

CareerMap: Visualizing Career Trajectory

Kan Wu, Jie Tang, Zhou Shao, Xinyi Xu, Bo Gao & Shu Zhao

Dept. Of Computer Science, Tsinghua University

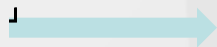


Challenges

Challenge 1:

Name ambiguity

solution

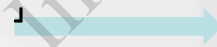


Unified Probabilistic Models [1]

Challenge 2:

Data incompleteness

solution

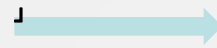


Spatial-Temporal Factor Graph Model (STFGM)

Challenge 3:

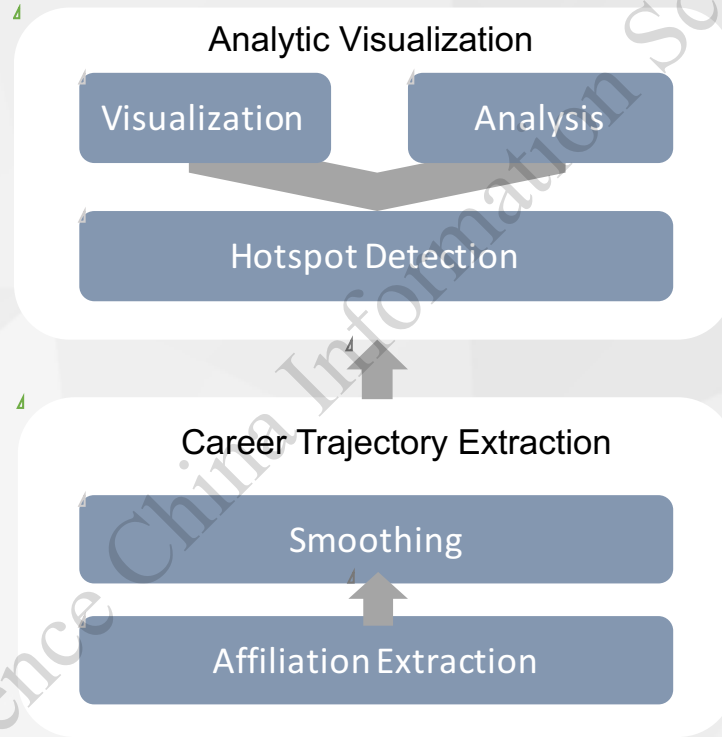
Visualize many scholars' merged trajectories on the map, e.g. 100 people move from Boston to New York

solution



Hotspot detection algorithm

Architecture

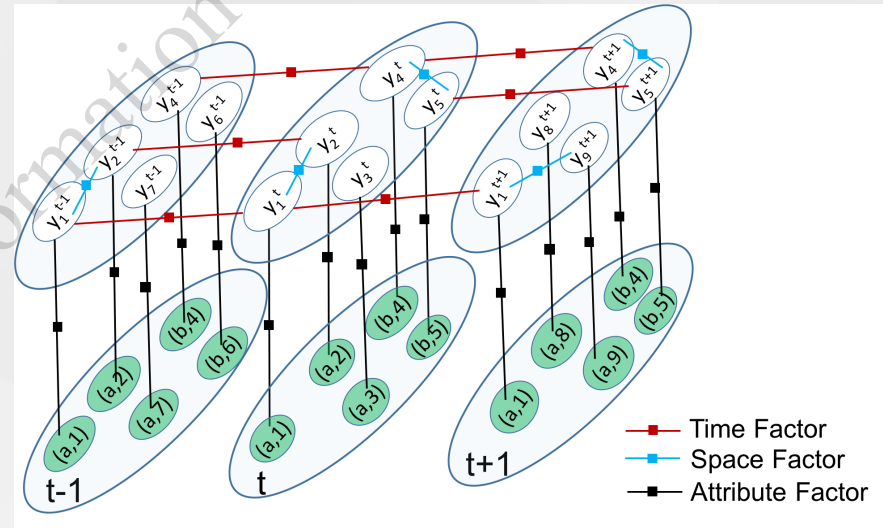


Spatial-Temporal Factor Graph Model (continued)

• The general idea

- try to find the affiliation-known coauthor who has the same affiliation as the target author with missing affiliation.

Each green point with common t outside, representing a tuple of $\langle \text{Time } t, \text{ Author } a_{i1}, \text{ Author } a_{i2} \rangle$, is an observation instance where a_{i1} is the target author and a_{i2} is a coauthor with known affiliation at t . Associated with each observation instance is a hidden binary-valued variable representing the affiliation similarity between the two authors. If they belong to the same affiliation at that time, the hidden value is 1, otherwise 0.



