



关联数据应用：知识发现

SinoPedia

陈涛

上海图书馆

2018-05-23

提纲

1

关联数据

2

SinoPedia知识库

3

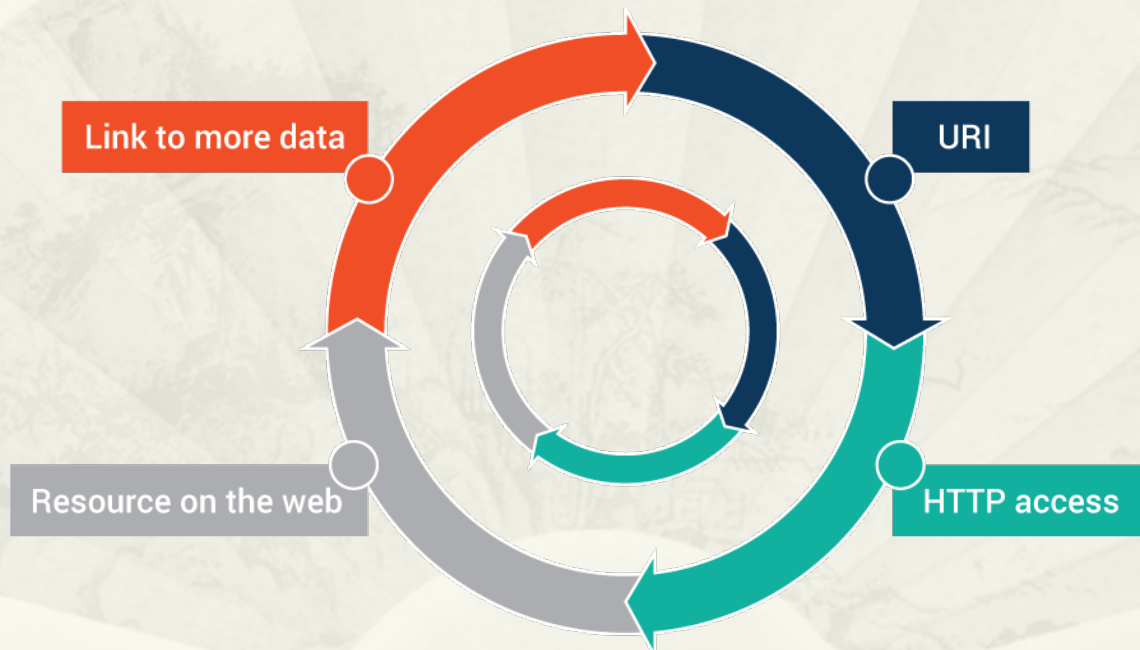
如何使用关联数据

4

知识图谱与知识发现

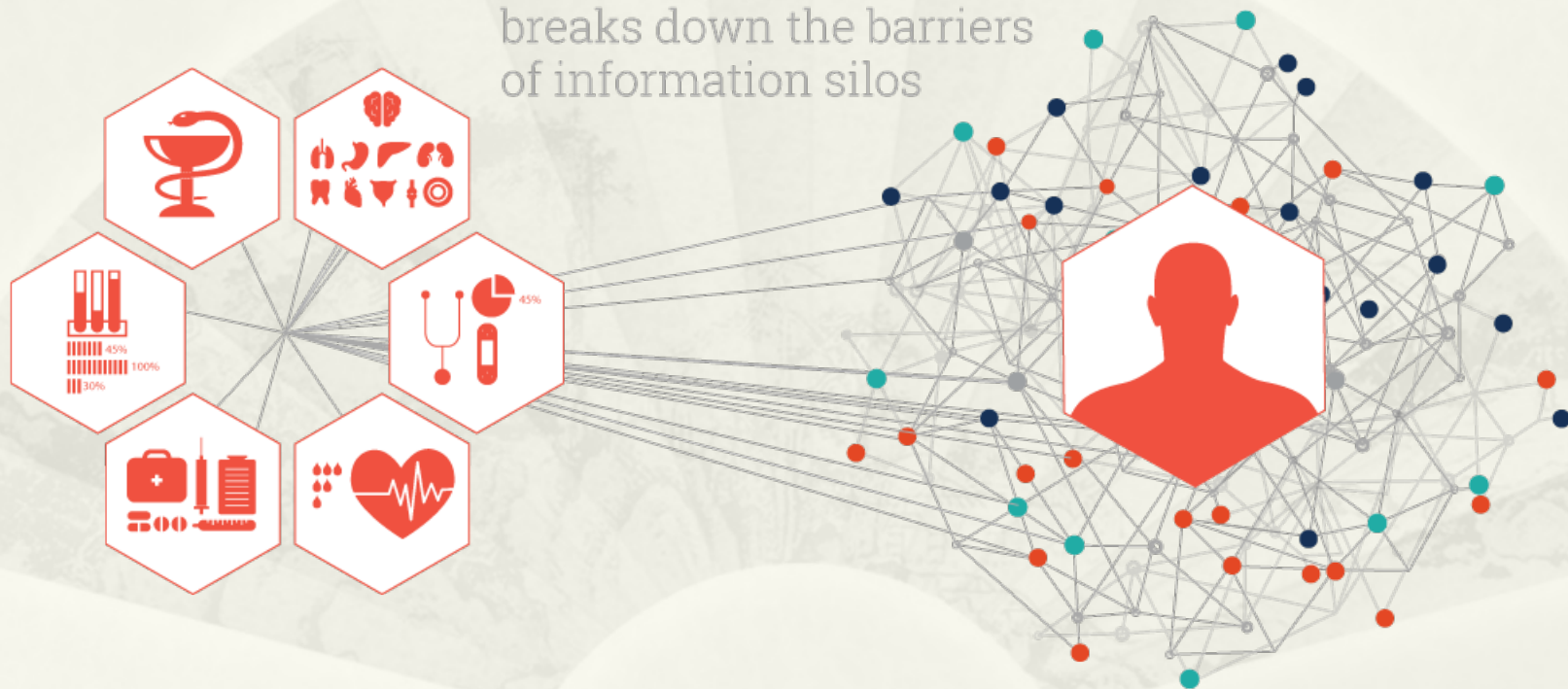
关联数据

