# The interrelationship of knowledge structure across language groups in communal data sets

Jisung Yoon[1,2], Jinhyuk Yun[3] and Woo-Sung Jung[1,4,5]

1. Dept. of Industrial and Management Engineering, POSTECH, Pohang, Korea
2. School of Informatics, Computing, and Engineering, Indiana University, Bloomington, IN, USA
3. Department of Scientometric Research, Korea Institute of Science and Technology Information(KISTI), Korea
-> Future Technology Analysis Center
4. Dept. of Physics, POSTECH, Pohang, Korea
5. Asia Pacific Center for Theoretical Physics(APCTP), Korea

Conference on Complex System 2019
Oct 1st, 2019

# Table of Contents

**Introduction: Knowledge Structure**

- Research background
- Research question

**Data and methods**

- Wikipedia data: communal data set
- Calculate similarity between knowledge structure

**Results**

- Community detection results
- Factor analysis

**Summary**

1. 발표자료가 전체적으로 밝혀줌
   → 발표자료로 어느정도 이해될 수 있도록
2. 관사. 단복수 한번 더 체크 …

**Knowledge** is a familiarity, awareness, or understanding
of someone or something

which is acquired through **experience** or **education**
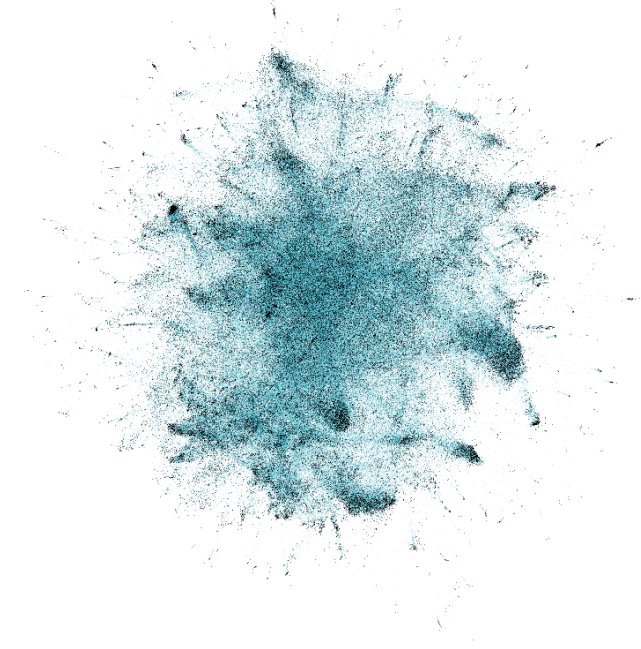by perceiving, discovering, or learning.

-Wikipedia

**Human understanding** is root of the general laws
of nature that organize all experience

-Immanuel Kant

# Knowledge Structure

- Knowledge structure can be varied by **personality**, **living country** or **linguistic profile** based on the social structure and education system
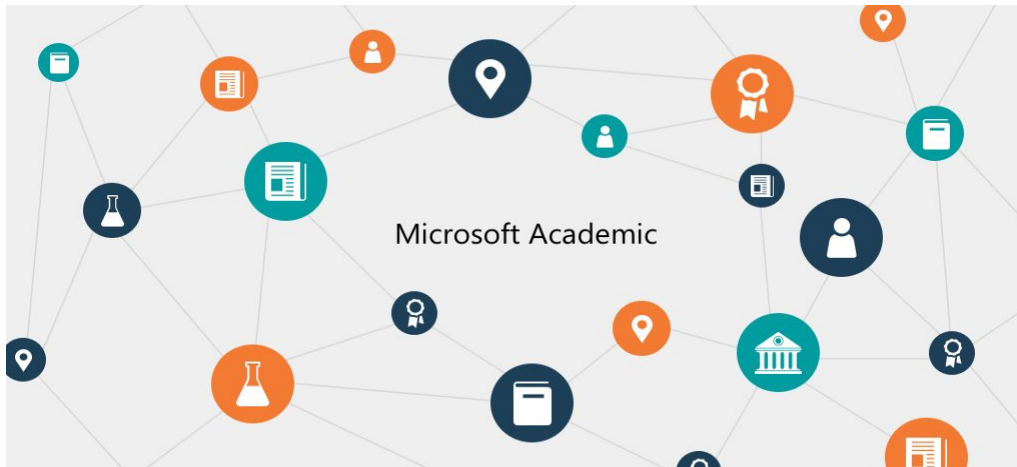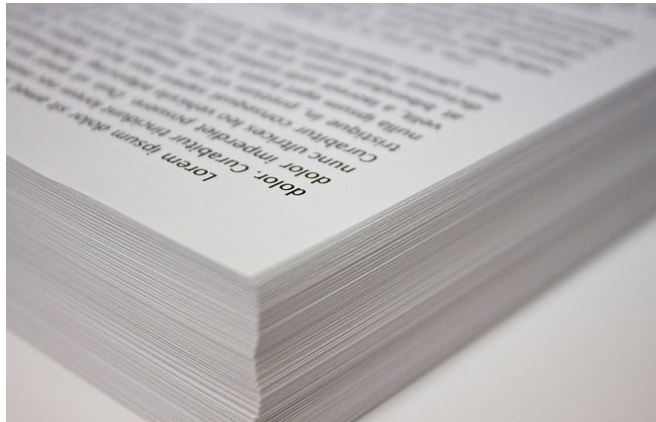
# Language Affect Knowledge

- Language Socialization (Schieffelin, B. B., & Ochs, E.,1986)
  - Socialization(acquiring knowledge) through the use of language
  - Schieffelin, B. B., & Ochs, E. (1986). Language socialization. Annual review of anthropology, 15(1), 163-191.

