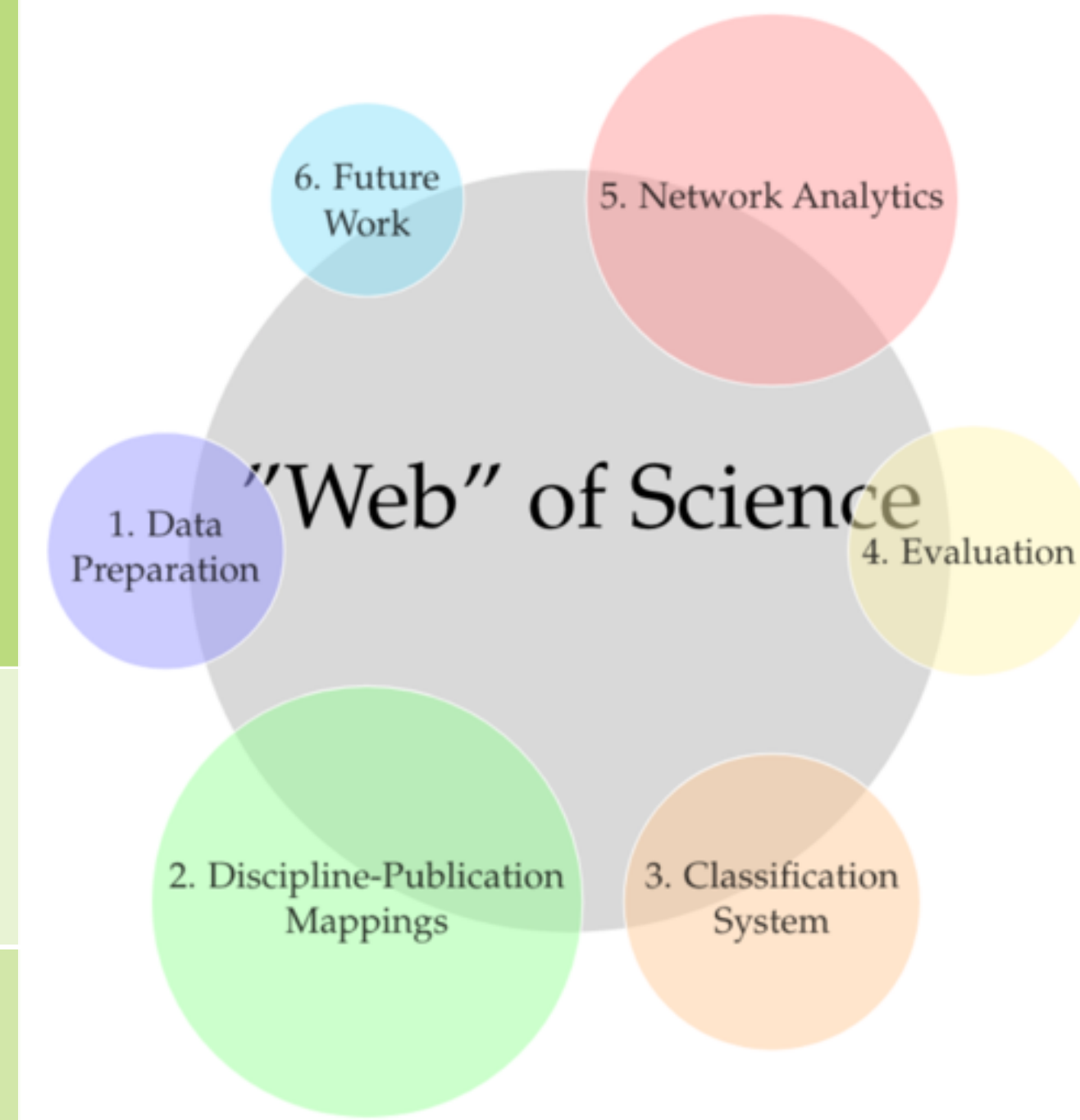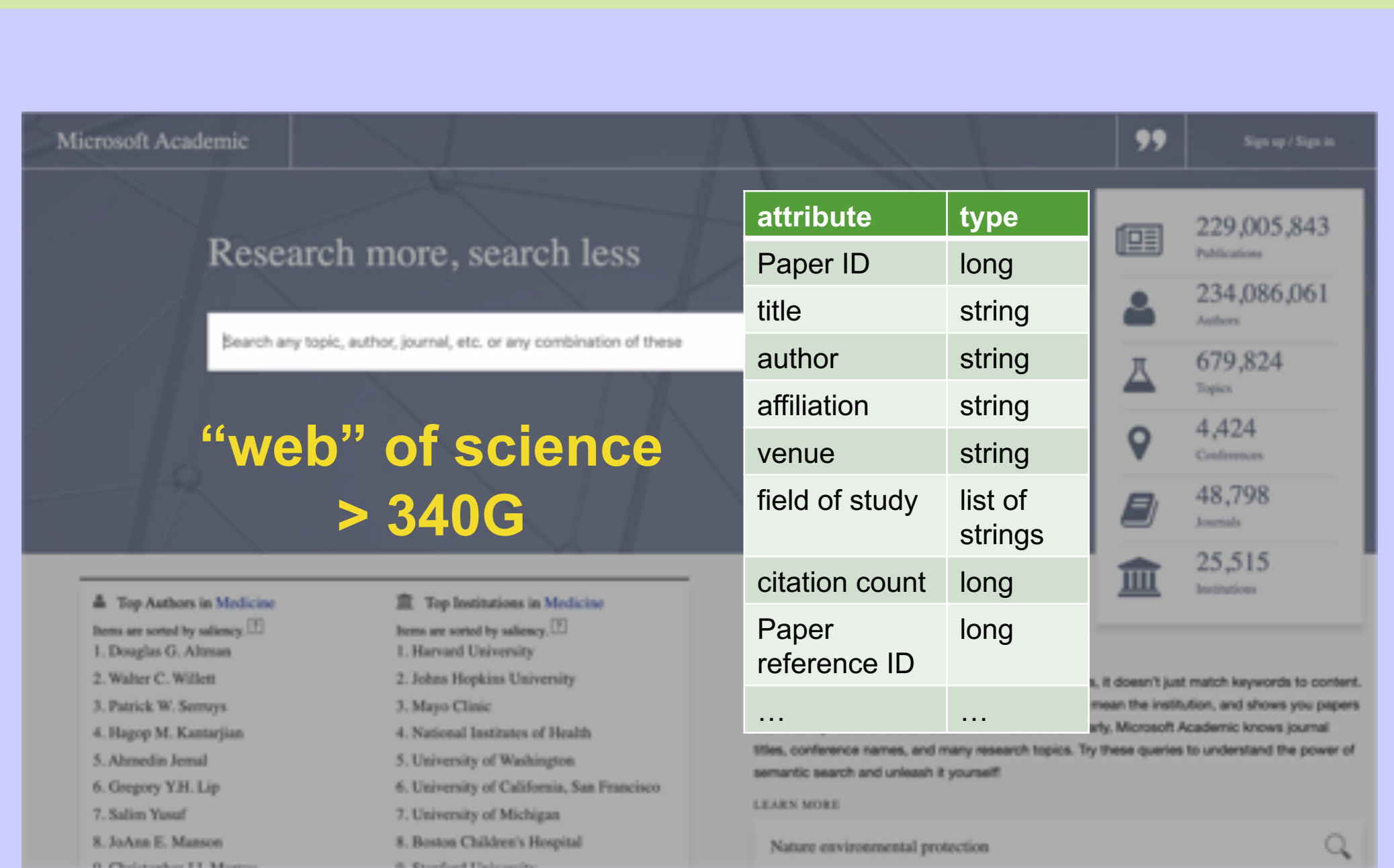# Understanding the Web of Science Using Deep Learning