关联数据是在网络上发布和共享数据的推荐方法

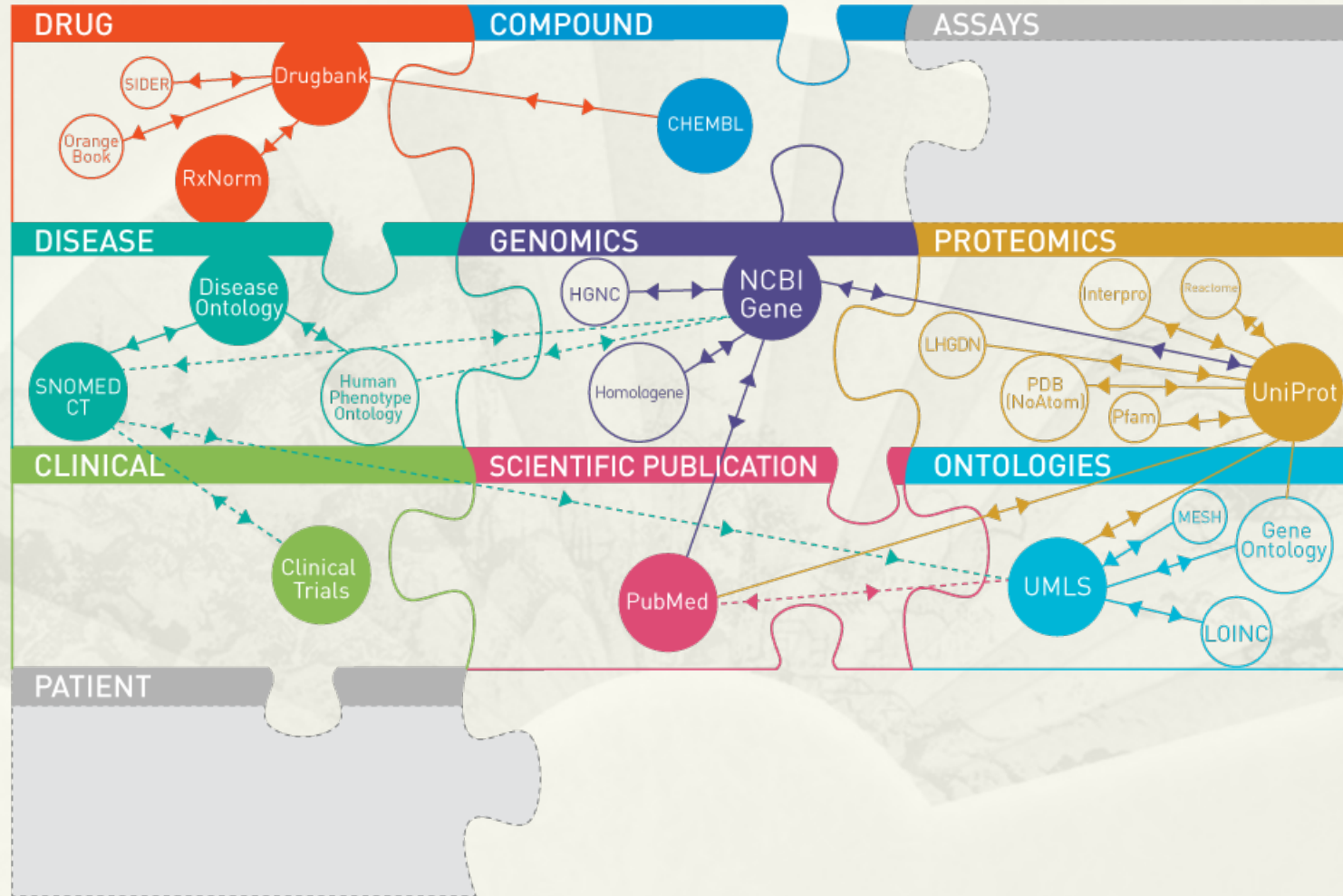


关联数据作用

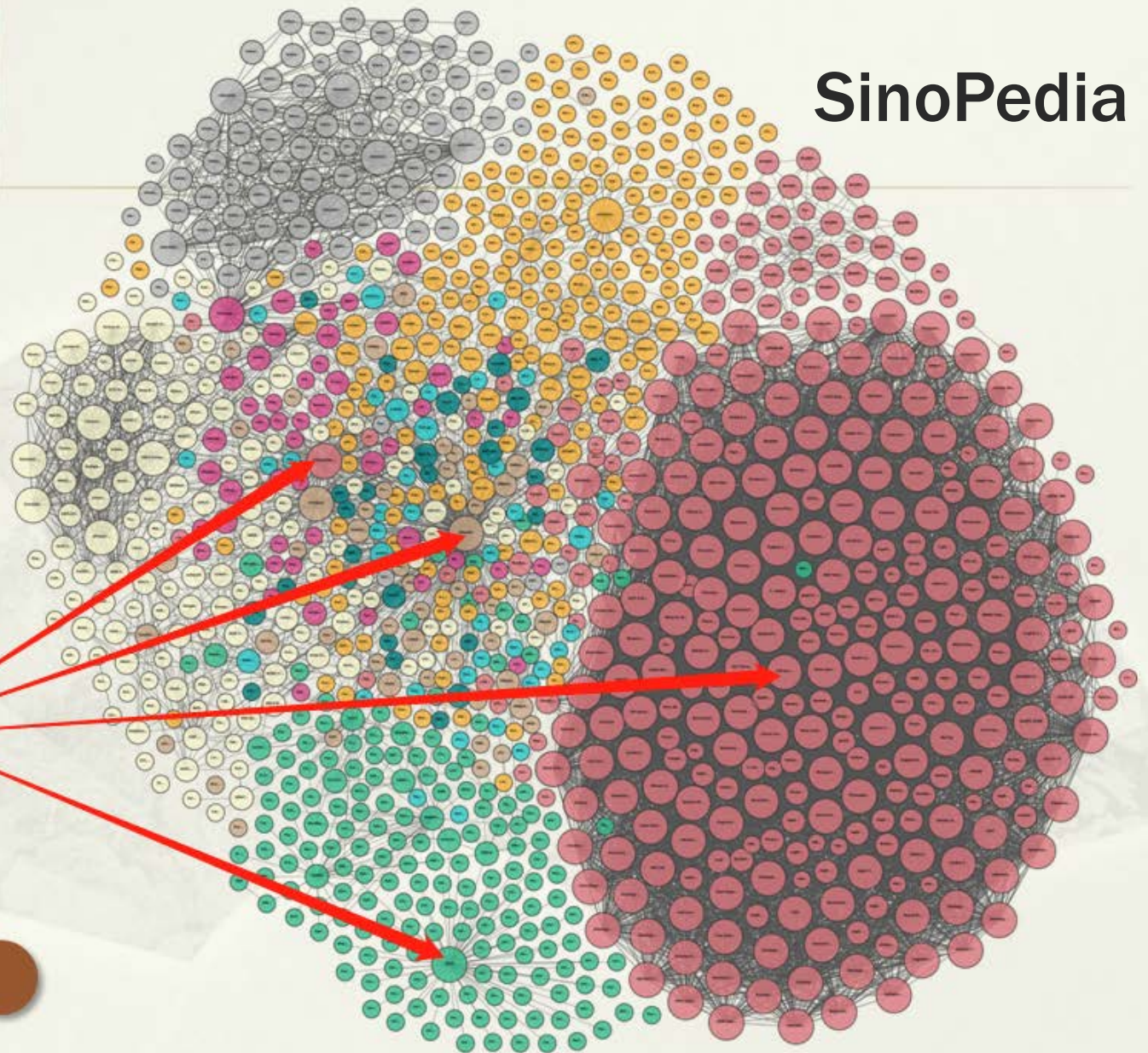