- Language and Knowledge (Code, L., 1980)
  - Language and knowledge are mutually influential
  - Perception and knowledge are organized by language from the flux of sensory experience.

# Research Questions

- What are major factors influencing the similarity of knowledge structure across the language group?

  1. How can we construct a knowledge structure of a language group?

  2. How can we compute similarity among the obtained knowledge structure? *derive*

  3. What are major factors influencing the similarity?

# Knowledge Database





Microsoft Academic



Not proper to construct knowledge structure of specific language group

# Communal data set - Wikipedia

- Internet encyclopedia is edited by users that use specific language
  - Result of a collective intelligence

- 294 active language editions (April, 2019)

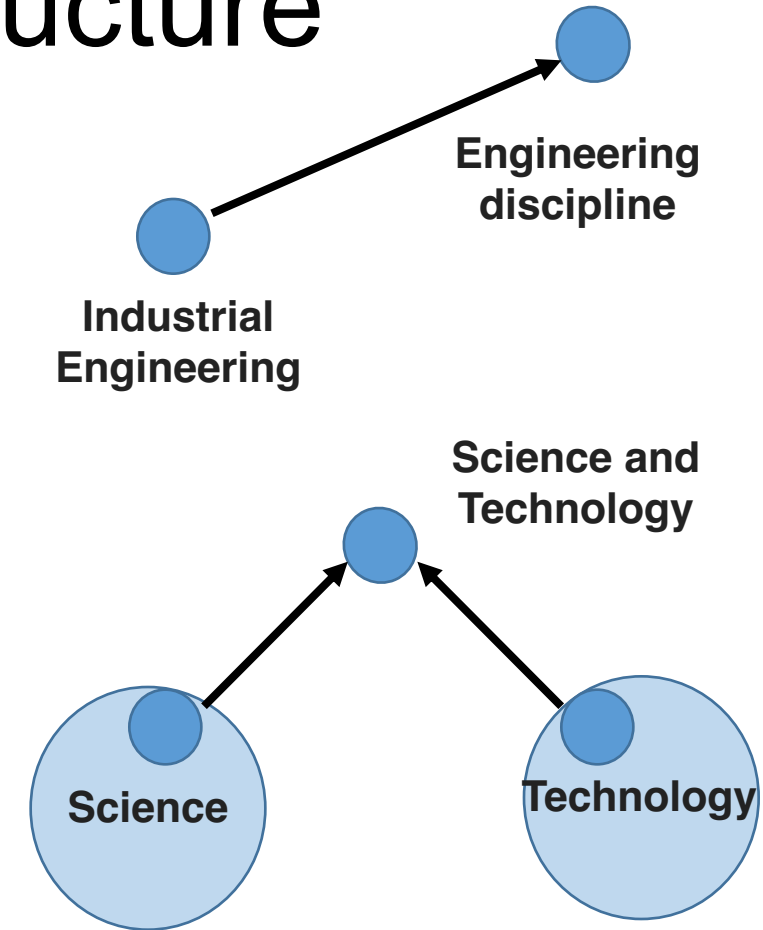- Possible to get **knowledge structure** of a **specific language group**

# Data

- Dump data of 59 different language editions of Wikipedia on August 20, 2018
  - Category-link data set: relation between a category and other items
    - For constructing a knowledge network

  - Language-link data set: bridge data between another language editions of Wikipedia items that same meaning
    - For comparing knowledge structure of different language edition

# Construct Knowledge Structure

- Knowledge Network
  - **One for each language**
  - **Node**: each category or page
  - **Link**: directed If node A refer node B, node A -> node B

  - Sub-network with artificial root assigned as a common parent node "**Science**" and "**Technology**"
    - To get fine-grained form of knowledge network
    - Science covers all branch of science
      - Applied sciences, Formal sciences, Natural Sciences, Social sciences



**Engineering discipline**

**Industrial Engineering**

**Science and Technology**

**Science**

**Technology**

## Complex system

From Wikipedia, the free encyclopedia

*"Complex systems" redirects here. For the journal, see Complex Systems (journal).*

Categories: Complex dynamics | Complex systems theory | Cybernetics | Emergence | Systems | Systems science | Mathematical modeling

# Construct Knowledge Structure

## Complex system

From Wikipedia, the free encyclopedia

*"Complex systems" redirects here. For the journal, see Complex Systems (journal).*

Categories: Complex dynamics | Complex systems theory | Cybernetics | Emergence | Systems | Systems science | Mathematical modeling

## 복잡계

위키백과, 우리 모두의 백과사전.

Mechanics

Statistical Physics ~~Physics~~     System Science

분류: 복잡계 이론 | 통계역학 | 시스템 | 시스템 과학

Complex system Theory     System

## Similar, but slightly different
## Then, how can we calculate similarity between knowledge structure?

# Calculate knowledge structure similarity

*→ mapping(?) / pairing(?)*
*→ translate는 조금 안 맞는듯...*

- Calculate subject similarity first.
  - **Characterize** with genealogy vector, and **Translate** to target language with **language-link data set,** and **compare!**
    - **[Characterize]** Genealogy vector of a given node as a Personalized Page Rank (Jeh, G., & Widom, 2003) of a subject in network.
    - **[Translate]**  Matching with language link data set
      - E.g.) Republic of Korea (en) => 한국 (ko)
    - **[Compare]**  Calculate between translated genealogy vector and target genealogy vector
      - We use 1  - Euclidian distance as similarity

- Then, knowledge structure similarity is average value of all subject similarity

*→ We define        the Similarity between two knowledge structure*
*as the average value of all subject similarity*
*between the languages.*

