

# Challenges in Chinese Knowledge Graph Construction

Chengyu Wang, Ming Gao, Xiaofeng He, Rong Zhang

Institute for Data Science and Engineering  
East China Normal University  
Shanghai, China



# Knowledge Graph -

## Modeling Knowledge as a Graph

### Entities

- Concepts
- Instances
- Values

**Nodes:** entities (concept, named entity, ...)

**Edges:** semantic relationships

### Knowledge Graph

### Relations

- IsA
- Co-occurrence
- Others



Google Knowledge Graph



Satori (Bing Search)

# Chinese Knowledge Graph

## Data Sources & Challenges

- Sources

- Heterogeneous data sources
- No public knowledge repositories or semantic networks

- Methods

- Machine translation: low quality
- Information extraction: difficult

### Chinese Wikis:



Chinese Wikipedia  
(0.8M+ articles)

维基百科  
自由的百科全书



Baidu Baike  
(10M+ articles)



Hudong Baike  
(11M+ articles)

# Data Sparsity

- Comparison between Chinese & English Wikipedias

|            | Chinese Wikipedia | English Wikipedia |                  |
|------------|-------------------|-------------------|------------------|
| #Articles  | ~0.8M             | ~4M               | <b>5 times!</b>  |
| #Infoboxes | ~0.1M             | ~1.6M             | <b>13 times!</b> |

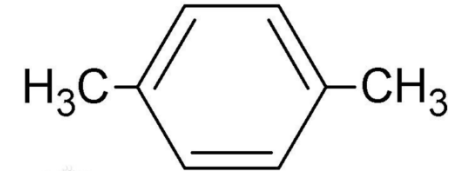
- Challenges

- Entities: extracting long-tailed entities
- Relations: construction of a “dense” KG

- Solution

- Data fusion from different sources

# Information Accuracy



- “Editing war” on PX (P-Xylene)
  - Polarized attitudes towards plan of a PX factory in city of Xiamen, China
  - Edited 76 times in total
  - Supporters: PX is slightly toxic.
  - Protesters: PX is extremely toxic!
- Challenges
  - Mining editing logs
  - Detecting inaccurate attributes

## Editing log on PX

|                  |                    |                            |
|------------------|--------------------|----------------------------|
| 2014-04-02 16:37 | <a href="#">查看</a> | <a href="#">溺水三千s</a>      |
| 2014-04-02 13:53 | <a href="#">查看</a> | <a href="#">1162007677</a> |
| 2014-04-02 12:27 | <a href="#">查看</a> | <a href="#">枫之群動</a>       |
| 2014-04-01 11:20 | <a href="#">查看</a> | <a href="#">cycrc7</a>     |
| 2014-03-31 21:56 | <a href="#">查看</a> | <a href="#">邓靖轩</a>        |
| 2014-03-31 20:41 | <a href="#">查看</a> | <a href="#">zhiyuanep1</a> |
| 2014-03-31 18:49 | <a href="#">查看</a> | <a href="#">道牙子没事</a>      |
| 2014-03-31 17:05 | <a href="#">查看</a> | <a href="#">sunbingame</a> |
| 2014-03-31 12:05 | <a href="#">查看</a> | <a href="#">道牙子没事</a>      |
| 2014-03-31 10:54 | <a href="#">查看</a> | <a href="#">道牙子没事</a>      |
| 2014-03-30 20:33 | <a href="#">查看</a> | <a href="#">847872000</a>  |
| 2014-03-30 18:41 | <a href="#">查看</a> | <a href="#">亚豆亚豆</a>       |
| 2014-03-30 00:09 |                    | <a href="#">幻想书生wjc</a>    |

