

Zhishi.me

Weaving Chinese Linking Open Data

Xing Niu¹, Xinruo Sun¹, Haofen Wang¹, Shu Rong¹, Guilin Qi² and Yong Yu¹

¹Shanghai Jiao Tong University

²Southeast University

2011.10.25

Agenda

- Introduction
- Semantic Data Extraction
- Web Access to Zhishi.me
- Data-level Mapping among Different Datasets
- Conclusions and Future Work



Introduction

- Why we need Chinese Linking Open Data (CLOD)?
- Linked Open Data contains very sparse Chinese knowledge at the present time.
 - *Freebase*: Chinese labels
 - *DBpedia*: Chinese labels and short abstract [1]
 - *UWN* (Universal WordNet): maps Chinese words to corresponding vocabulary entities in WordNet [2]

Introduction (con't)

- In order to attract more efforts to publish Chinese semantic data linked to the CLOD, we prefer building a hub:
 - Covers all branches of knowledge
 - Contains abundant data

- We extracted knowledge from textual articles of other independent data sources
 - DBpedia is a representative one
 - Chinese Web-based collaborative encyclopedias together contain even more articles than the largest English one: Wikipedia

您已登录到该词条，更多详情请查看上海。

上海

Labels

百科名片



上海，中国大陆第一大城市，四个中央直辖市之一；是中国大陆的经济、金融、贸易和航运中心。上海创造了和打破了中国世界纪录协会多项世界之最、中国之最。上海位于中国最大河流长江中部的长江口，拥有中国最大的外贸港口、最大的工业基地。有超过2000万人居住和生活在上海地区，其中大部分属汉族江浙民系，通行吴语上海话。上海又是一座新兴的旅游目的地，具有深厚的近代城市文化底蕴和众多的历史古迹。如今上海已发展成为全球国际化大都市，并致力于在2020年建设成为国际金融中心

Table with 2 columns: 中文名称, 外文名称, 别名, etc. for Shanghai.

Internal Links

- 历史沿革, 上海概况, 行政区划, etc.

历史沿革

春秋时属吴，吴灭后入越，越王无疆大败于楚，后成为楚国置郢地（故上海别称“申”），前223年秦灭越后设会稽郡（治所在今苏州）并辖有郢县（今金山），由会稽（今嘉兴）和海盐县等。郢县包括今上海市大部（其余尚未成陆），前207年郢县改名郢县。晋时松江（吴淞江）居民创造一种竹编捕鱼工具“罾”，后又因流入海处称“濬”，松江下游一带被称为“濬”，后又改名为“沪”。唐天宝十年（751年），吴郡刺史赵居贞上书，奏请划昆山南境、嘉兴东境、海盐北境，设华亭县。宋淳化二年（991年）松江淤浅，船舶无法上海湾。南宋咸亨三年（1267年）华亭县在上海西岸设“上海镇”。元初至元十四年（1277年），华亭县升格为府，翌年更名松江府，辖华亭县（一府一厅）。元初至元十九年（1292年），析华亭县部分，设上海县，均隶于松江府。上海县属治在今黄浦区，这是上海建城的开始。

扩展阅读

- 上海市人民政府门户网站, 上海市人民政府新闻办公室, etc.

开放分类

中国, 直辖市, 上海, 申城, 沪上

上海 相关词条

伦敦 纽约 巴黎 东京 芝加哥 法兰克福 香港 洛杉矶 米兰 旧金山 悉尼 多伦多 苏黎世 布鲁塞尔 马德里

上海市

上海市的概述, 上海市 - 概述, 上海市位于北纬31度14分, 东经121度29分. 地处长江三角洲南缘, 东濒东海, 南临杭州湾, 西接江苏、浙江两省, 北界长江入海口, 正当中国南北海岸线的中部, 交通便利, 腹地广阔, 地理位置优越, 是一个良好的江海港口.



上海, 中国大陆第一大城市, 四个中央直辖市之一; 是中国大陆的经济、金融、贸易和航运中心。上海创造了和打破了中国世界纪录协会多项世界之最、中国之最。上海位于我国大陆海岸线中部的长江口, 拥有中国最大的外贸港口、最大的工业基地。有超过2000万人居住和生活在上海地区, 其中大部分属汉族江浙民系, 通行吴语上海话。上海又是一座新兴的旅游目的地, 具有深厚的近代城市文化底蕴和众多的历史古迹。如今上海已发展成为全球国际化大都市, 并致力于在2020年建设成为国际金融中心

目录: 1 概述, 2 历史沿革, 3 行政区划, 4 地理气候, 5 人口

上海市 - 概述

上海市位于北纬31度14分, 东经121度29分. 地处长江三角洲南缘, 东濒东海, 南临杭州湾, 西接江苏、浙江两省, 北界长江入海口, 正当中国南北海岸线的中部, 交通便利, 腹地广阔, 地理位置优越, 是一个良好的江海港口.

地形 除西南部有少数丘陵外, 上海境内全为坦荡低平的平原, 是长江三角洲冲积平原的一部分.

能源 煤、石油、水力的储藏和生产, 所需能源都靠其他省市的支援. 但是, 上海具有一定数量和较高质量的“二次能源”. 产品主要是电力、石油油品、煤炭和煤气 (包括液化石油气). 其他可以利用开发的能源还有沼气、风能、潮汐能及太阳能.