Susie Xi Rao[1,2], Ce Zhang[1], Peter H. Egger[2]
[1] DS3Lab, ETH Zurich;  [2] KOF Swiss Economic Institute
Systems Group Retreat 2020
rao@kof.ethz.ch

6. Future Work
5. Network Analytics
"Web" of Science
1. Data Preparation
4. Evaluation
2. Discipline-Publication Mappings
3. Classification System

**Research Goal:**
1) A system that takes a growing amount of scholarly publications from each discipline and tells you to which discipline it belongs to.
2) Creation of networks that depict the publication behaviors of authors and institutions and how they impact the innovations and economic developments around the globe.

**1 Data Preparation (X)**
Microsoft Academic Service, frequent update, topics on different granularity for each discipline

Microsoft Academic

Research more, search less

Search any topic, author, journal, etc. or any combination of these

"web" of science > 340G

| attribute | type |
|---|---|
| Paper ID | long |
| title | string |
| author | string |
| affiliation | string |
| venue | string |
| field of study | list of strings |
| citation count | long |
| Paper reference ID | long |
| … | |

229,005,843
234,086,061
679,824
4,424
48,798
25,515

**2 Discipline-Publication Mappings (YL1 – YL2 – YL3)**

56 discipline by Wikipedia

2 Humanities
  2.1 Anthropology
    2.1.1 Archaeology
  2.2 History
  2.3 Linguistics and languages
  2.4 Philosophy
  2.5 Religion
  2.6 The arts
    2.6.1 Culinary arts
    2.6.2 Literature
    2.6.3 Performing arts
    2.6.4 Visual arts
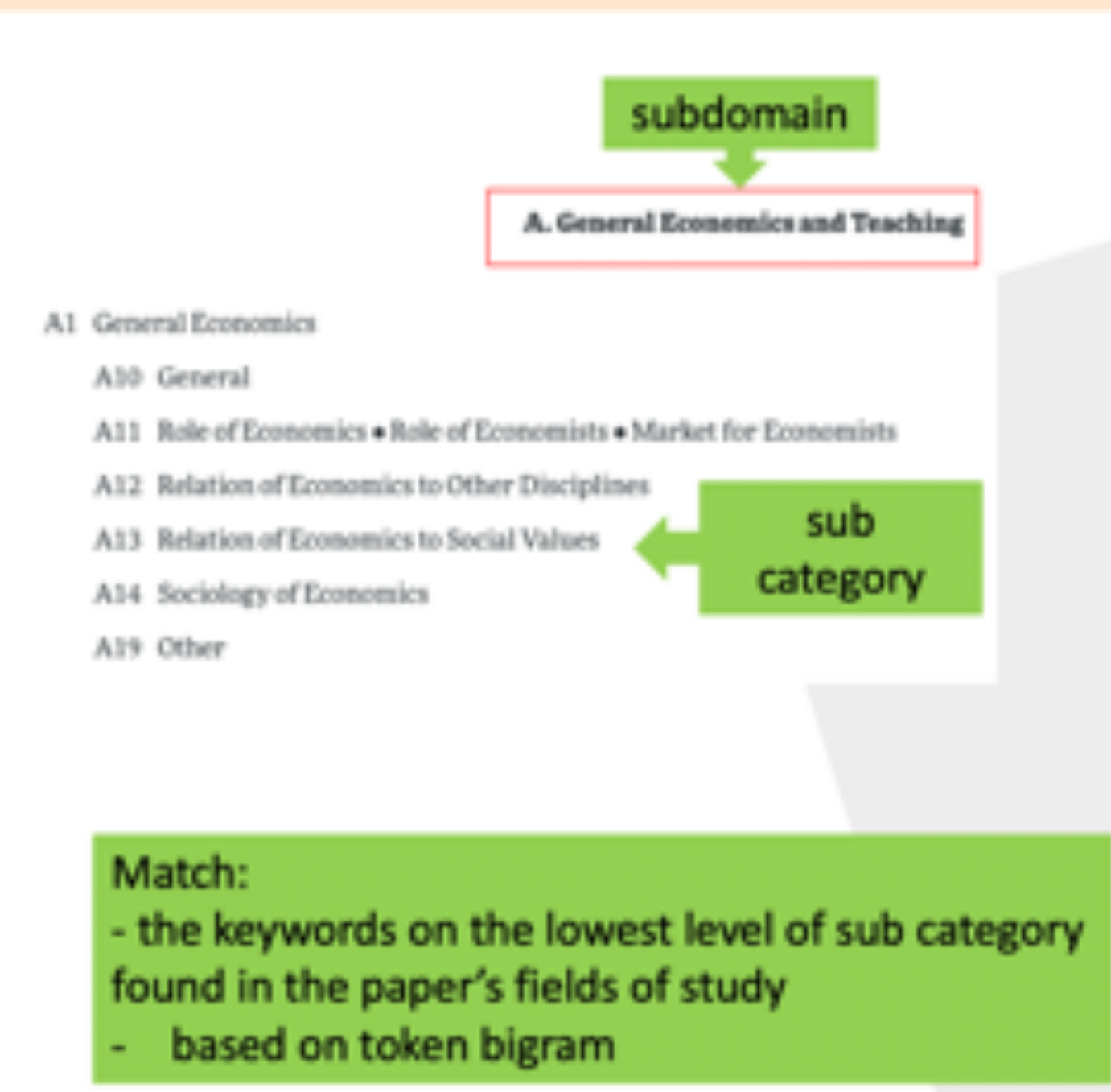
ACM classification for CS

