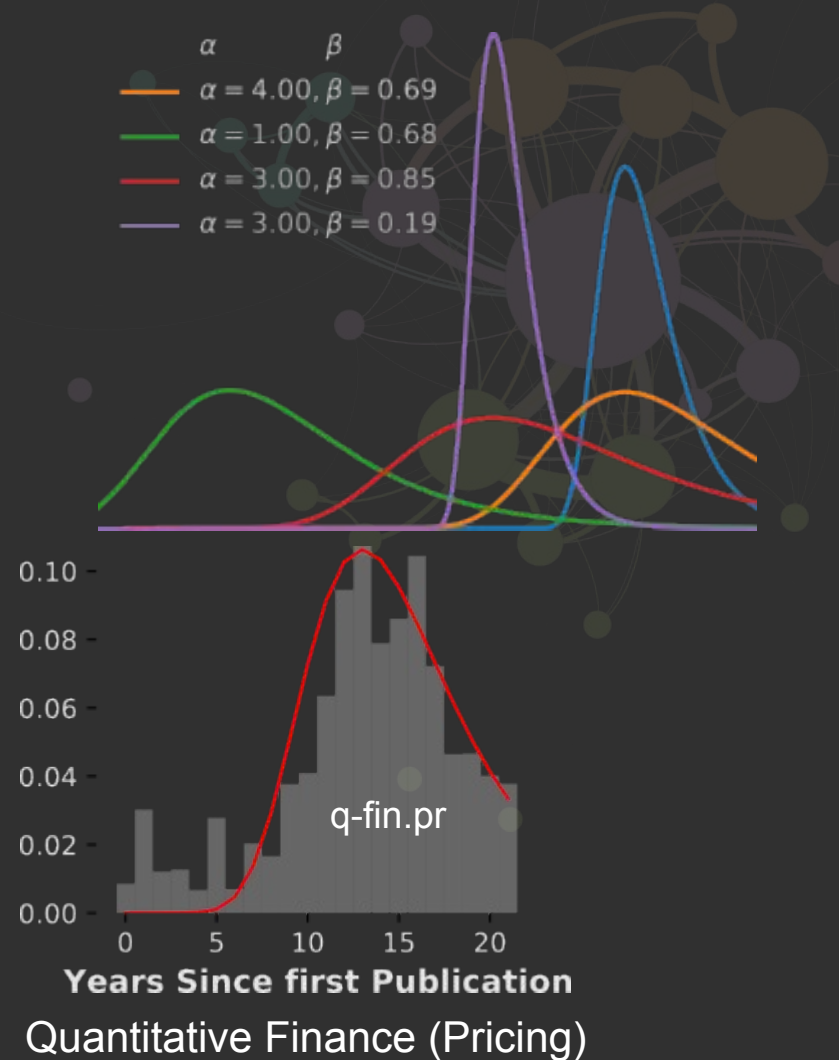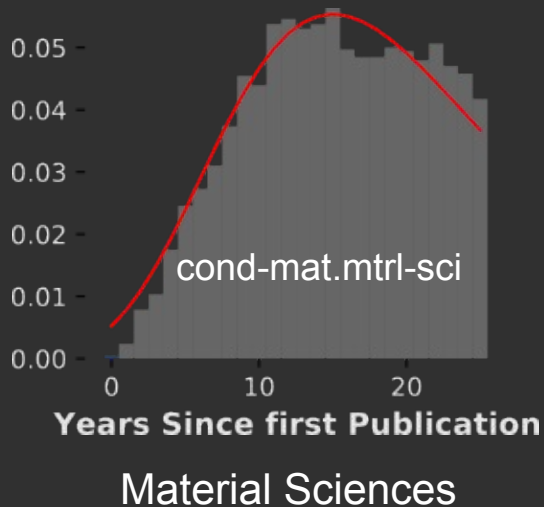$$CD_{ij} = min(\sum_e \frac{1}{W_e})$$