生物 上海濒临东海, 有丰富的水产资源, 据统计, 东海、南海的水产资源有200多种. 此外, 上海地处长江口, 这里江面宽阔, 海淡水交汇, 是鱼类繁衍、繁殖、栖息的场所, 有各种鱼类200多种, 其中经济鱼类有20多种. 上海有众多的天然湖泊, 蝴蝶礁等底栖生物资源非常丰富. 掘泥的水网, 为淡水养殖提供了良好的条件.

航运 至2006年末, 上海拥有16条国际集装箱班轮航线; 上海港每月集装箱航班数已达1867班次, 其中国际航班42班次.

远洋航线: 从上海始可以分别抵达香港、台湾 (经第三地)、韩国、日本、东南亚、澳大利亚、以色列、中东、西北欧、南非、南美、美国东西岸等地;

沿海航线: 可通达从北到市沿海主要港口;

长江航线: 可直达长江中下游各港口;

内河航线

开放分类: 中国城市, 中国直辖市, 中国省级行政区, 五角场

城市 (地区)

上海市的概述, 上海市 - 概述, 上海市位于北纬31度14分, 东经121度29分. 地处长江三角洲南缘, 东濒东海, 南临杭州湾, 西接江苏、浙江两省, 北界长江入海口, 正当中国南北海岸线的中部, 交通便利, 腹地广阔, 地理位置优越, 是一个良好的江海港口.

上海市 (4) 发音 (帮助·信息); 吴语: Zaoonhe), 简称沪, 别称申, 中华人民共和国直辖市、国家中心城市。上海位于中国南北弧形海岸线中部, 长江三角洲东部, 全市最北部为处于长江入海口中的崇明岛, 东向东海, 并隔海与日本九州岛相望, 南濒杭州湾, 西部与江苏、浙江两省相接。

上海市 (4) 发音 (帮助·信息); 吴语: Zaoonhe), 简称沪, 别称申, 中华人民共和国直辖市、国家中心城市。上海位于中国南北弧形海岸线中部, 长江三角洲东部, 全市最北部为处于长江入海口中的崇明岛, 东向东海, 并隔海与日本九州岛相望, 南濒杭州湾, 西部与江苏、浙江两省相接。

上海市 (4) 发音 (帮助·信息); 吴语: Zaoonhe), 简称沪, 别称申, 中华人民共和国直辖市、国家中心城市。上海位于中国南北弧形海岸线中部, 长江三角洲东部, 全市最北部为处于长江入海口中的崇明岛, 东向东海, 并隔海与日本九州岛相望, 南濒杭州湾, 西部与江苏、浙江两省相接。

上海市 (4) 发音 (帮助·信息); 吴语: Zaoonhe), 简称沪, 别称申, 中华人民共和国直辖市、国家中心城市。上海位于中国南北弧形海岸线中部, 长江三角洲东部, 全市最北部为处于长江入海口中的崇明岛, 东向东海, 并隔海与日本九州岛相望, 南濒杭州湾, 西部与江苏、浙江两省相接。

上海市 (4) 发音 (帮助·信息); 吴语: Zaoonhe), 简称沪, 别称申, 中华人民共和国直辖市、国家中心城市。上海位于中国南北弧形海岸线中部, 长江三角洲东部, 全市最北部为处于长江入海口中的崇明岛, 东向东海, 并隔海与日本九州岛相望, 南濒杭州湾, 西部与江苏、浙江两省相接。

上海市 (4) 发音 (帮助·信息); 吴语: Zaoonhe), 简称沪, 别称申, 中华人民共和国直辖市、国家中心城市。上海位于中国南北弧形海岸线中部, 长江三角洲东部, 全市最北部为处于长江入海口中的崇明岛, 东向东海, 并隔海与日本九州岛相望, 南濒杭州湾, 西部与江苏、浙江两省相接。

上海市 (4) 发音 (帮助·信息); 吴语: Zaoonhe), 简称沪, 别称申, 中华人民共和国直辖市、国家中心城市。上海位于中国南北弧形海岸线中部, 长江三角洲东部, 全市最北部为处于长江入海口中的崇明岛, 东向东海, 并隔海与日本九州岛相望, 南濒杭州湾, 西部与江苏、浙江两省相接。

上海市 (4) 发音 (帮助·信息); 吴语: Zaoonhe), 简称沪, 别称申, 中华人民共和国直辖市、国家中心城市。上海位于中国南北弧形海岸线中部, 长江三角洲东部, 全市最北部为处于长江入海口中的崇明岛, 东向东海, 并隔海与日本九州岛相望, 南濒杭州湾, 西部与江苏、浙江两省相接。

上海市 (4) 发音 (帮助·信息); 吴语: Zaoonhe), 简称沪, 别称申, 中华人民共和国直辖市、国家中心城市。上海位于中国南北弧形海岸线中部, 长江三角洲东部, 全市最北部为处于长江入海口中的崇明岛, 东向东海, 并隔海与日本九州岛相望, 南濒杭州湾, 西部与江苏、浙江两省相接。