Spatial-Temporal Factor Graph Model (continued)

• Attribute factor

- captures the features of each tuple
<Time t , Author a_{i1} , Author a_{i2} >

$$f(\mathbf{x}_i^t, \mathbf{y}_i^t) \triangleq \frac{1}{Z_\omega} \exp \{ \omega^T \Phi(\mathbf{x}_i^t, \mathbf{y}_i^t) \} \quad (1)$$

• Space factor

- captures the correlation between the hidden variables in the same time
- \mathcal{N}_S denotes all the space relations

$$\mathcal{S}(y_i^t, \mathcal{N}_S(y_i^t)) \triangleq \frac{1}{Z_\beta} \exp \left\{ \sum_{y_j^t \in \mathcal{N}_S(y_i^t)} \beta^T \Psi(y_i^t, y_j^t) \right\} \quad (2)$$

• Time factor

- captures the correlation between the hidden variables in the same time
- \mathcal{N}_T denotes all the time relations

$$\mathcal{T}(y_i^t, \mathcal{N}_T(y_i^t)) \triangleq \frac{1}{Z_\gamma} \exp \left\{ \sum_{y_i^{t'} \in \mathcal{N}_T(y_i^t)} \gamma^T \Omega(y_i^t, y_i^{t'}) \right\} \quad (3)$$

Spatial-Temporal Factor Graph Model

• Model Learning

- Maximize the likelihood of the observed data
- $\theta \triangleq (\omega^T, \beta^T, \gamma^T)^T$ is the parameters to be learned of the model

$$\begin{aligned} P(Y|X, \theta) &= \prod_t \prod_i f(\mathbf{x}_i^t, y_i^t) \mathcal{S}(y_i^t, \mathcal{N}_S(y_i^t)) \mathcal{T}(y_i^t, \mathcal{N}_T(y_i^t)) \\ &= \frac{1}{Z_\omega Z_\beta Z_\gamma} \exp \left\{ (\omega^T, \beta^T, \gamma^T) \sum_t \sum_i g(y_i^t) \right\} \\ &= \frac{1}{Z_\theta} \exp \left\{ \theta^T \mathcal{G}(Y) \right\} \end{aligned} \tag{4}$$

$$\theta^* = \arg \max_{\theta} \mathcal{O}(\theta) = \arg \max_{\theta} \log P(Y^L | X, \theta) \tag{5}$$

Smoothing

- The general idea

- Use weight to reflect confidence of an affiliation at a time.
- Leverage the number of papers with the affiliation at time t as the weight.
- Denoting the weights at t_1 and t_2 are w_1 and w_2 respectively, the weight center t_c can be computed from:

$$\frac{t_c - t_1}{t_2 - t_c} = \frac{w_1}{w_2}$$

- If information between t_1 and t_2 is missing,
- $\forall t (t_1 < t < t_c)$, Affiliation(a, t) = Affiliation(a, t_1)
- $\forall t (t_c < t < t_2)$, Affiliation(a, t) = Affiliation(a, t_2)

Example of scholar career trajectory extraction

Scholar Career Trajectory

name:

affiliation:

[Aminer Search](#)

[Heatmap](#)

- 2015 - 2015 : w3c and mit boston ma usa
- 2013 - 2014 : computer science and ai laboratory massachusetts institute of technology vassar street cambridge ma 02139 usa
- 2012 - 2012 : university of southampton
- 2007 - 2011 : w3c and mit boston ma usa
- 2005 - 2006 : computer science and ai laboratory massachusetts institute of technology vassar street cambridge ma 02139 usa
- 2002 - 2004 : w3c and mit boston ma usa
- 1997 - 2001 : massachusetts institute of technology's laboratory
- 1994 - 1996 : cern geneva switzerland
- 1992 - 1993 : cern
- 1990 - 1991 : w3c and mit boston ma usa
- 1988 - 1989 : cern geneva switzerland



Hotspot detection algorithm

- The general idea

- The heat centers have more neighbors than surrounding points.
- The heat centers "absorb" their surrounding points as their neighbors. If a point is "absorbed" by a heat center, then its neighbors are emptied.
- Finally, the points left out with nonempty neighbors are heat centers.