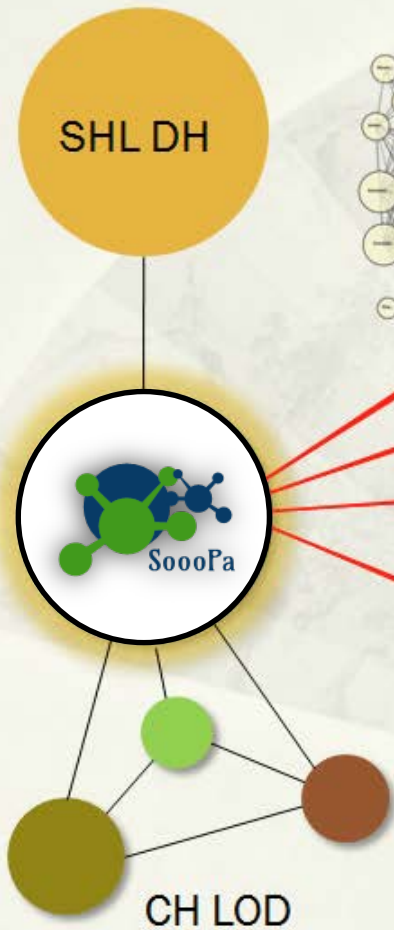
Linked Data
breaks down the barriers
of information silos