上海市 (4) 发音 (帮助·信息); 吴语: Zaoonhe), 简称沪, 别称申, 中华人民共和国直辖市、国家中心城市。上海位于中国南北弧形海岸线中部, 长江三角洲东部, 全市最北部为处于长江入海口中的崇明岛, 东向东海, 并隔海与日本九州岛相望, 南濒杭州湾, 西部与江苏、浙江两省相接。

上海市 (4) 发音 (帮助·信息); 吴语: Zaoonhe), 简称沪, 别称申, 中华人民共和国直辖市、国家中心城市。上海位于中国南北弧形海岸线中部, 长江三角洲东部, 全市最北部为处于长江入海口中的崇明岛, 东向东海, 并隔海与日本九州岛相望, 南濒杭州湾, 西部与江苏、浙江两省相接。

上海市 (4) 发音 (帮助·信息); 吴语: Zaoonhe), 简称沪, 别称申, 中华人民共和国直辖市、国家中心城市。上海位于中国南北弧形海岸线中部, 长江三角洲东部, 全市最北部为处于长江入海口中的崇明岛, 东向东海, 并隔海与日本九州岛相望, 南濒杭州湾, 西部与江苏、浙江两省相接。

上海市 (4) 发音 (帮助·信息); 吴语: Zaoonhe), 简称沪, 别称申, 中华人民共和国直辖市、国家中心城市。上海位于中国南北弧形海岸线中部, 长江三角洲东部, 全市最北部为处于长江入海口中的崇明岛, 东向东海, 并隔海与日本九州岛相望, 南濒杭州湾, 西部与江苏、浙江两省相接。

上海市

上海市 (4) 发音 (帮助·信息); 吴语: Zaoonhe), 简称沪, 别称申, 中华人民共和国直辖市、国家中心城市。上海位于中国南北弧形海岸线中部, 长江三角洲东部, 全市最北部为处于长江入海口中的崇明岛, 东向东海, 并隔海与日本九州岛相望, 南濒杭州湾, 西部与江苏、浙江两省相接。

上海市 (4) 发音 (帮助·信息); 吴语: Zaoonhe), 简称沪, 别称申, 中华人民共和国直辖市、国家中心城市。上海位于中国南北弧形海岸线中部, 长江三角洲东部, 全市最北部为处于长江入海口中的崇明岛, 东向东海, 并隔海与日本九州岛相望, 南濒杭州湾, 西部与江苏、浙江两省相接。

上海市 (4) 发音 (帮助·信息); 吴语: Zaoonhe), 简称沪, 别称申, 中华人民共和国直辖市、国家中心城市。上海位于中国南北弧形海岸线中部, 长江三角洲东部, 全市最北部为处于长江入海口中的崇明岛, 东向东海, 并隔海与日本九州岛相望, 南濒杭州湾, 西部与江苏、浙江两省相接。

上海市 (4) 发音 (帮助·信息); 吴语: Zaoonhe), 简称沪, 别称申, 中华人民共和国直辖市、国家中心城市。上海位于中国南北弧形海岸线中部, 长江三角洲东部, 全市最北部为处于长江入海口中的崇明岛, 东向东海, 并隔海与日本九州岛相望, 南濒杭州湾, 西部与江苏、浙江两省相接。

上海市 (4) 发音 (帮助·信息); 吴语: Zaoonhe), 简称沪, 别称申, 中华人民共和国直辖市、国家中心城市。上海位于中国南北弧形海岸线中部, 长江三角洲东部, 全市最北部为处于长江入海口中的崇明岛, 东向东海, 并隔海与日本九州岛相望, 南濒杭州湾, 西部与江苏、浙江两省相接。

上海市 (4) 发音 (帮助·信息); 吴语: Zaoonhe), 简称沪, 别称申, 中华人民共和国直辖市、国家中心城市。上海位于中国南北弧形海岸线中部, 长江三角洲东部, 全市最北部为处于长江入海口中的崇明岛, 东向东海, 并隔海与日本九州岛相望, 南濒杭州湾, 西部与江苏、浙江两省相接。

上海市 (4) 发音 (帮助·信息); 吴语: Zaoonhe), 简称沪, 别称申, 中华人民共和国直辖市、国家中心城市。上海位于中国南北弧形海岸线中部, 长江三角洲东部, 全市最北部为处于长江入海口中的崇明岛, 东向东海, 并隔海与日本九州岛相望, 南濒杭州湾, 西部与江苏、浙江两省相接。

上海市 (4) 发音 (帮助·信息); 吴语: Zaoonhe), 简称沪, 别称申, 中华人民共和国直辖市、国家中心城市。上海位于中国南北弧形海岸线中部, 长江三角洲东部, 全市最北部为处于长江入海口中的崇明岛, 东向东海, 并隔海与日本九州岛相望, 南濒杭州湾, 西部与江苏、浙江两省相接。

上海市 (4) 发音 (帮助·信息); 吴语: Zaoonhe), 简称沪, 别称申, 中华人民共和国直辖市、国家中心城市。上海位于中国南北弧形海岸线中部, 长江三角洲东部, 全市最北部为处于长江入海口中的崇明岛, 东向东海, 并隔海与日本九州岛相望, 南濒杭州湾, 西部与江苏、浙江两省相接。