Algorithm 1 Hotspot detection

Input : set of points $\{v_1, v_2, \dots, v_N\}$, radius of heatCenters R

Output: heatCenters and points in them

for $i \leftarrow 1$ to N **do**

 | neighbors(v_i) $\leftarrow \{v_i\}$

end

for $i \leftarrow 1$ to $N - 1$ **do**

 | **for** $j \leftarrow i + 1$ to N **do**

 | **if** distance(v_i, v_j) $< 2R$ **then**

 | neighbors(v_i) \leftarrow neighbors(v_i) $\cup \{v_j\}$

 | neighbors(v_j) \leftarrow neighbors(v_j) $\cup \{v_i\}$

 | **end**

 | **end**

end

$\{v'_1, v'_2, \dots, v'_N\} \leftarrow$ sort $\{v_1, v_2, \dots, v_N\}$ according to neighbor size in desc order

for $i \leftarrow 1$ to $N - 1$ **do**

 | **for** $v_j \in$ neighbors(v'_i) **do**

 | neighbors(v_j) $\leftarrow \emptyset$

 | **end**

end

for $i \leftarrow 1$ to N **do**

 | **if** neighbors(v'_i) $\neq \emptyset$ **then**

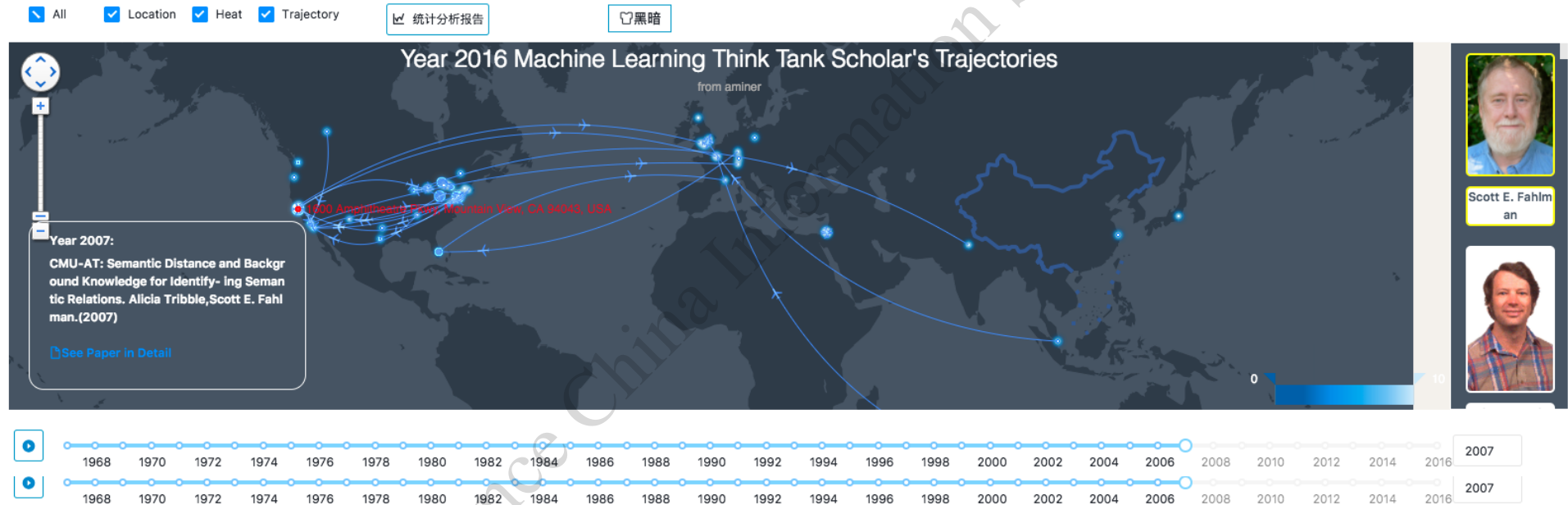
 | heatCenters(v'_i) = neighbors(v'_i)