# Link Quality

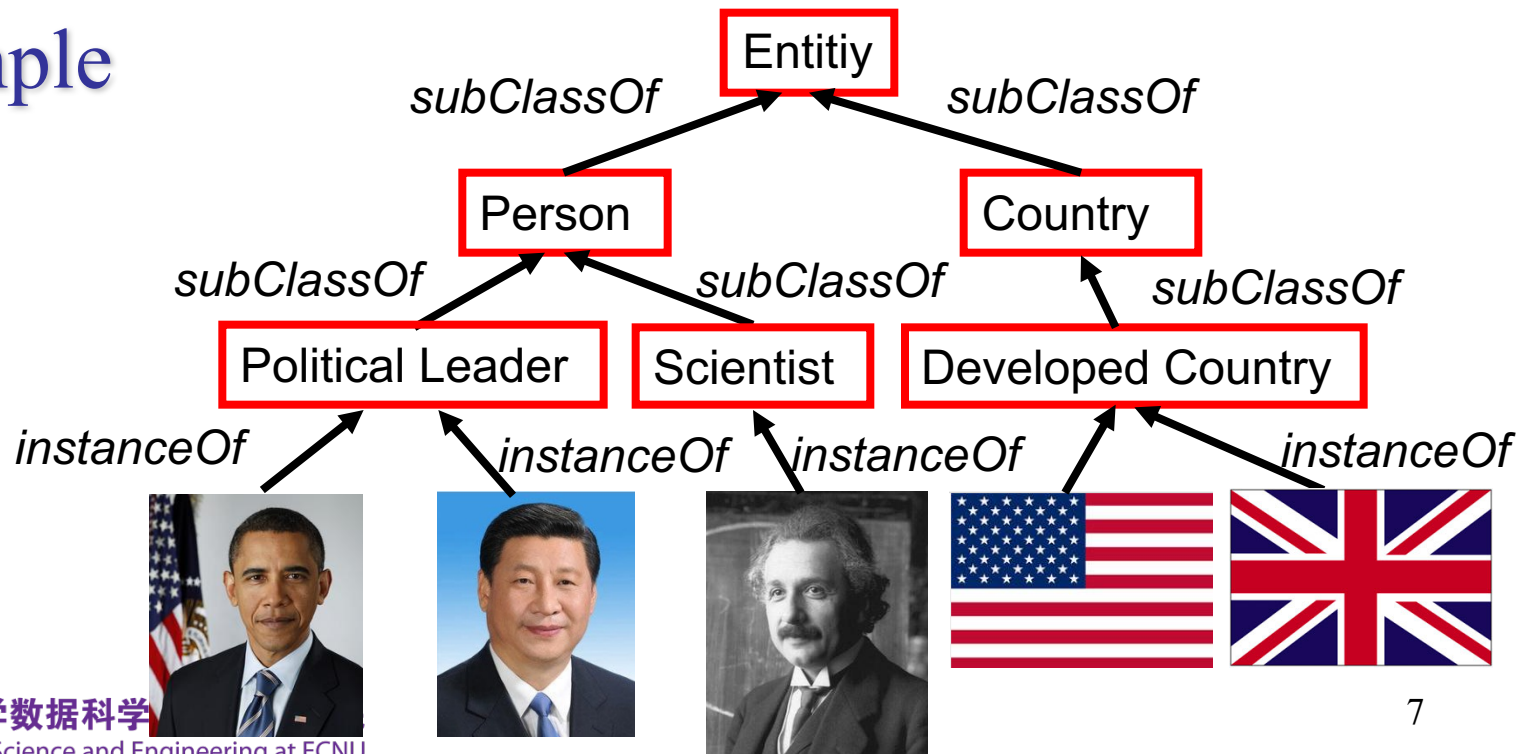
- **Hyperlinks in Wikipedia**
  - Link entity mentions in texts with corresponding Wikipedia pages
  - Serve as evidence to perform entity linking  
Barack Hussein Obama II is the [44th](#) and [current President](#) of the [United States](#), and the [first African American](#) to hold the office.
- **Wrongly annotated links in Chinese Wikipedia**
  - **Wu Mei** (Prof of Peking Univ.) in page [May Fourth Movement](#) linked to [Wu Mei](#) (dubbing actress in Hong Kong)
  - Automatic detection of error links in Wikipedia