上海市 (4) 发音 (帮助·信息); 吴语: Zaoonhe), 简称沪, 别称申, 中华人民共和国直辖市、国家中心城市。上海位于中国南北弧形海岸线中部, 长江三角洲东部, 全市最北部为处于长江入海口中的崇明岛, 东向东海, 并隔海与日本九州岛相望, 南濒杭州湾, 西部与江苏、浙江两省相接。

上海市 (4) 发音 (帮助·信息); 吴语: Zaoonhe), 简称沪, 别称申, 中华人民共和国直辖市、国家中心城市。上海位于中国南北弧形海岸线中部, 长江三角洲东部, 全市最北部为处于长江入海口中的崇明岛, 东向东海, 并隔海与日本九州岛相望, 南濒杭州湾, 西部与江苏、浙江两省相接。

上海市 (4) 发音 (帮助·信息); 吴语: Zaoonhe), 简称沪, 别称申, 中华人民共和国直辖市、国家中心城市。上海位于中国南北弧形海岸线中部, 长江三角洲东部, 全市最北部为处于长江入海口中的崇明岛, 东向东海, 并隔海与日本九州岛相望, 南濒杭州湾, 西部与江苏、浙江两省相接。

上海市 (4) 发音 (帮助·信息); 吴语: Zaoonhe), 简称沪, 别称申, 中华人民共和国直辖市、国家中心城市。上海位于中国南北弧形海岸线中部, 长江三角洲东部, 全市最北部为处于长江入海口中的崇明岛, 东向东海, 并隔海与日本九州岛相望, 南濒杭州湾, 西部与江苏、浙江两省相接。

上海市 (4) 发音 (帮助·信息); 吴语: Zaoonhe), 简称沪, 别称申, 中华人民共和国直辖市、国家中心城市。上海位于中国南北弧形海岸线中部, 长江三角洲东部, 全市最北部为处于长江入海口中的崇明岛, 东向东海, 并隔海与日本九州岛相望, 南濒杭州湾, 西部与江苏、浙江两省相接。

Abstracts

上海市 (4) 发音 (帮助·信息); 吴语: Zaoonhe), 简称沪, 别称申, 中华人民共和国直辖市、国家中心城市。上海位于中国南北弧形海岸线中部, 长江三角洲东部, 全市最北部为处于长江入海口中的崇明岛, 东向东海, 并隔海与日本九州岛相望, 南濒杭州湾, 西部与江苏、浙江两省相接。

上海市 (4) 发音 (帮助·信息); 吴语: Zaoonhe), 简称沪, 别称申, 中华人民共和国直辖市、国家中心城市。上海位于中国南北弧形海岸线中部, 长江三角洲东部, 全市最北部为处于长江入海口中的崇明岛, 东向东海, 并隔海与日本九州岛相望, 南濒杭州湾, 西部与江苏、浙江两省相接。

上海市 (4) 发音 (帮助·信息); 吴语: Zaoonhe), 简称沪, 别称申, 中华人民共和国直辖市、国家中心城市。上海位于中国南北弧形海岸线中部, 长江三角洲东部, 全市最北部为处于长江入海口中的崇明岛, 东向东海, 并隔海与日本九州岛相望, 南濒杭州湾, 西部与江苏、浙江两省相接。



上海市市花白玉兰

上海市 (4) 发音 (帮助·信息); 吴语: Zaoonhe), 简称沪, 别称申, 中华人民共和国直辖市、国家中心城市。上海位于中国南北弧形海岸线中部, 长江三角洲东部, 全市最北部为处于长江入海口中的崇明岛, 东向东海, 并隔海与日本九州岛相望, 南濒杭州湾, 西部与江苏、浙江两省相接。

上海市 (4) 发音 (帮助·信息); 吴语: Zaoonhe), 简称沪, 别称申, 中华人民共和国直辖市、国家中心城市。上海位于中国南北弧形海岸线中部, 长江三角洲东部, 全市最北部为处于长江入海口中的崇明岛, 东向东海, 并隔海与日本九州岛相望, 南濒杭州湾, 西部与江苏、浙江两省相接。

上海市 (4) 发音 (帮助·信息); 吴语: Zaoonhe), 简称沪, 别称申, 中华人民共和国直辖市、国家中心城市。上海位于中国南北弧形海岸线中部, 长江三角洲东部, 全市最北部为处于长江入海口中的崇明岛, 东向东海, 并隔海与日本九州岛相望, 南濒杭州湾, 西部与江苏、浙江两省相接。

Semantic Data Extraction

- Wikipedia uses the wikitext language, a lightweight markup language
 - database backup dumps
- Baidu Baike and Hudong Baike provide the WYSIWYG (what you see is what you get) HTML editors
 - HTML file archives
- 12 types of article content are considered
 - abstracts, aliases, categories, disambiguation, external links, images, infobox properties, internal links (pagelinks), labels, redirects, related pages and resource ids.
- Totally, **124,573,857** RDF triples are extracted