# Calculate knowledge structure similarity

- For example, *Complex System* and 복잡계 (English to Korean)

  - **Characterize**

$$X_{\text{Complex System}} = [\frac{1}{5}, \frac{1}{5}, \frac{1}{5}, 0, \frac{1}{5}, \frac{1}{5}, \dots] \in \mathrm{R}^{\mathrm{N_E}}$$

$$Y_{복잡계} = \left[\frac{1}{4}, \frac{1}{4}, 0, \frac{1}{4}, \frac{1}{4}, \dots\right] \in R^{N_k}$$

  - **Translate**

$$Y_{Complex\ System} = \left[\frac{1}{5}, \frac{1}{5}, 0, 0, \frac{1}{5}, \frac{1}{5}, \dots\right] \in \mathrm{R}^{\mathrm{N_K}}$$

  - **Compare**

$$S_{Complex\ System\ -복잡계} = 1 - d\ (Y_{Complex\ System}, Y_{복잡계})$$

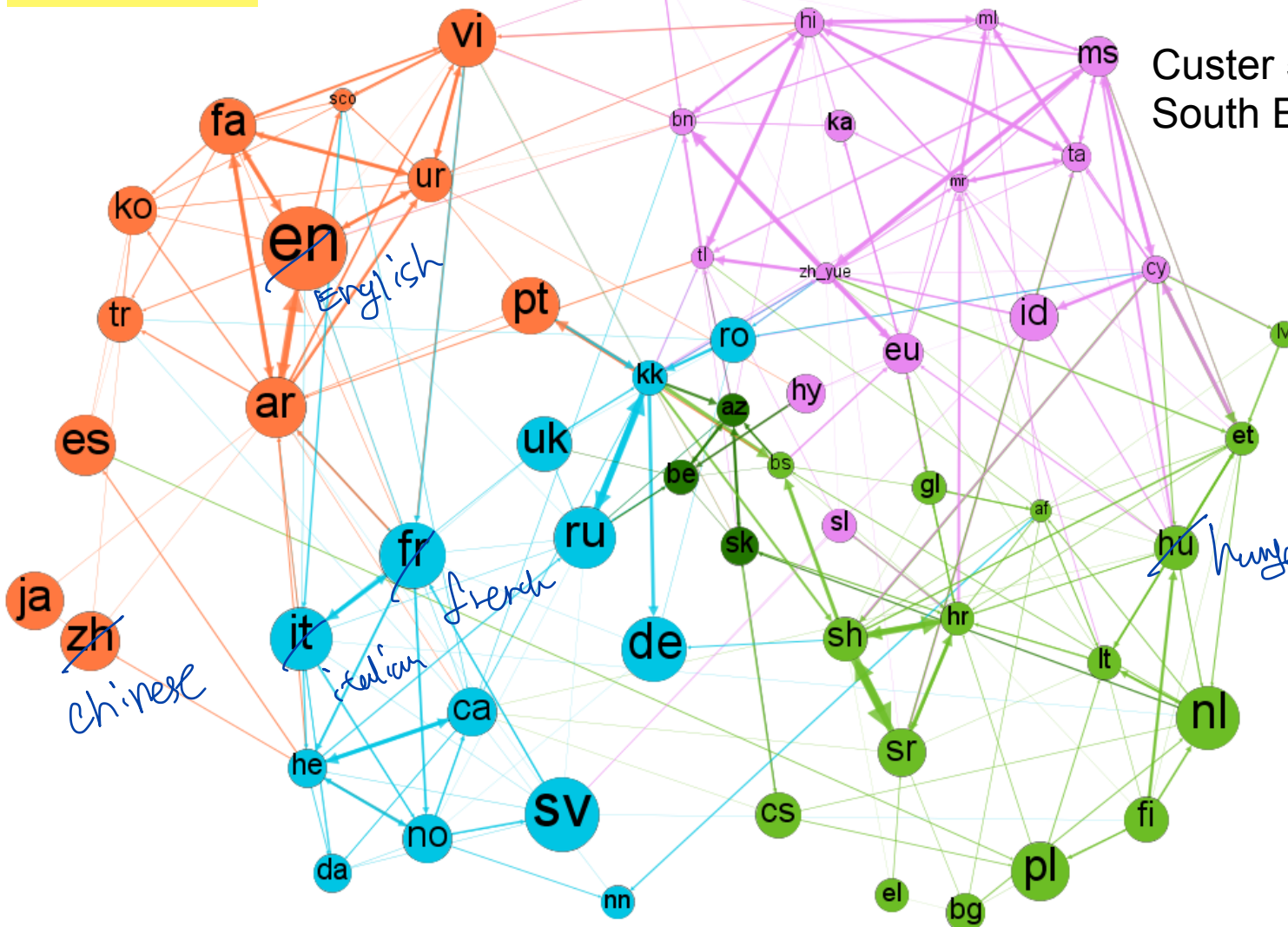- ==Knowledge structure similarity== English to Korean can calculate with averaging over all the subject.