$e$ : edge on shortest path
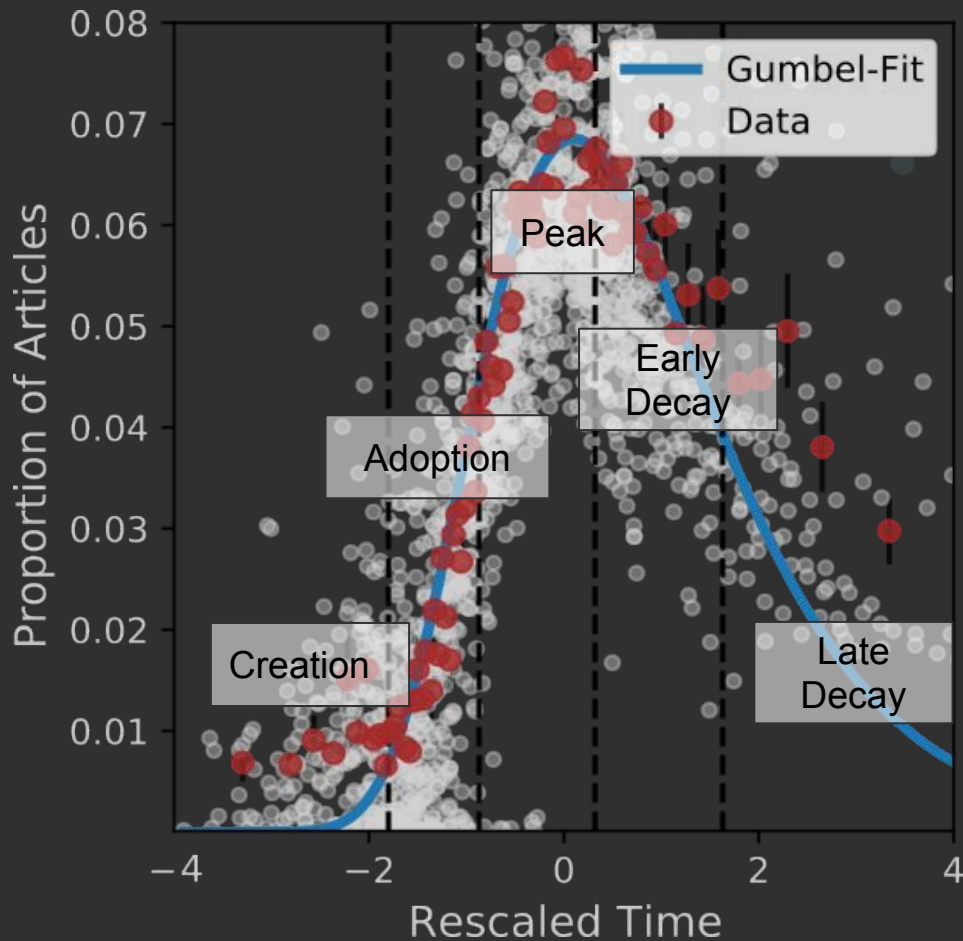
Gumbel distribution function

$$G = \frac{1}{\beta} e^{\frac{-(x-\alpha)}{\beta}} e^{-e^{\frac{-(x-\alpha)}{\beta}}}$$



| $\alpha$ | $\beta$ |
|----------|---------|
| $\alpha = 4.00, \beta = 0.69$ | |
| $\alpha = 1.00, \beta = 0.68$ | |
| $\alpha = 3.00, \beta = 0.85$ | |
| $\alpha = 3.00, \beta = 0.19$ | |

cond-mat.mtrl-sci

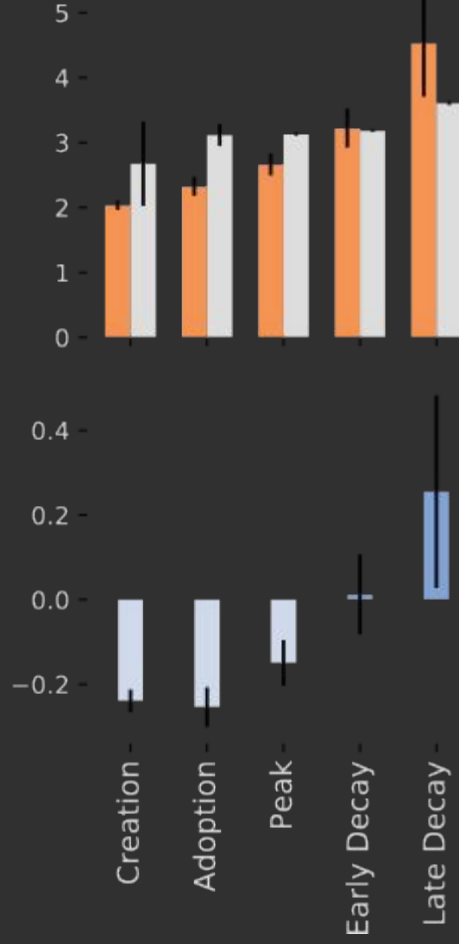Material Sciences

q-fin.pr

Quantitative Finance (Pricing)