 | **end**

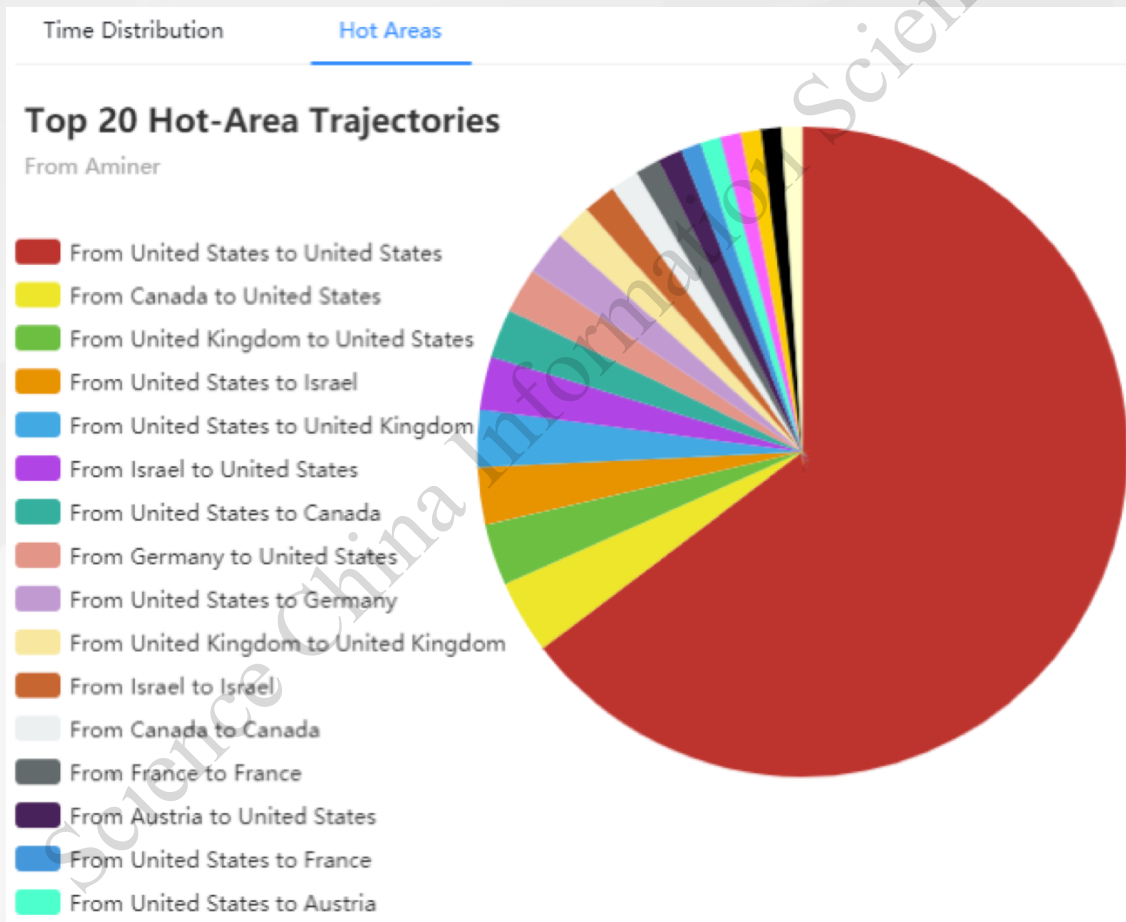
end

Trajectory map generated by Career Map

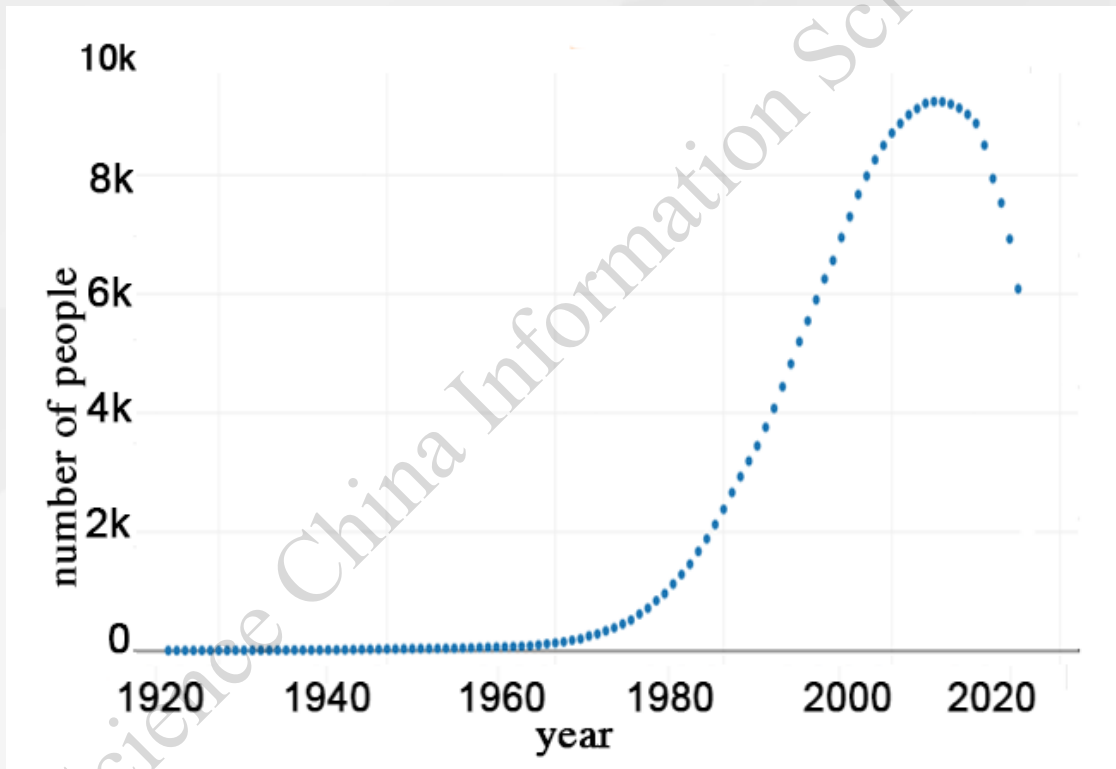
Scholars' Trajectories and Heat Map



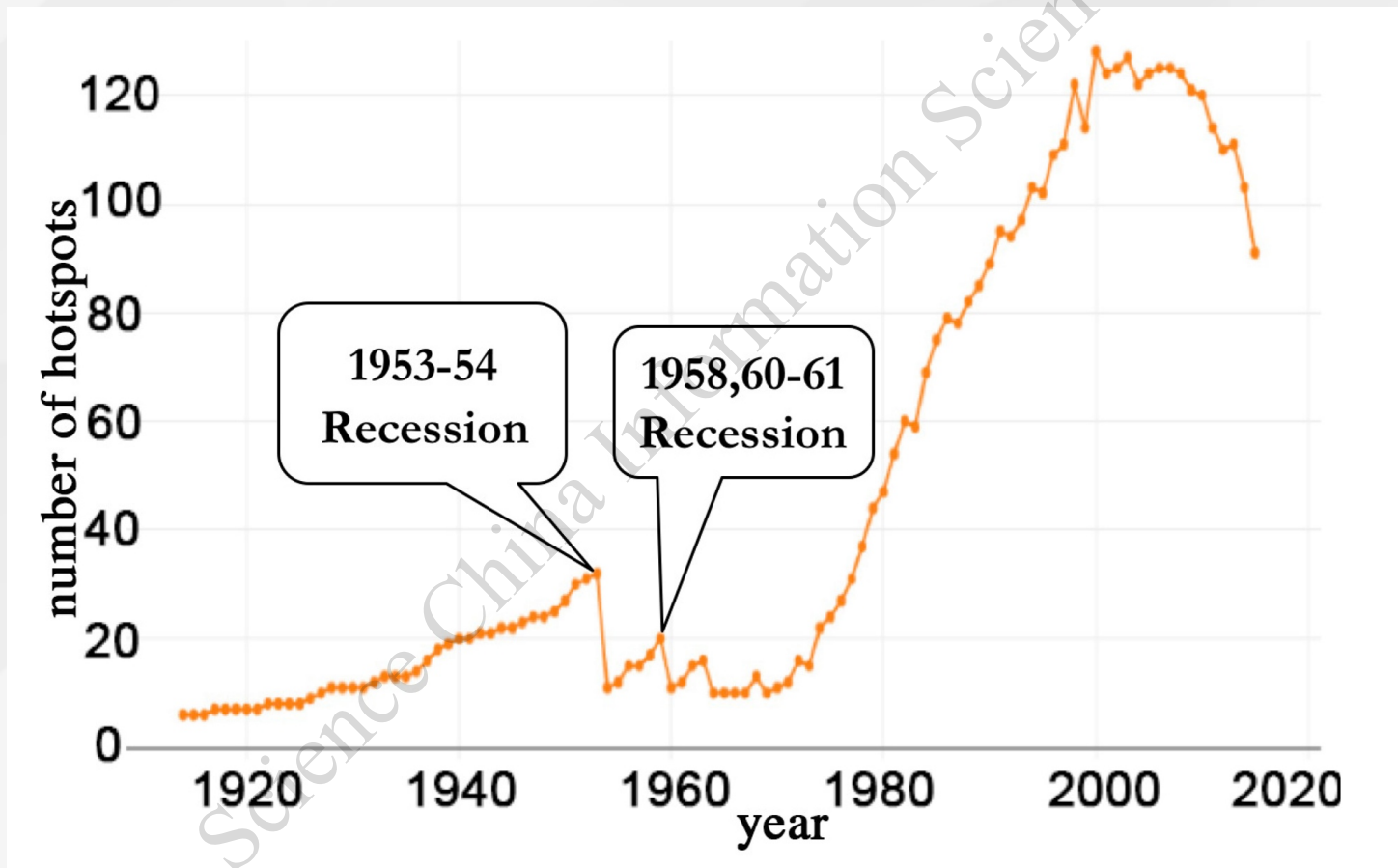
Analytic Visualization