↳ Structural similarity of knowledge from

Custer 1:
Transnational Cluster

Custer 5:
South East Asia Cluster

Custer 4:
Northeastern Europe Cluster

$$r_{ij} = Strength_{total} * \frac{s_{ij}}{strength_{out} \; of \; i \; * \; strngth_{in} \; of \; j}$$
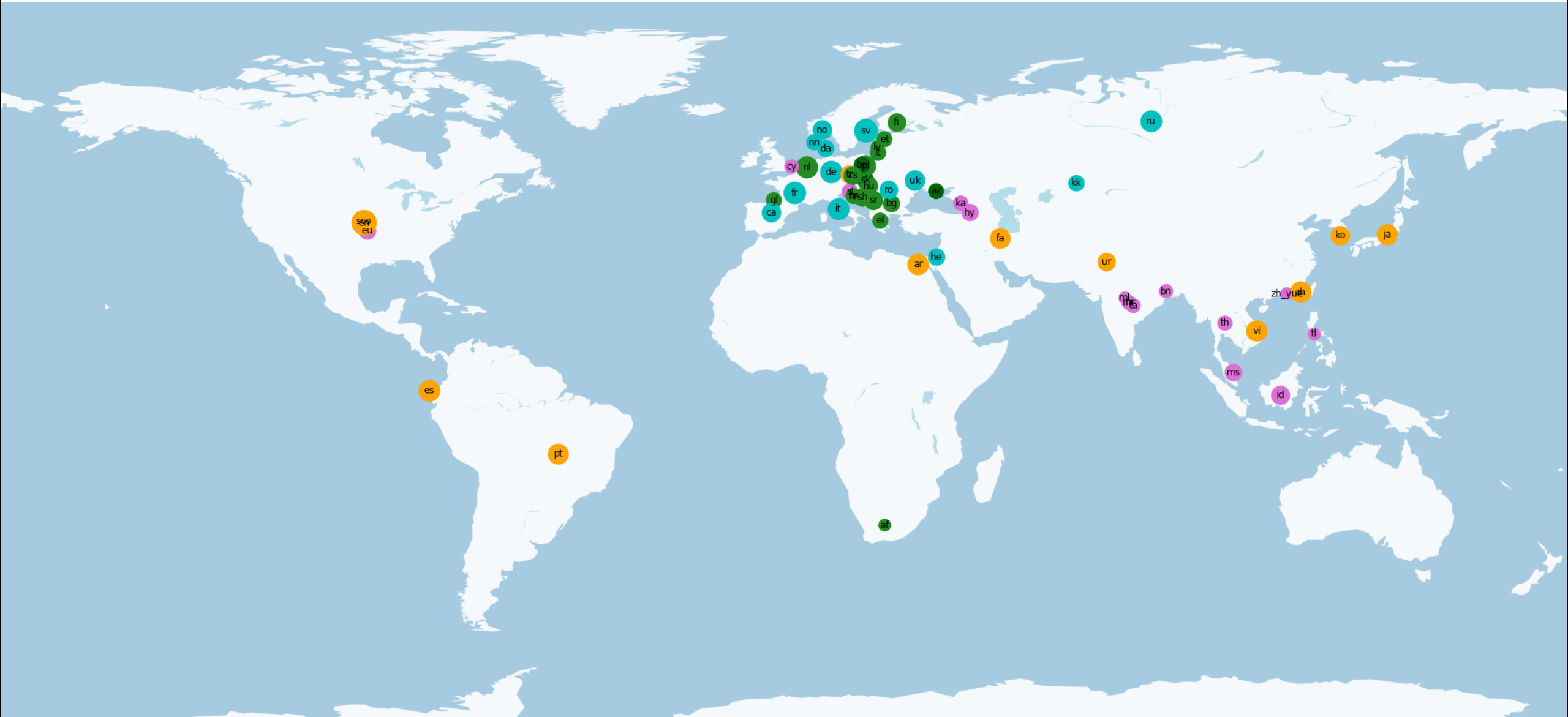
Custer 2:
Western Europe Cluster