Overall Statistics on Extraction Results

Items	Baidu Baike	Hudong Baike	Chinese Wikipedia
Resources	3,234,950	2,765,833	559,402
~ that have abstracts	393,094 12.2%	469,009 17.0%	324,627 58.0%
~ that have categories	2,396,570 74.1%	912,627 33.0%	314,354 56.2%
~ that have infoboxes	56,762 1.8%	197,224 7.1%	24,398 4.4%
Categories	516,309	38,446	93,191
Properties	13,226	474	2,304
	per res.	per res.	per res.
Article Categories	6,774,442 2.09	2,067,349 0.75	796,679 1.42
External Links	2,529,364 0.78	827,145 0.30	573,066 1.02
Images	2,593,856 0.80	1,765,592 0.64	221,171 0.40
Infobox Properties	477,957 0.14	1,908,368 0.69	120,509 0.22
Internal Links	15,462,699 4.78	19,141,664 6.92	9,359,108 16.73
Related Pages	2,397,416 0.74	17,986,888 6.50	— —
Aliases	—	—	362,495
Disambiguation Links	28,937	13,733	40,015
Redirects	97,680	37,040	190,714

Overall Statistics on Extraction Results

Items	Baidu Baike		Hudong Baike		Chinese Wikipedia	
Resources	3,234,950		2,765,833		559,402	
~ that have abstracts	393,094	12.2%	469,009	17.0%	324,627	58.0%
~ that have categories	2,396,570	74.1%	912,627	33.0%	314,354	56.2%
~ that have infoboxes	56,762	1.8%	197,224	7.1%	24,398	4.4%
Categories	516,309		38,446		93,191	
Properties	13,226		474		2,304	
		per res.		per res.		per res.
Article Categories	6,774,442	2.09	2,067,349	0.75	796,679	1.42
External Links	2,529,364	0.78	827,145	0.30	573,066	1.02
Images	2,593,856	0.80	1,765,592	0.64	221,171	0.40
Infobox Properties	477,957	0.14	1,908,368	0.69	120,509	0.22
Internal Links	15,462,699	4.78	19,141,664	6.92	9,359,108	16.73
Related Pages	2,397,416	0.74	17,986,888	6.50	—	—
Aliases	—		—		362,495	
Disambiguation Links	28,937		13,733		40,015	
Redirects	97,680		37,040		190,714	

- 1. Linked Data
- The Chinese characters are non-ASCII, so we choose IRIs
- IRI Patterns
 - `http://zhishi.me/[DataSource]/resource/[Label]`
 - e.g. `http://zhishi.me/hudongbaike/resource/北京`
- IRIs are incompatible with HTML4 [3], we have to encode non-ASCII characters in some cases
 - e.g.
`http://zhishi.me/hudongbaike/resource/%E5%8C%97%E4%BA%AC`

Label

hudong:北京

foaf:page <http://www.hudong.com/wiki/%E5%8C%97%E4%B8%A4%E5%8C%97%E4%BA%AC>
 zhishi:resourceID 757341

- Index
- [zhishi:abstract](#)
 - [infobox](#)
 - [dcterm:subject](#)
 - [zhishi:thumbnail](#)
 - [zhishi:relatedPage](#)
 - [zhishi:externalLink](#)
 - [zhishi:internalLink](#)

owl:sameAs

- [hudong:北京 \(this\)](#)
- [baidu:北京](#)
- [zhwiki:北京](#)

[MERGE PAGE](#)

Abstract

zhishi:abstract

北京是中华人民共和国的首都，四个直辖市之一，是中国的政治、文化、科技、教育和国际交往中心，也是中国最大的航空交通枢纽和中国北方的经济中心。北京市位于华北平原北端，东南局部地区与天津市相连，其余为河北省所环绕。以市区人口数相比，北京为次于上海的中国第二大城市。北京是一座有二千余年历史的文化名城，历史上共有五个朝代曾在此定都，北京荟萃了自元明清以来的中华文化，拥有诸多名胜古迹和人文景观。北京于2008年成功举办第29届夏季奥林匹克运动会（2008年8月8日至8月24日）及残疾人奥林匹克运动会（2008年9月6日至9月17日）。

>>TOP

Infobox

infobox

所属地区	中国
GDP	11,865.9亿元人民币(2009年)
主要街道	长安街、平安大街
面积	16800平方公里
中文名	北京
下辖地区	海淀区 西城区 顺义区 大兴区
友好城市	纽约、东京、巴黎
知名企业	首钢集团、北京科技园置地有限公司、北京京宝公司、北京珠联公司
特产	烤鸭、豌豆黄、艾窝窝、豆汁
人口	1,755万人(2009年)
电话区码	010
英文名	Municipality of Beijing

>>TOP

Subjects

dcterm:subject

水游地名 文化 首都 各国首都 植物 历史 互联网 地理 中国 中国城市

>>TOP

Thumbnails

zhishi:thumbnail



>>TOP

Other properties

zhishi:relatedPage

圆明园 [SHOW MORE \(10\)](#)

>>TOP

zhishi:externalLink

<http://www.cq6.cc/index.php%3Fdoc-innerlink-%25E5%258C%2597%25E4%25BA%25AC>
[SHOW MORE \(14\)](#)

>>TOP

zhishi:internalLink

福建省 [SHOW MORE \(798\)](#)

>>TOP



■ 2. Lookup Service (<http://zhishi.me/lookup/>)