Rescaled time

$$t' = \frac{t - \alpha}{\beta}$$

Field stages are defined at 2.5%, 16%,50% and 84% of the fit curve (blue). These number are borrowed from the **diffusion of innovation** literature
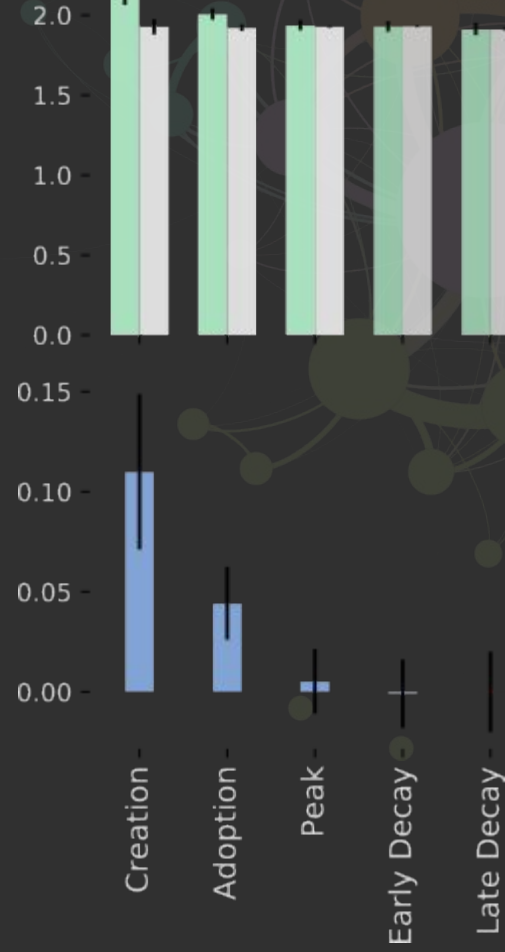
**7x more enrichment with network based measure**

Researchers in creation phase connect distant fields together (**explore**) whereas in later phases they focus on closely related fields (**exploit**)

- Scientific fields exhibit a universal rise and fall process allowing to define standardised stages of development

- Early stages are enriched with small teams of interdisciplinary authors

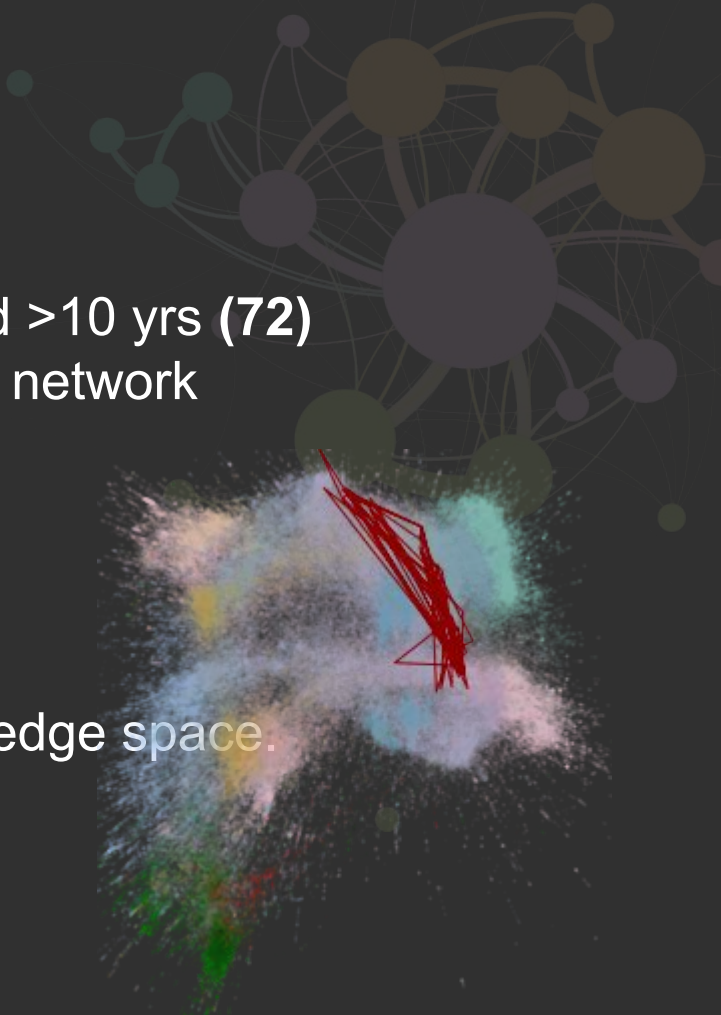- **Network based cognitive distance is a strong marker of early innovation (exploration)**