Custer 3:
Eastern Europe Cluster

**Threshold = 1.04**
**Resolution parameter = 1**
**Modularity = 0.42**

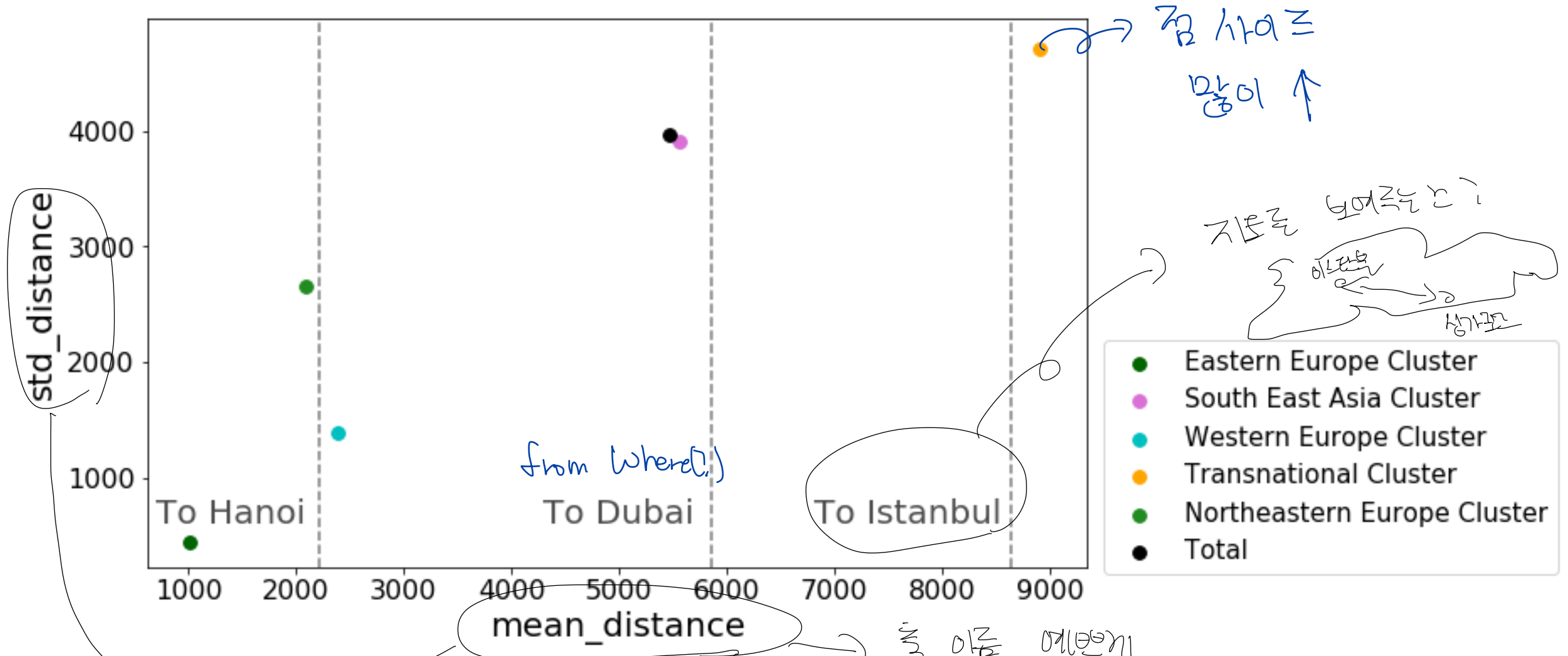# Community result on map

Custer 1:
Transnational Cluster

Custer 2:
Western Europe Cluster

Custer 3:
Eastern Europe Cluster

Custer 4:
Northeastern Europe Cluster

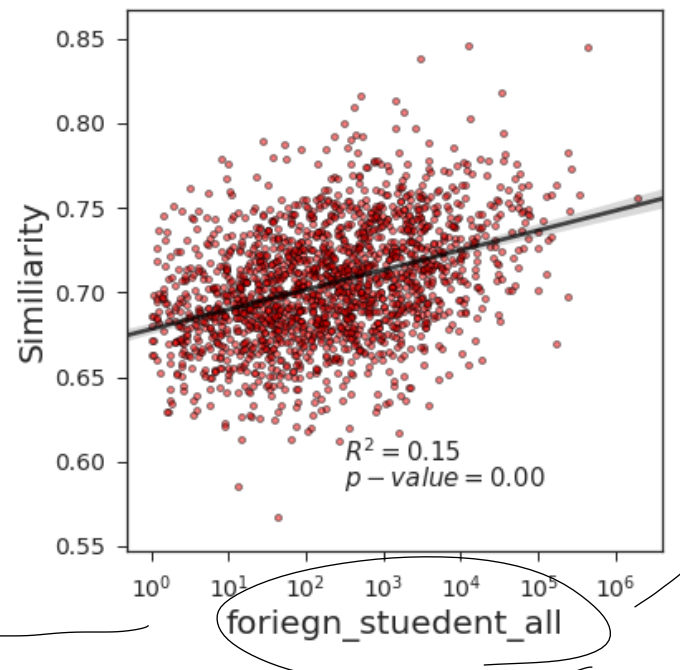Custer 5:
South East Asia Cluster
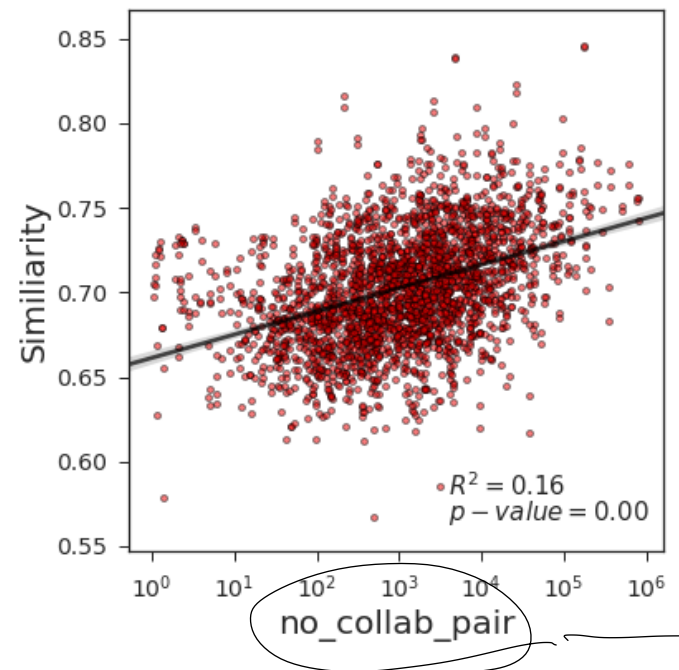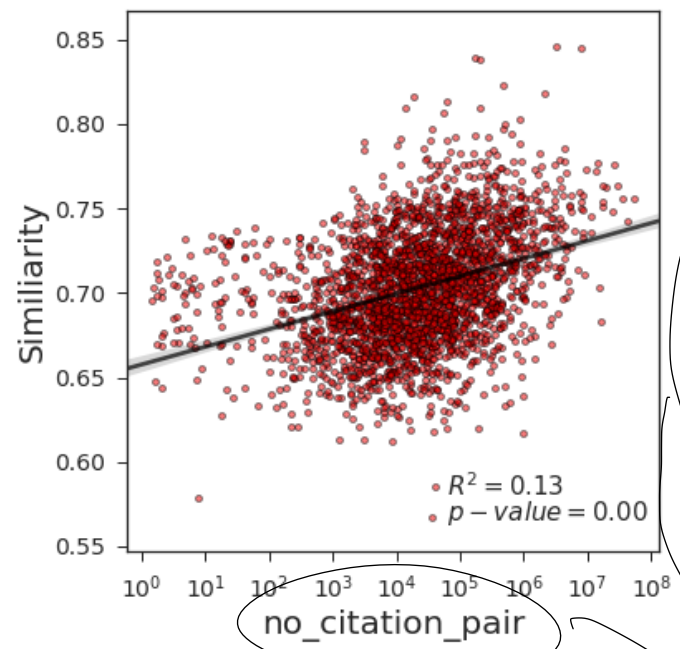
# Community result on map



Distance matters, but not significant

# Factor Analysis

- 4 Factors to analysis
    1. **Physical Distance** - distance between language
    2. **Weak knowledge transfer**\* – number of citation (paper)
    3. **Strong knowledge transfer**\* – number of collaboration (paper)
    4. **Soft Power Mobility**\* – number of foreign students

↳ Movement of the soft power.