太平洋

baidu:太平洋
 hudong:太平洋
 hudong:《太平洋》
 hudong:太平洋[海洋]
 hudong:太平洋[股票]
 zhwiki:太平洋

■ 3. SPARQL Endpoint (<http://zhishi.me/sparql/>)

- AllegroGraph RDFStore is used to store the extracted triples and provide querying capabilities.

hudong:北京

foaf:page http://www.hudong.com/wiki/%E5%8C%97%E4%...J%E5%8C%97%E4%BA%AC
zhishi:resourceID 757341

- Index
- zhishi:abstract
- infobox
- dcterms:subject
- zhishi:thumbnail
- zhishi:relatedPage
- zhishi:externalLink
- zhishi:internalLink

zhishi:abstract

北京是中华人民共和国的首都，四个直辖市之一，是中国的政治、文化、科技、教育和国际交往中心，也是中国最大的陆空交通枢纽和中国北方的经济中心。北京市位于华北平原北端，东南局部地区与天津市相连，其余为河北省所环绕。以市区人口数相比，北京为次于上海的中国第二大城市。北京是一座有三千余年历史的文化名城，历史上共有五个朝代曾在此定都，北京荟萃了自元明清以来的中华文化，拥有诸多名胜古迹和人文景观。北京于2008年成功举办了第29届夏季奥林匹克运动会（2008年8月8日至8月24日）及残疾人奥林匹克运动会（2008年9月6日至9月17日）。

>>TOP

infobox

所属地区	中国
GDP	11,865.9亿元人民币(2009年)
主要街道	长安街、平安大街
面积	16800平方公里
中文名	北京
下辖地区	海淀区 西城区 顺义区 大兴区
友好城市	纽约、东京、巴黎
知名企业	首钢集团、北京科技园置地有限公司、北京京宝公司、北京珠联公司
特产	烤鸭、豌豆黄、艾窝窝、豆汁
人口	1,755万人(2009年)
电话区码	010
英文名	Municipality of Beijing

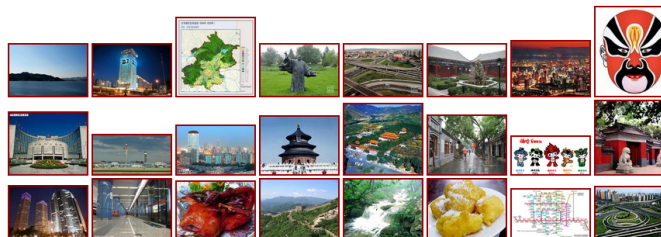
>>TOP

dcterms:subject

水游地名 文化 首都 各国首都 植物 历史 互联网 地理 中国 中国城市

>>TOP

zhishi:thumbnail



>>TOP

zhishi:relatedPage

圆明园 SHOW MORE (10)

>>TOP

zhishi:externalLink

http://www.cq6.cc/index.php%3Fdoc-innerlink-%25E5%258C%2597%25E4%25BA%25AC
SHOW MORE (14)

>>TOP

zhishi:internalLink

福建省 SHOW MORE (798)

>>TOP

owl:sameAs

- hudong:北京 (this)
- baidu:北京
- zhwiki:北京

[MERGE PAGE](#)

owl:sameAs

- hudong:北京 (this)
- baidu:北京
- zhwiki:北京

[MERGE PAGE](#)



Data-level Mapping among Different Datasets

- We utilize distributed MapReduce [4] framework to accomplish the large-scale data-level mapping task.
- Map phase:
 - Resources are indexed
- Reduce phase:
 - Resources with the same index term (match candidates) gather together and further comparisons are made
- The improved version of our mapping framework participates in IM@OAEI2011
 - “Zhishi.links Results for OAEI 2011” @OM2011

Data-level Mapping (con't)

- In this paper, we focused on resolving the problem occurs in Map phase (index generating step):
 - different resources have the same label
 - the same resource have different labels
- We proposed three reasonable but not complex strategies to generate the index:
 - Using Original Labels
 - Punctuation Cleaning
 - Extending Synonyms

Data-level Mapping (con't)

- Using Original Labels
- We extract different meanings of homonyms as different resources
 - e.g. “Jupiter (mythology)” != “Jupiter Island”
 - it is impossible to find two resources that have different meanings with the same label if all homonyms are recognized