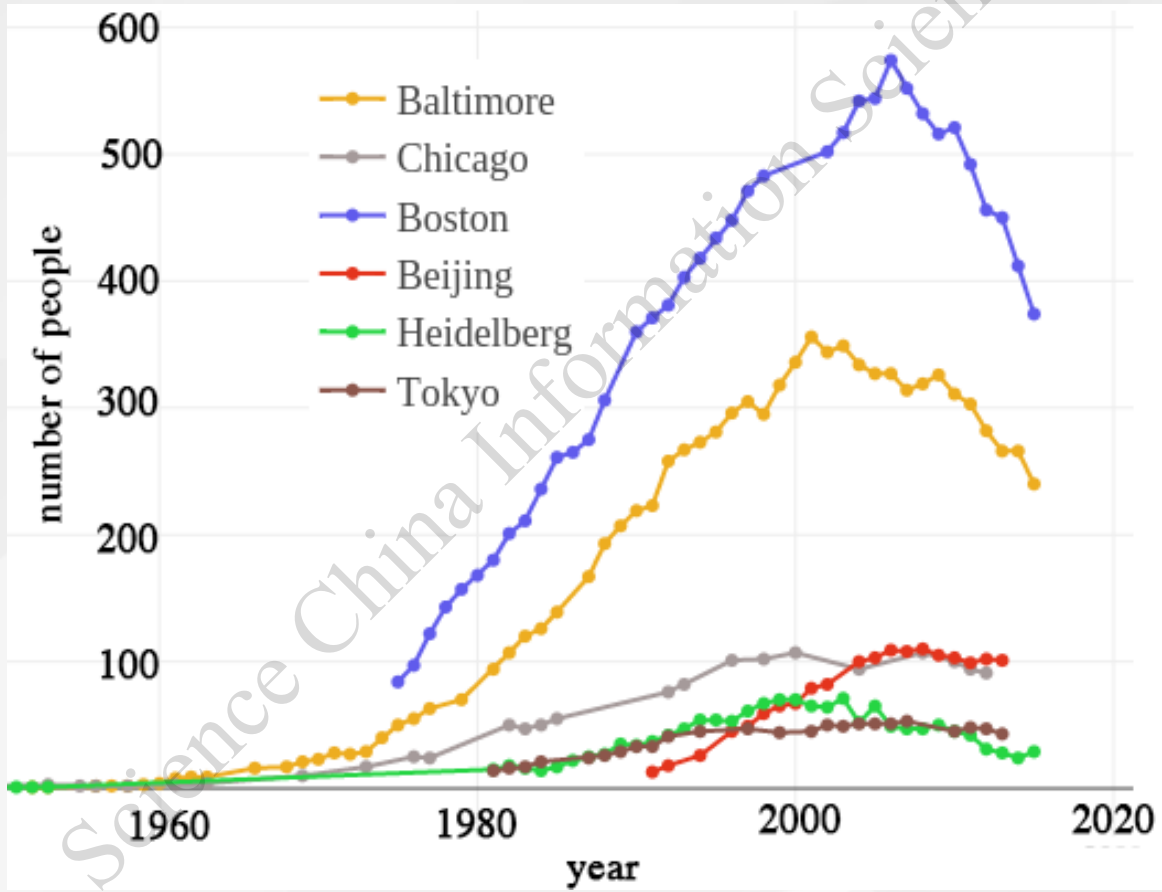
Some Interesting Case Study (continued)



Some Interesting Case Study (continued)



Some Interesting Case Study



Summary

Some interesting
case studies

3

2

Architecture, technologies
and main features of the
system

1

We introduce the
challenges of building
CareerMap, a system for
visualizing scholars'
career trajectory

Science China Information Sciences

Thanks
Q&A

Science China Information Sciences