\* These data are from SCOPUS and OECD. Basically, there are county to country data. We projected to language to language dimension with country to language data set (Ronen et al, 2014).

**Negative**
    Physical Distance

**Positive**
    Weak knowledge Transfer
    Strong Knowledge Transfer
    Soft Power Movement

이건 당연한 뭐뭐

이건데 3가는 associativeness이나

힘이줄 표기기

# Summary

- We use Wikipedia dump data of 59 different language editions and construct the Knowledge Network to compare the knowledge structure between languages.

- We found 5 geo-locational clusters, but physical distance is not a significant for some clusters.

- We conduct factor analysis to identify knowledge structure similarity between languages, and find a pattern ~~a pattern~~ correlations for :
    1. **Physical Distance** - distance between language
    2. **Weak knowledge transfer\*** – number of citation (paper)
    3. **Soft Power Mobility\*** – number of foreign students
    4. **Strong knowledge transfer\*** – number of collaboration (paper)

- It helps to understand fact that knowledge structure has been affected by language groups.

⑤ 나라 여기 표든 이상해 준이나요...

To do에 대해 언급하시고 ...

# References

- Schieffelin, B. B., & Ochs, E. (1986). Language socialization. Annual review of anthropology, 15(1), 163-191.

- Code, L. (1980). Language and knowledge. word, 31(3), 245-258.

- Jeh, G., & Widom, J. (2003, May). Scaling personalized web search. In *Proceedings of the 12th international conference on World Wide Web* (pp. 271-279). Acm

- Ronen, S., Gonçalves, B., Hu, K. Z., Vespignani, A., Pinker, S., & Hidalgo, C. A. (2014). Links that speak: The global language network and its association with global fame. *Proceedings of the National Academy of Sciences*, *111*(52), E5616-E5622
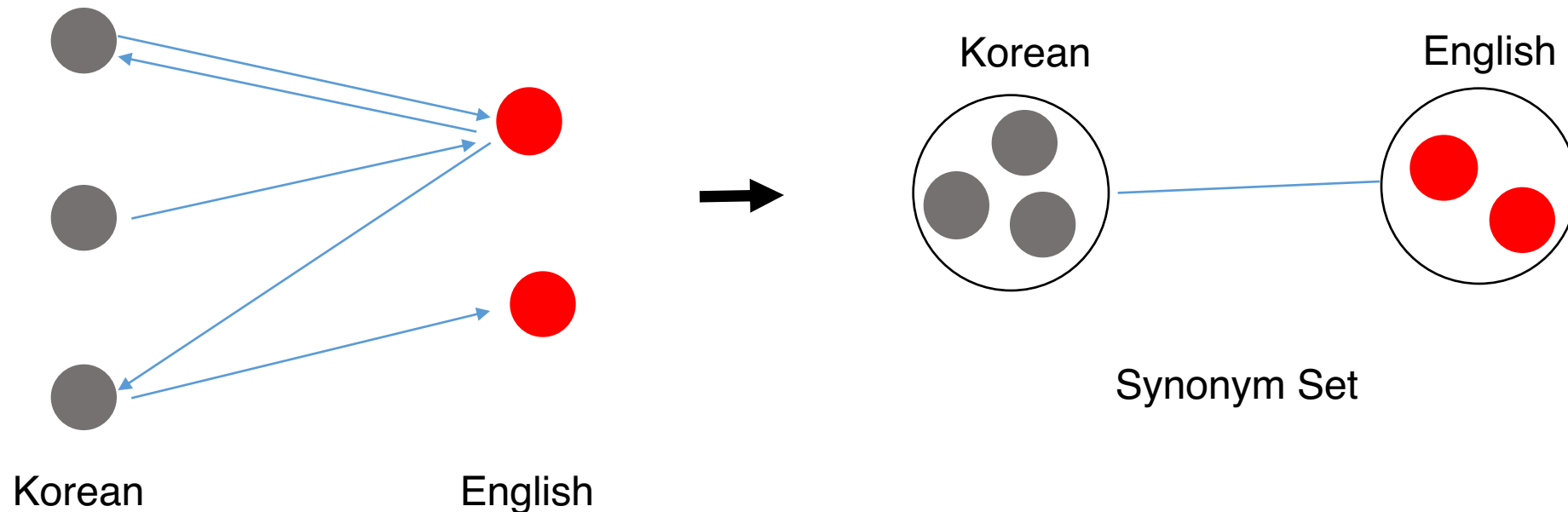
# Thank you for attention

Jisung.yoon@postech.ac.kr

# Appendix 1. Literatures on Wikipedia

- Dynamics on editing Wikipedia
  - Dynamics and pattern of modification (Yasseri et al., 2012a; Yasseri et al., 2012b),
  - Mechanistic model for intellectual interchanges (Yun et al., 2016)


- Credibility of Wikipedia data
  - TBA


- Data analysis with Wikipedia data
  - Extracting knowledge structure of Wikipedia (Ponzetto andNavigli;2009;Gabella,2017 )
  - Clustering of languages across the wikipedia growth (Ban, 2017)