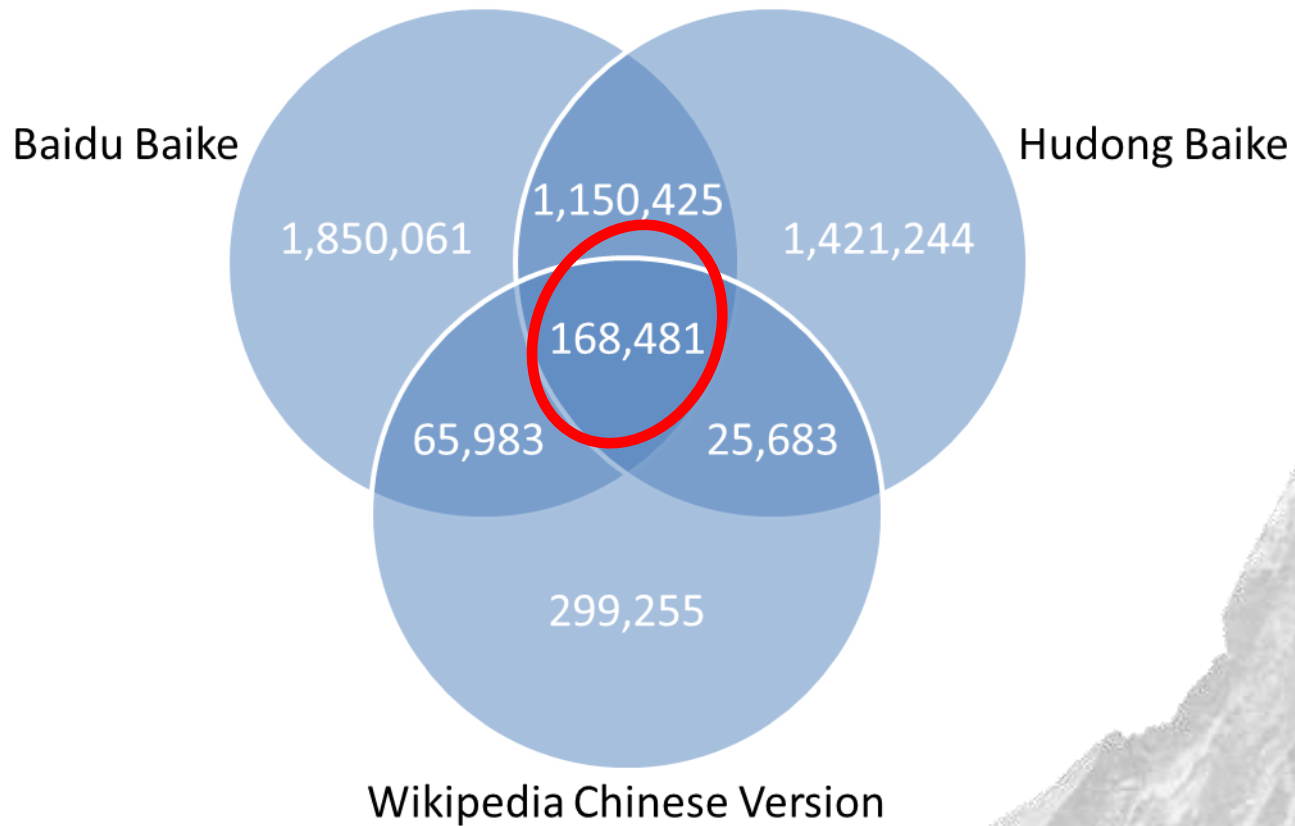
Data-level Mapping (con't)

- Punctuation Cleaning
- It is one of the most efficient methods we used to discover mappings between resources with different labels
 - Just one case of Punctuation Cleaning is given here:
 - In Chinese, we often insert an interpoint (·) or a hyphen (-) between two personal name components. In some certain cases, people may just adjoin these components.
 - e.g. a comet named after two persons
 - 海尔波普彗星 = 海尔·波普彗星 = 海尔-波普彗星

Data-level Mapping (con't)

- Extending Synonyms
- Making use of high quality synonym relations obtained from redirects information
 - A redirects to B means A and B are synonyms
 - Label(A) and Label(B) are both used as index terms for a single resource

Data-level Mapping (con't)



- Intersections of Our Three Data Sources Shown in Venn Diagram

Data-level Mapping (Web Access)

owl:sameAs

- hudong:北京 (this)
- baidu:北京
- zhwiki:北京

MERGE PAGE



- zhwiki:北京市
- hudong:北京
- baidu:北京

resources

Index

- zhishi:abstract
- infobox
- dcterms:subject
- zhishi:thumbnail
- zhishi:relatedPage
-

<owl:sameAs> check box

owl:sameAs

- zhwiki:北京市 (this)
- baidu:北京 (this)
- zhwiki:北京
- hudong:北京 (this)
- dbpedia:Beijing
- hudong:北京市
- baidu:北京市

MERGE PAGE

zhishi:abstract ■ ■ ■

北京有着三千余年的建城史和八百五十余年的建都史，最初见于记载的名字为“蓟”。民国时期，称北平。新中国成立后，是中华人民共和国的首都，简称“京”，现为中国四个中央直辖市之一，全国第二大城市及政治、交通和文化中心。北京位于华北平原北端，东南局部地区与天津相连，其余为河北省所环绕。它荟萃了元、明、清以来的中华文化，拥有众多名胜古迹和人文景观，是世界上拥有世界文化遗产最多的城市，每年有超过1亿4700万的旅客。

infobox

市花

■ 月季、菊花

a merged statement

政府驻地

■ 东城区正义路2号

下辖地区

■ 海淀区 西城区 顺义区 大兴区 ■ 东城、西城、海淀、丰台、朝阳等

时区

■ UTC+8 (东八区)