- H.: Information Systems
  o H.0: GENERAL
  o H.1: MODELS AND PRINCIPLES
    - H.1.0: General
    - H.1.1: Systems and Information Theory
    - H.1.2: User/Machine Systems
    - H.1.m: Miscellaneous
  o H.2: DATABASE MANAGEMENT
    - H.2.0: General
    - H.2.1: Logical Design
    - H.2.2: Physical Design
    - H.2.3: Languages
    - H.2.4: Systems
    - H.2.5: Heterogeneous Databases
    - H.2.6: Database Machines
    - H.2.7: Database Administration
    - H.2.8: Database Applications
    - H.2.m: Miscellaneous

JEL classification for econ

**3 Classification Systems (Generating training/testing material)**

Match:
- the keywords on the lowest level of sub category found in the paper's fields of study
- based on token bigram

subdomain
A. General Economics and Teaching
sub category

**L1: FNN/CNN/RNN**
Input Layer — Hidden Layer — Output High Level

$x_0$
$x_1$
$x_i$
$x_{f-1}$
$x_f$

$y_0$
$y_1$
$y_i$
$y_{m-1}$
$y_m$

Economics
Computer Science
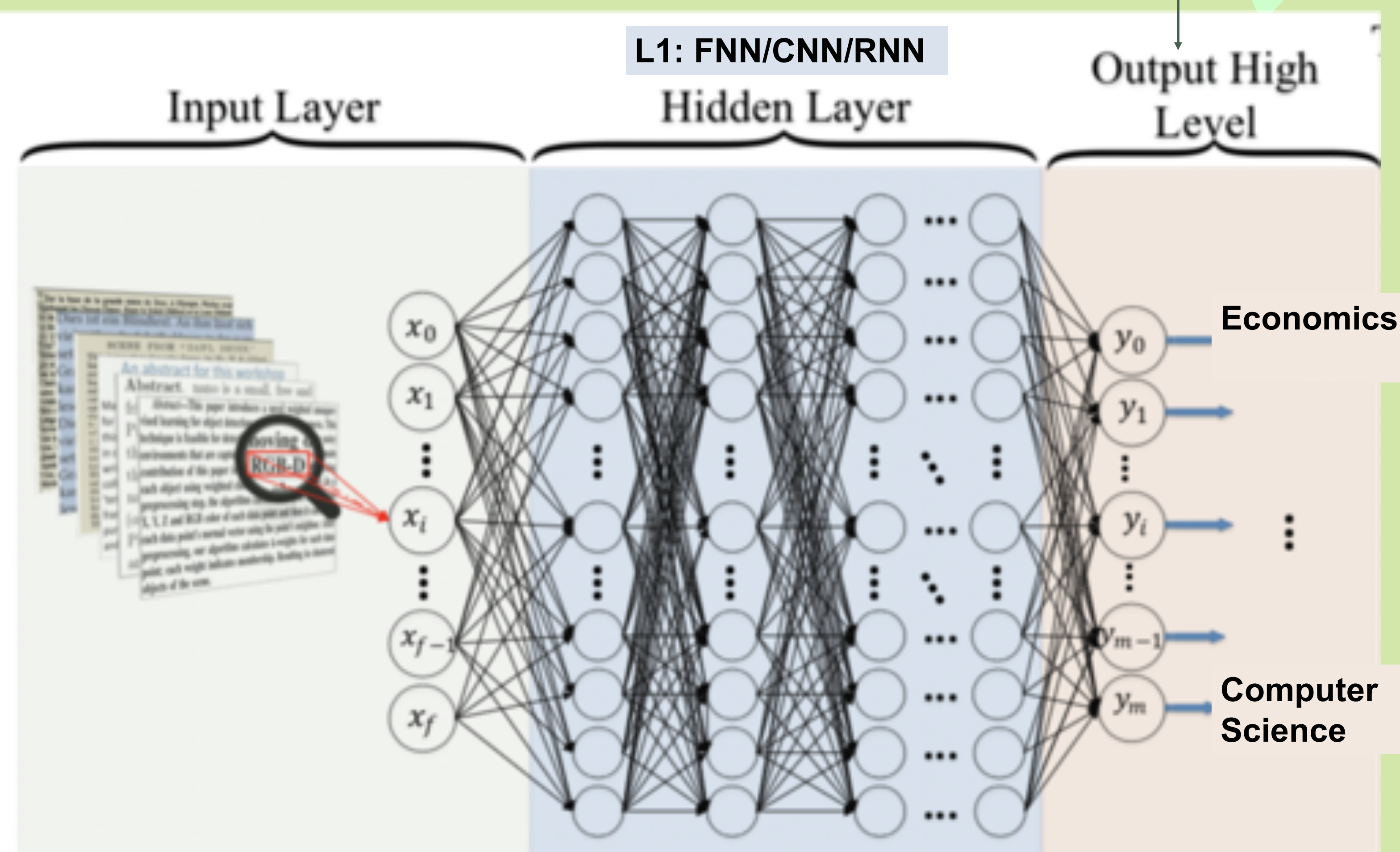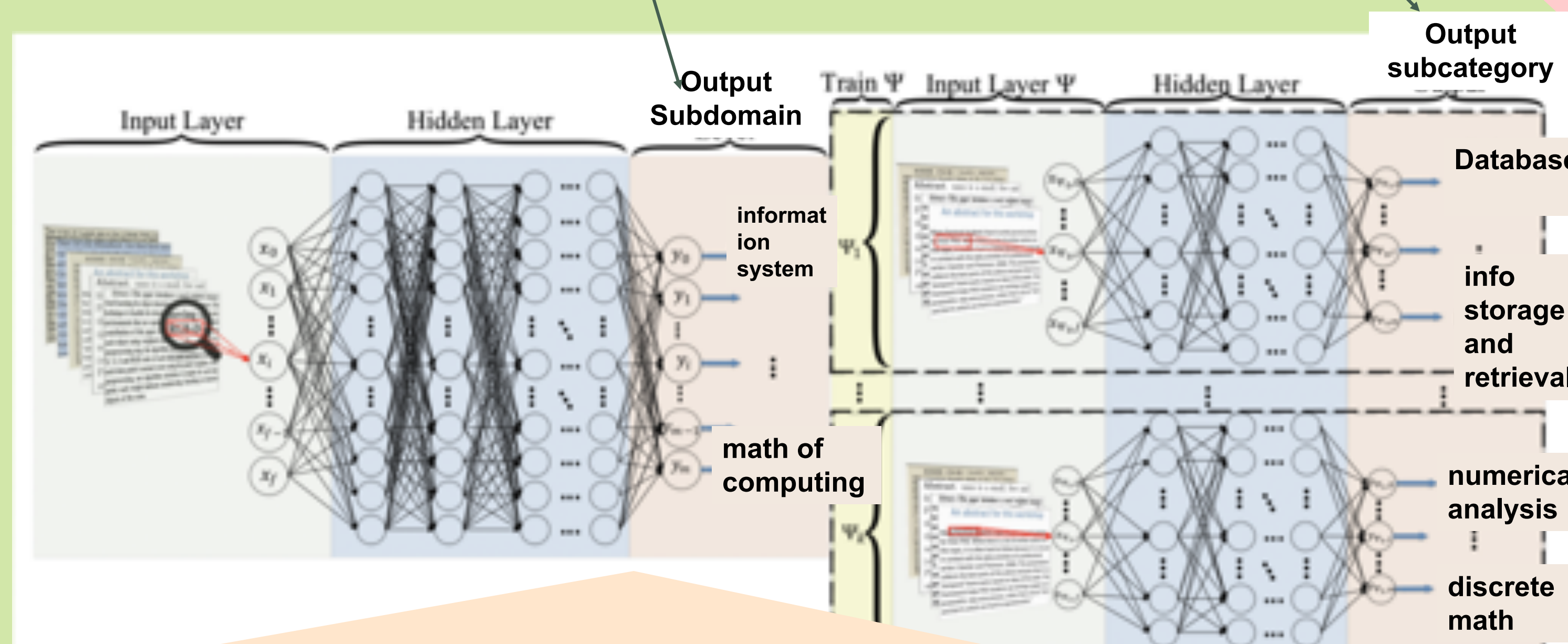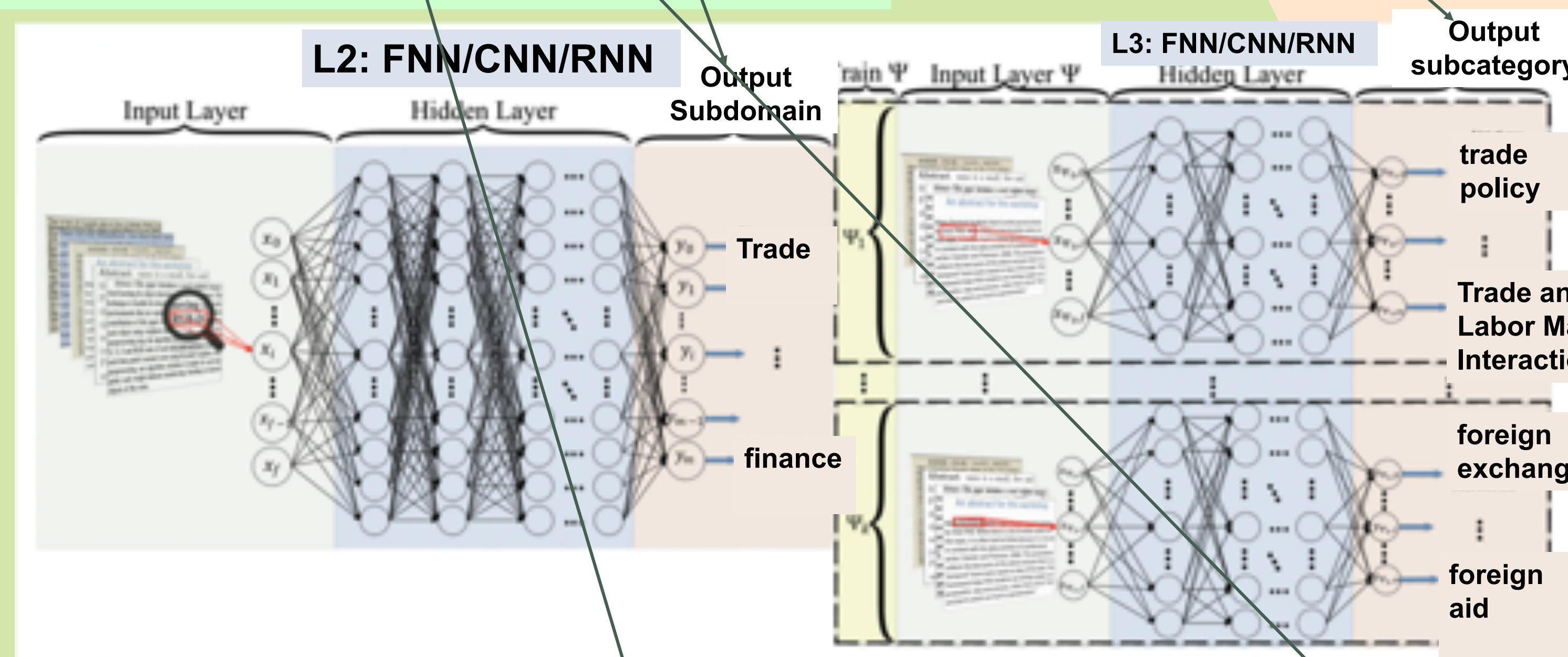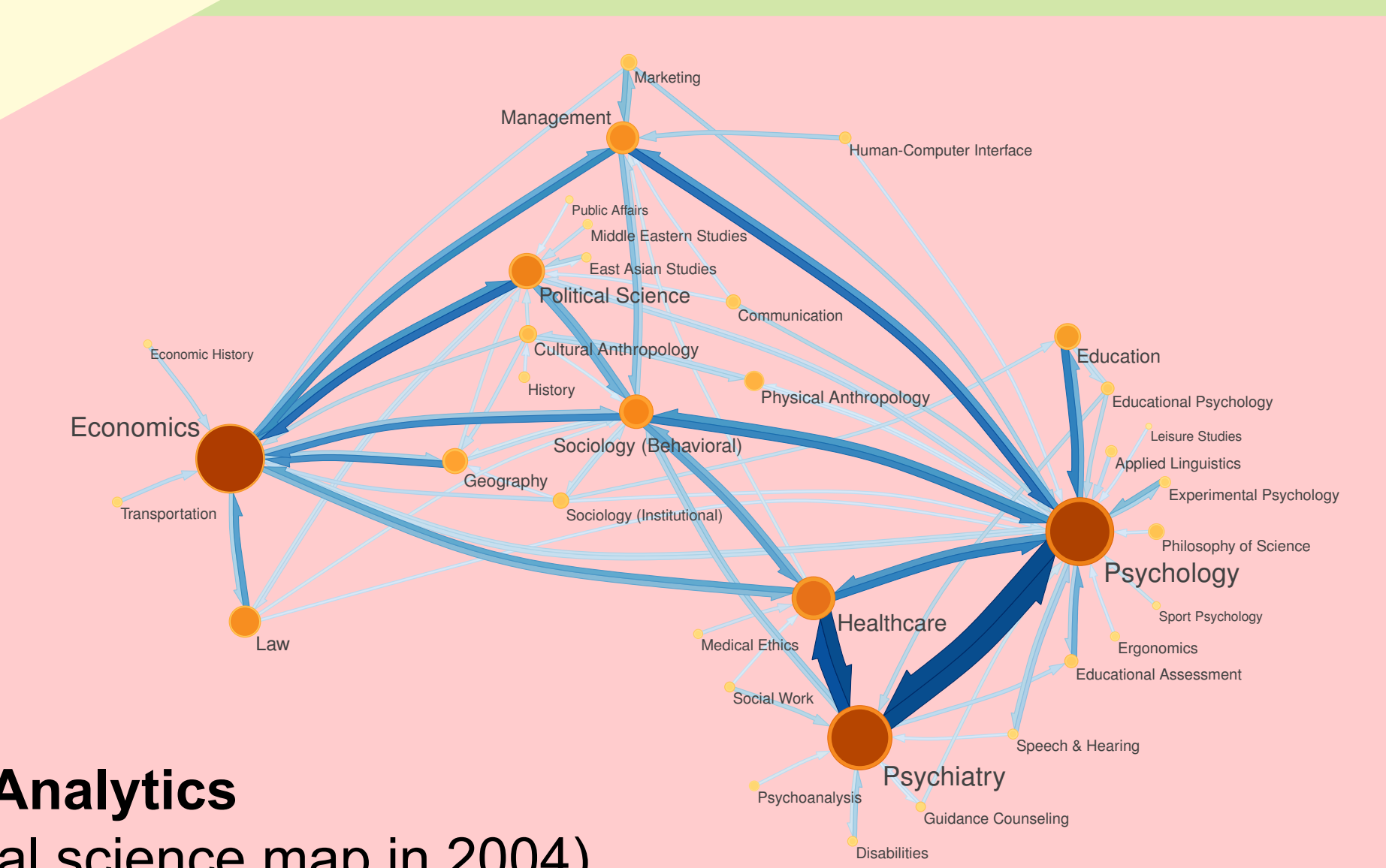
Fig. 1: HDLTex: Hierarchical Deep Learning for Text Classification. This is our Deep Neural Network (DNN) approach for text classification. The left figure depicts the parent-level of our model, and the right figure depicts child-level models defined by $\Psi_r$ as input documents in the parent level.