# Appendix 2. Matching with language link data set

- For more general case (Many to Many)
  - Pairwise between two different language editions
  - **Node**: each category or page
  - **Link**: directed If A is connected as same documents B, A->B
    - E.g.) Republic of Korea (en) => 한국(ko)
  - Remove direction and merge after construction (likes synonym set)
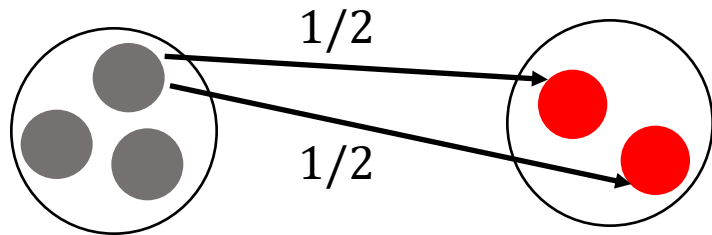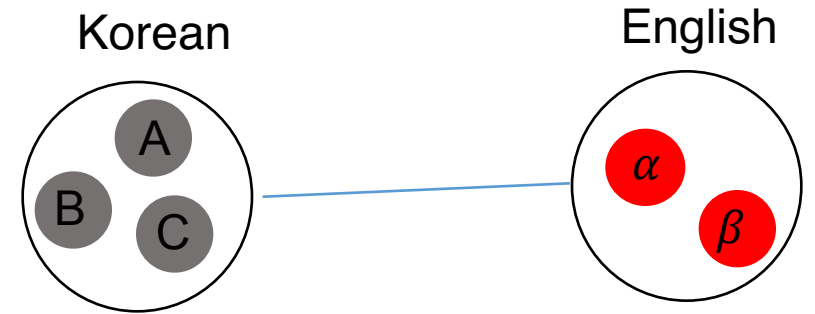


Korean          English

Korean          English

Synonym Set

# Appendix 3. Calculate subject distance

- For general case, many to many

Korean genealogy vector of node A, $X_A \in \boldsymbol{R^{N_k}}$
English genealogy vector of node B, $Y_\alpha \in \boldsymbol{R^{N_E}}$
Transition matrix, $\boldsymbol{T_{K \to E}} \in \boldsymbol{R^{N_k * N_E}}$
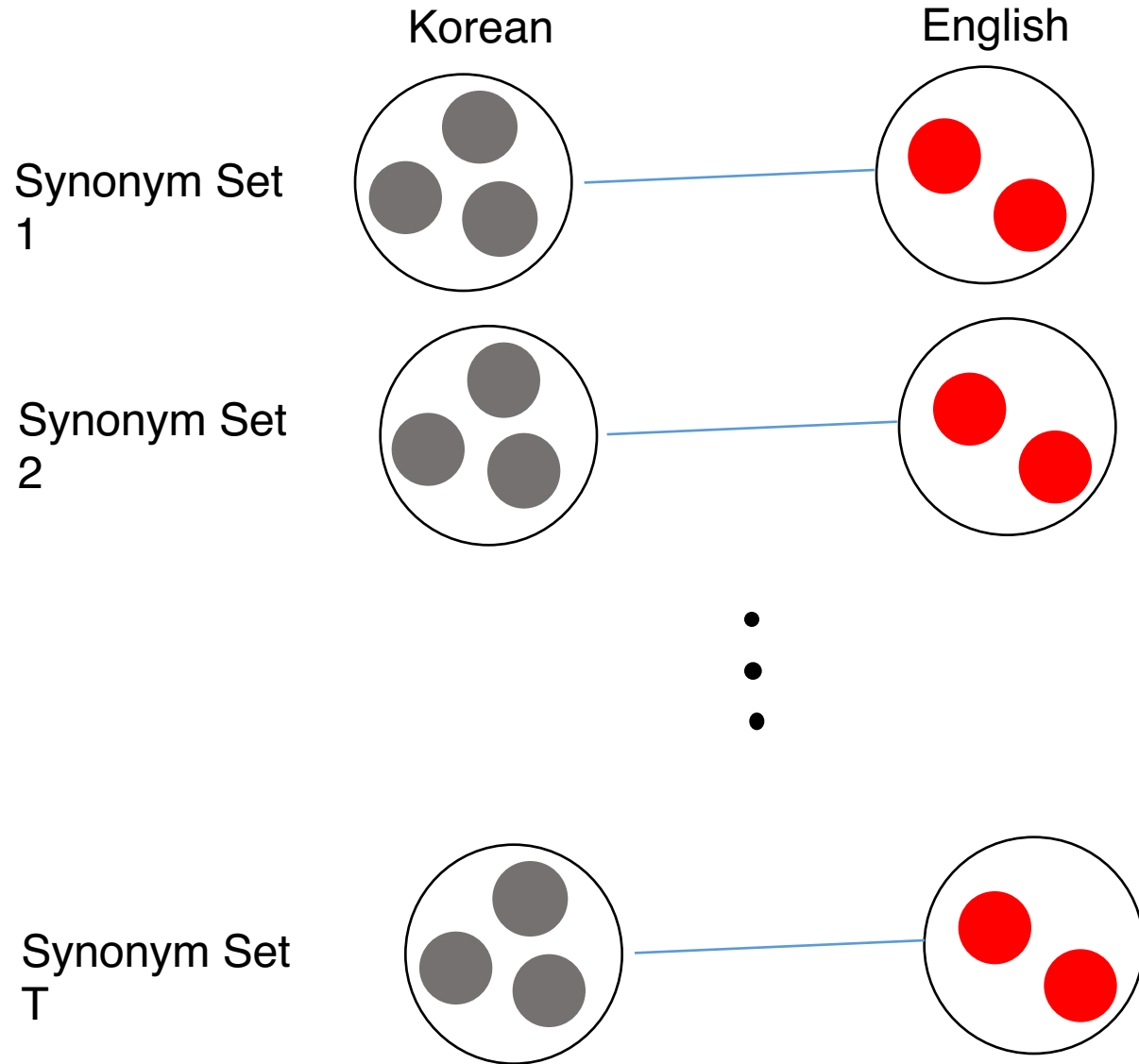
Korean