行政区类别

■ 直辖市

友好城市

■ 纽约、东京、巴黎

dcterms:subject

■ 中华人民共和国直辖市 ■ 京津冀城市群 ■ 中国城市 ■ 首都 ■ 燕京 ■ 中国

zhishi:thumbnail

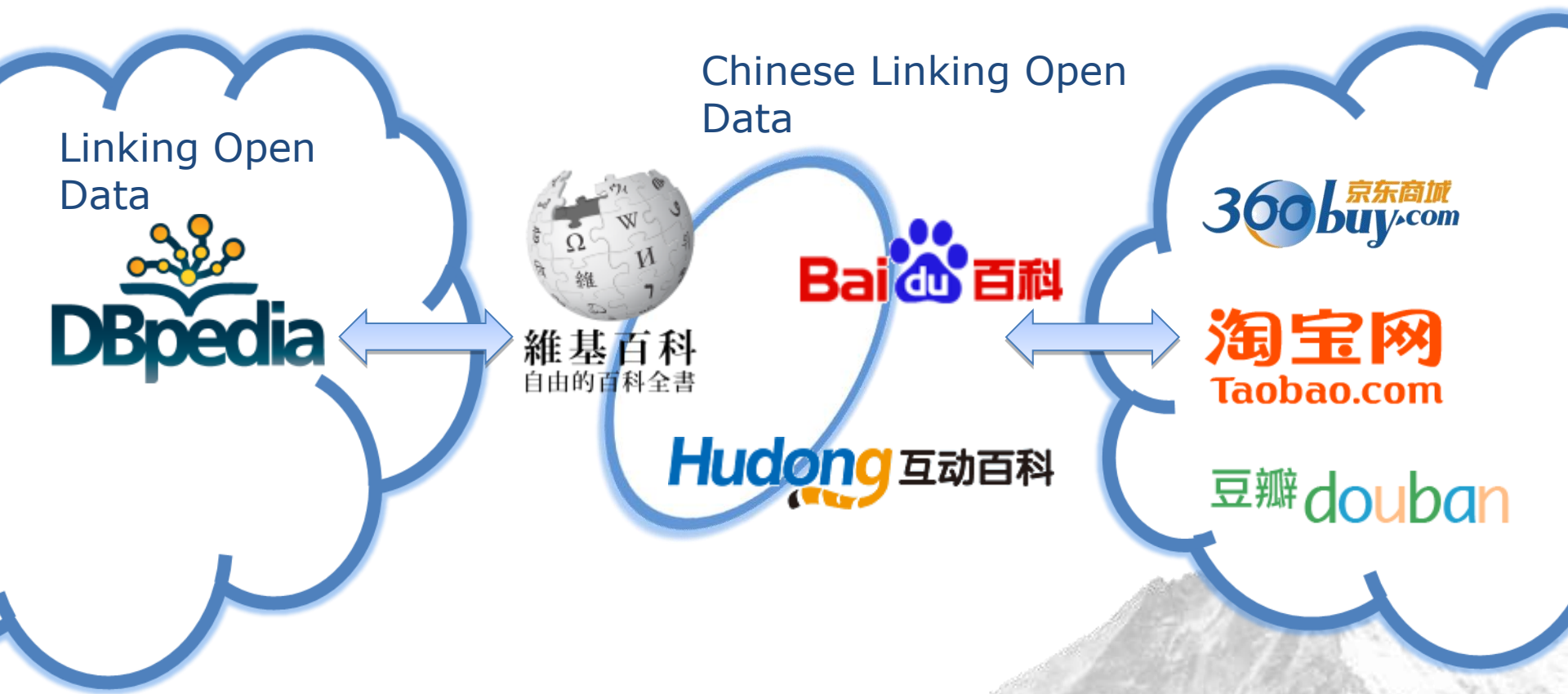


all other statements

zhishi:relatedPage

■ 圆明园 [SHOW MORE \(182\)](#)

Data-level Mapping (con't)



Conclusions and Future Work

■ Conclusions

- Zhishi.me is the first effort to build Chinese Linking Open Data
- We extracted semantic data from three Web-based free-editable encyclopedias
- Three heuristic strategies were adopted to discover `<owl:sameAs>` links between equivalent resources
- We provided Web access entries to our knowledge base for both professional and non Semantic Web community users.

■ Future Work

- Several Chinese non-encyclopedia data sources will be accommodated in our knowledge;
- We will improve instance matching strategies and provide necessary evaluations of matching quality;
- Refine extracted properties and building a general but consistent ontology automatically.

Thanks!



References

- [1] Bizer, C., Lehmann, J., Kobilarov, G., Auer, S., Becker, C., Cyganiak, R., Hellmann, S.: DBpedia - A crystallization point for the Web of Data. *J. Web Sem.* 7(3), 154–165 (2009)
- [2] de Melo, G., Weikum, G.: Towards a universal wordnet by learning from combined evidence. In: Cheung, D.W.L., Song, I.Y., Chu, W.W., Hu, X., Lin, J.J. (eds.) *CIKM*, pp. 513–522. ACM (2009)
- [3] Raggett, D., Hors, A.L., Jacobs, I.: *HTML 4.01 Specification - Appendix B: Performance, Implementation, and Design Notes*. W3C Recommendation (December 1999), <http://www.w3.org/TR/html4/appendix/notes.html>
- [4] Dean, J., Ghemawat, S.: *MapReduce: Simplified Data Processing on Large Clusters*. In: *OSDI*. pp. 137–150 (2004)