- Input: article abstracts (ca. 300 words)
- Output: three labels (discipline, subdomain, subcategory)

**L2: FNN/CNN/RNN**
Input Layer — Hidden Layer — Output Subdomain
Train $\Psi$  Input Layer $\Psi$  Hidden Layer — Output subcategory

Trade — trade policy
Trade and Labor Market Interactions
finance — foreign exchange — foreign aid

**L3: FNN/CNN/RNN**

Output Subdomain
information system — info storage and retrieval
Database
math of computing — numerical analysis — discrete math

**4 Evaluation (architecture search)**
- For L1 + L2, econ + CS: RNN + CNN works the best, > 90% accuracies for output subcategories
- For L1 + L2 + L3, econ + CS: CNN + CNN + CNN works the best, to be reported

**5 Network Analytics**
(e.g., a social science map in 2004)
Other interesting questions:
- Rise and fall of fields: Using our classification, what are the fields that became important sources of spillovers to other fields.
- Rise and fall of institutions
- We could "invent" a taxonomy of break-through innovations?

**3 Classification Systems (System optimization)**
- Capability to take in large datasets: precomputation of word index for large training corpus using MapReduce (6 min. for 24 mio. abstracts), precomputation of one-hot encodings of the abstracts in each discipline and arbitrary combination of training sets across disciplines are possible, precomputation of training/validation partitions on all the levels.
- Distributed learning using *tensorflow*.
- A hierarchical classification system application to other sorts of hierarchies.
- Softmax function at each classification step renders probabilities of discipline membership given one abstract → to what extent one publication belongs to one discipline/subdomain/subcategory.

References
[1] KHABSA, M., AND GILES, C. L. The number of scholarly documents on the public web. PloS one 9, 5 (2014), e93949.
[2] KOWSARI, K., BROWN, D. E., HEIDARYSAFA, M., MEIMANDI, K. J., GERBER, M. S., AND BARNES, L. E. HDLTex: Hierarchical deep learning for text classification. In Machine Learning and Applications (ICMLA), 2017 16th IEEE International Conference on (2017), IEEE, pp. 364–371.
[3] SHEN, Z., MA, H., AND WANG, K. A web-scale system for scientific knowledge exploration. arXiv preprint arXiv:1805.12216 (2018).
[4] SINHA, A., SHEN, Z., SONG, Y., MA, H., EIDE, D., HSU, B.-J. P., AND WANG, K. An overview of Microsoft academic service (MAS) and applications. In Proceedings of the 24th international conference on world wide web (2015), ACM, pp. 243–246.
[5] TANG, J., ZHANG, J., YAO, L., LI, J., ZHANG, L., AND SU, Z. Arnetminer: extraction and mining of academic social networks. In Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining (2008), ACM, pp. 990–998.
[6] WIKIPEDIA, List of academic fields, 2018.
[7] GOYAL, P., AND FERRARA, E. Graph embedding techniques, applications, and performance: A survey. Knowledge-Based Systems (2018), 151, 78-94.
[8] a social science map in 2004, see http://www.eigenfactor.org/map/maps.php

ETH zürich  Systems @ ETHzürich  DS3  KOF