English

Synonym Set 1

1/2

1/2

$$d_{A\alpha}^{\boldsymbol{k \to E}} = D(X_A * T_{k \to E}, Y_\alpha)$$

Calculate Euclidean distance between English version of
Korean genealogy vector and English genealogy vector

$$s_1^{\boldsymbol{k \to E}} = \boldsymbol{1} - \frac{1}{N * M} * \sum_{i=1}^{N} \sum_{j=1}^{M} d_{ij}^{\boldsymbol{k \to E}}$$

# Appendix 4. Calculate overall similarity

Korean            English

Synonym Set 1

Synonym Set 2

Synonym Set T

$$s^{K \to E} = \frac{1}{T} * \sum_{i=1}^{T} s_i^{K \to E}$$

$$s^{E \to K} = \frac{1}{T} * \sum_{i=1}^{T} s_i^{E \to K}$$

# Appendix 5. Extracting Backbone of similarity network

- Extracting Backbone of similarity network
  - Relative similarity

$$r_{ij} = \frac{\dfrac{s_{ij}}{\sum_j s_{ij}}}{\dfrac{\sum_j s_{ij}}{\sum_i \sum_j s_{ij}}} = \sum_i \sum_j s_{ij} * \frac{s_{ij}}{\sum_j s_{ij} * \sum_j s_{ij}}$$

$$r_{ij} = Strength_{total} * \frac{s_{ij}}{strength_{out} \ of \ i \ * strngth_{in} \ of \ j}$$

- Select edges that higher than threshold
  - we select 1.04 which network fully connected to one weakly-connected component
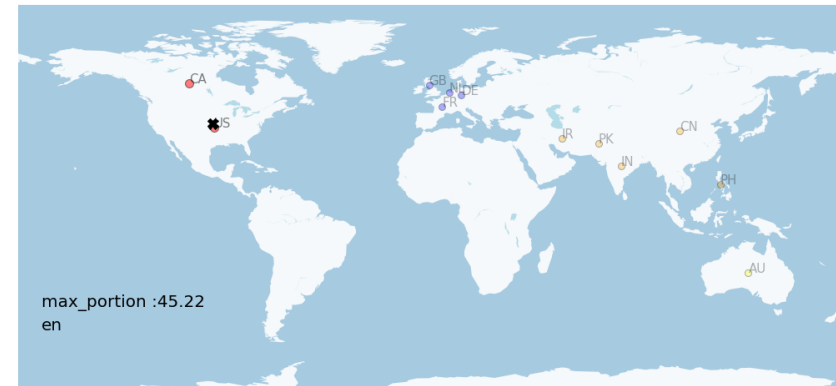
# Appendix 6. Language Code

| Code | Name | Code | Name | Code | Name | Code | Name | Code | Name |
|------|------|------|------|------|------|------|------|------|------|
| af | Afrikaans | el | Greek | hr | Croatian | ms | Malay | sr | Serbian |
| ar | Arabic | en | English | hy | Armenian | nl | Dutch | sv | Swedish |
| az | Azerbaijani | es | Spanish | id | Indonesian | nn | Norwegian nynorsk | ta | Tamil |
| be | Belarussian | et | Estonian | it | Italian | no | Norwegian | th | Thai |
| bg | Bulgarian | eu | Basque | ja | Japanese | pl | Polish | tl | Tagalog |
| bn | Bangla | fa | Persian | ka | Georgian | pt | Portuguese | tr | Turkish |
| bs | Bosnian | fi | Finnish | kk | Kazakh | ro | Romanian | uk | Ukrainian |
| ca | Catalan | fr | French | ko | Korean | ru | Russian | ur | Urdu |
| cs | Czech | gl | Galician | lt | Lithuanian | sco | Scots | vi | Vietnamese |
| cy | Welsh | he | Hebrew | lv | Latvian | sh | Serbo-croatian | zh | Chinese |
| da | Danish | hi | Hindi | ml | Malayalam | sk | Slovak | zh_yue | Cantonese |
| de | German | hu | Hungarian | mr | Marathi | sl | Slovenian | | |

# Appendix 7. Location of Language

- Each Wikipedia has a page view by country statistics.

  1. Get centroid locations of each country

  2. Conduct geo-location clustering, and get max portion cluster

     - To reduce noise

  3. Get weighted centroid of max portion cluster

Page views by country

| Page views | Name |
|---|---|
| 3B | United States of America |
| 744M | United Kingdom |
| 682M | India |
| 325M | Canada |
| 231M | Australia |
| 190M | Germany |
| 190M | Iran, Islamic Republic of |

max_portion :45.22
en

# Appendix 8. Language Projection Method

- Basically, socio-economic data are county to country data.

- For our analysis, we develop a method that projects county to country data to language to language data.

- Language projection method

  - $Y_{l \to L} = A_{L \to c}^{T} * X_{C \to C} * A_{L \to C}$, Language projected data

    - $X_{C \to C} \in R^{N_C * N_C}$, Country to country socio-economic data

    - $A_{L \to C} \in R^{N_C * N_L}$, Country to language matching matrix (Ronen et al, 2014)

      - e.g.) South Korea → 100% Korean

      - e.g.) United States → 82.1% English, 10.7% Spanish