**Singh, Chakresh**, Emma Barme, Robert Ward, Liubov Tupikina, and Marc Santolini. "**Quantifying the rise and fall of scientific fields.**" *arXiv preprint arXiv:2107.03749* (2021). (In Review)
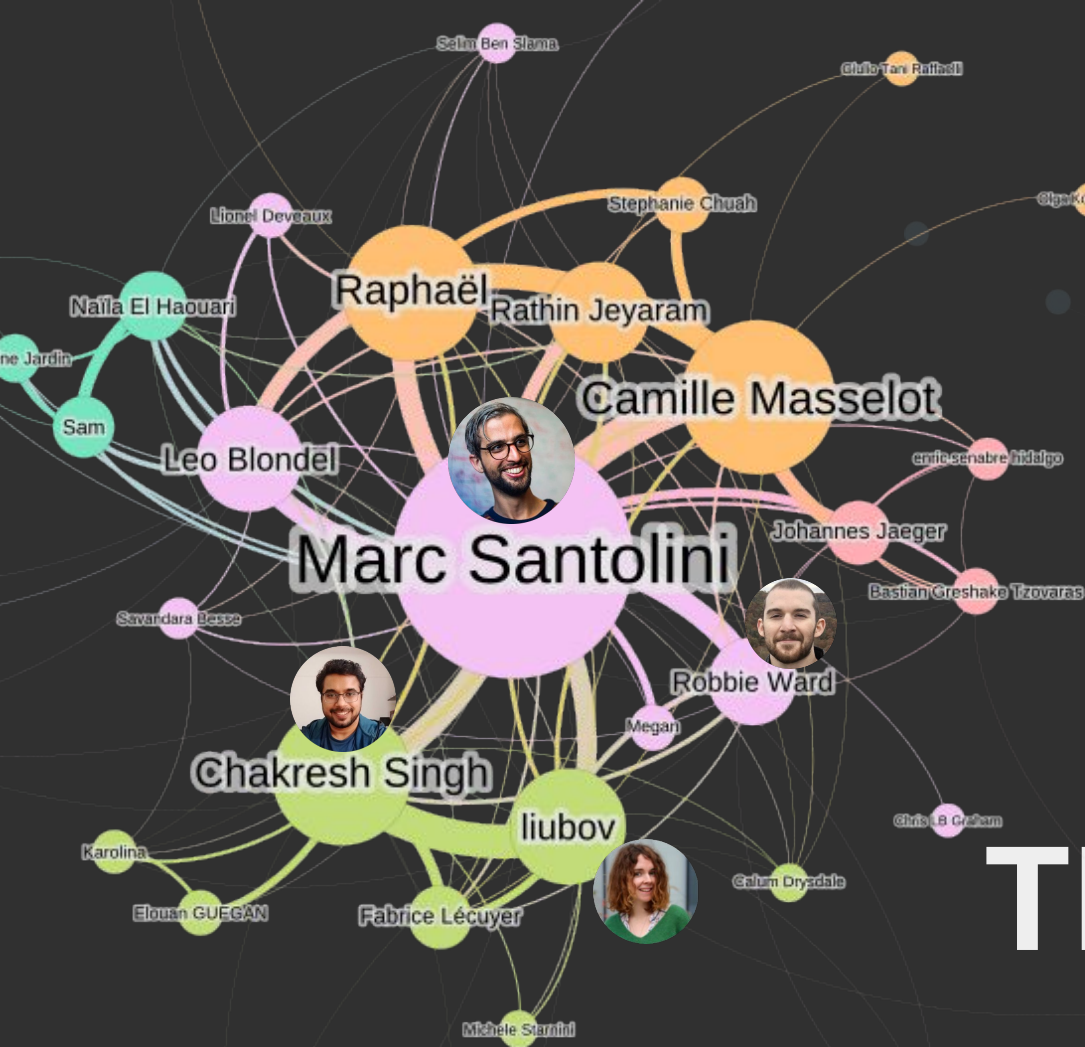
# Limitations -

- We focus on fields that were unimodal and >10 yrs **(72)**
- We assume the co-tag network as a static network

# Perspectives -

- **Knowledge foraging**:
    Researchers' trajectories in the knowledge space.
- **Model** the knowledge discovery process
- **Large Scale** Data-sets

https://interactiondatalab.com/

@InteractionData

THANK YOU

# Supplementary Slides