# Taxonomy Derivation

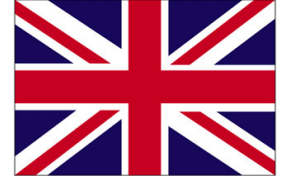
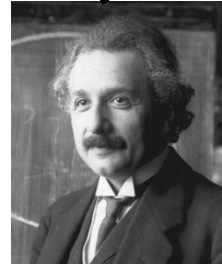
- Taxonomy: a hierarchical type system for KGs
  - **subClassOf** relations (subject: class, object: class)
  - **instanceOf** relations (subject: entity, object: class)

- Example

**Classes**



**Entities**



# Taxonomy Derivation

- Challenges in Chinese taxonomy derivation
  - Lack of resources (No Chinese equivalent of WordNet)
  - Hard to map entities to their categories

## Research directions

- Language patterns
- Classification
- Machine translation
- Complete taxonomy construction

习近平 锁定 ★ 收藏 👍 81502 📄 40593

习近平，男，汉族，1953年6月生，陕西富平人，1969年1月参加工作，1974年1月加入中国共产党，清华大学人文社会学院马克思主义理论与思想政治教育专业毕业，在职研究生学历，法学博士学位。

现任中国共产党中央委员会总书记，中共中央军事委员会主席，中华人民共和国主席，中华人民共和国中央军事委员会主席。<sup>[1]</sup>

[人物关系](#) 编辑

## Xi Jinping (Chinese President)

[图片](#) [习近平图册](#)

|      |         |      |       |
|------|---------|------|-------|
| 中文名  | 习近平     | 毕业院校 | 清华大学  |
| 别名   | 习大大     | 信仰   | 共产主义  |
| 国籍   | 中华人民共和国 | 入党时间 | 1974年 |
| 民族   | 汉族      | 籍贯   | 陕西富平  |
| 出生日期 | 1953年6月 | 学位   | 博士    |

**词条统计**  
浏览次数: 36939626次  
编辑次数: 109次 [历史版本](#)  
最近更新: 2015-02-12  
创建者: vusef

词条标签: 人物, 政治人物, 政治, 官员

Labels: Person, Politician, Politics, Official

relatedTo? topicOf?

subClassOf? instanceOf?



# IsA Extraction

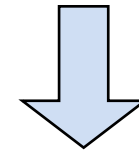
- Hearst patterns (Hearst. COLING'92)

- such NP as NP,\* or|and NP
- NP such as NP, NP, ..., and|or NP
- NP, including NP,\* or | and NP
- ...

- Chinese IsA patterns

- Poor NLP analysis in Chinese Web text
- Lack of explicit high-quality isA patterns
- Implicit expressions of isA relations

**Countries** such as **China**,  
**France** and **Germany**



**China** isA **Country**  
**France** isA **Country**  
**Germany** isA **Country**

**ProBase**

Largest taxonomy in  
English

- 2.6M+ concepts
- 20M+ isA pairs

# General Relation Extraction

- Relation extraction systems

Snowball (SIGMOD'01)

KnowItAll (WWW'04)

LELIA (KDD'06)

TextRunner (IJCAI'07)

StatSnowball (WWW'09)

Many others...

- Focus on English language

- Chinese relation extraction

- Extract knowledge from semi-structured and structured data

- Design statistical and NLP-based features for Chinese text

- Use facts of high precision to supervise RE process (distant supervision)

# Conclusion

- Web-scale Chinese KG construction
  - Quality of data sources: data fusion and cleaning
  - Taxonomy derivation: study on taxonomic relations in Chinese
  - Knowledge harvesting: isA patterns, Chinese RE systems

**Lots of challenges**

**Lots to do!**

**Thanks!**

Questions & Answers