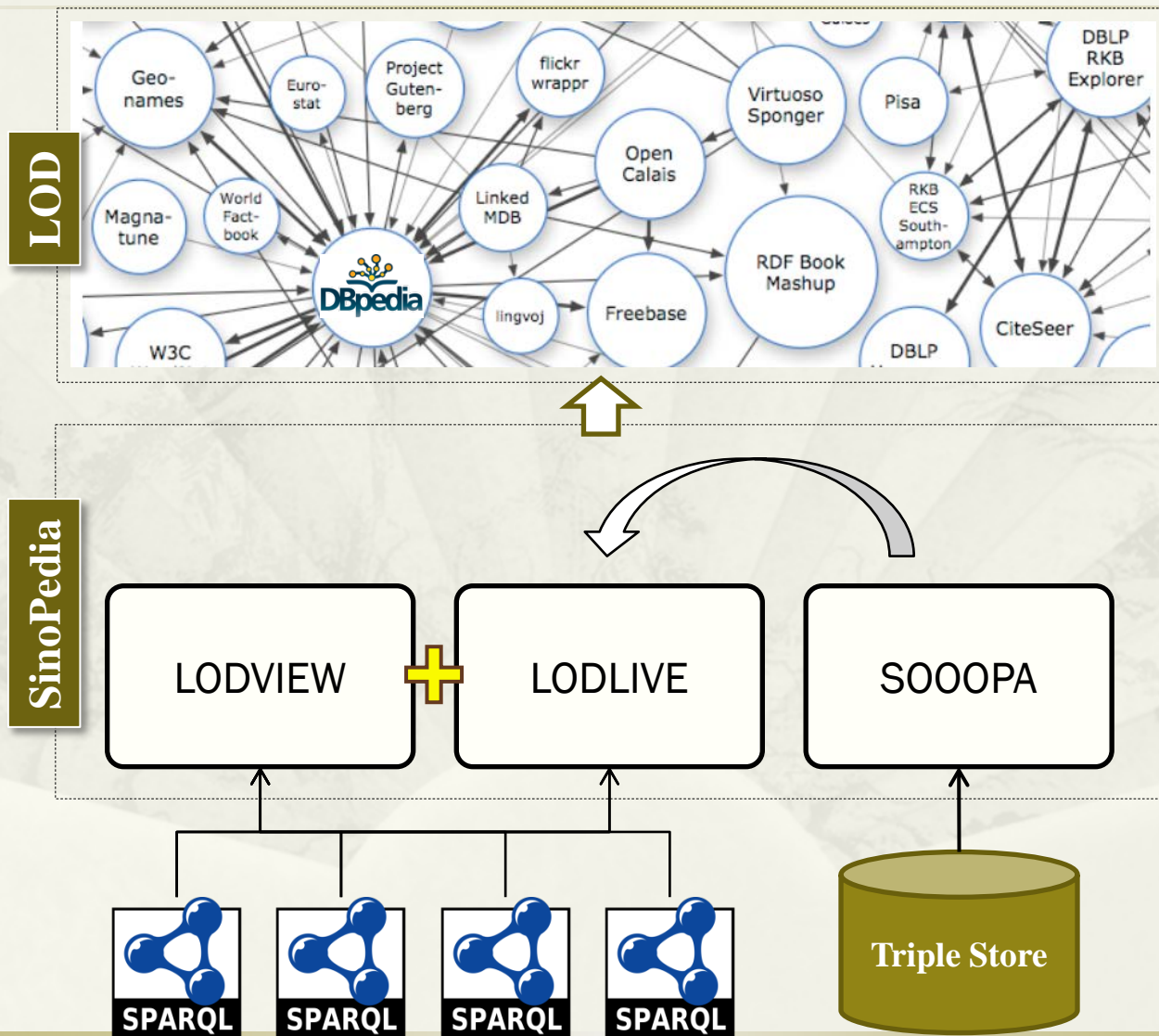
知识关联



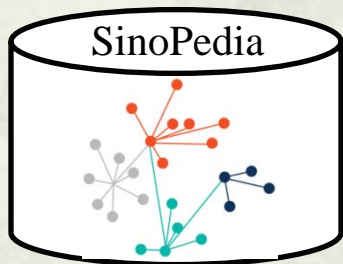
SinoPedia



SinoPedia框架



SinoPedia数据来源



数据载入

站点接入



享受关联数据发布\内容协商\知识图谱服务

开放资源

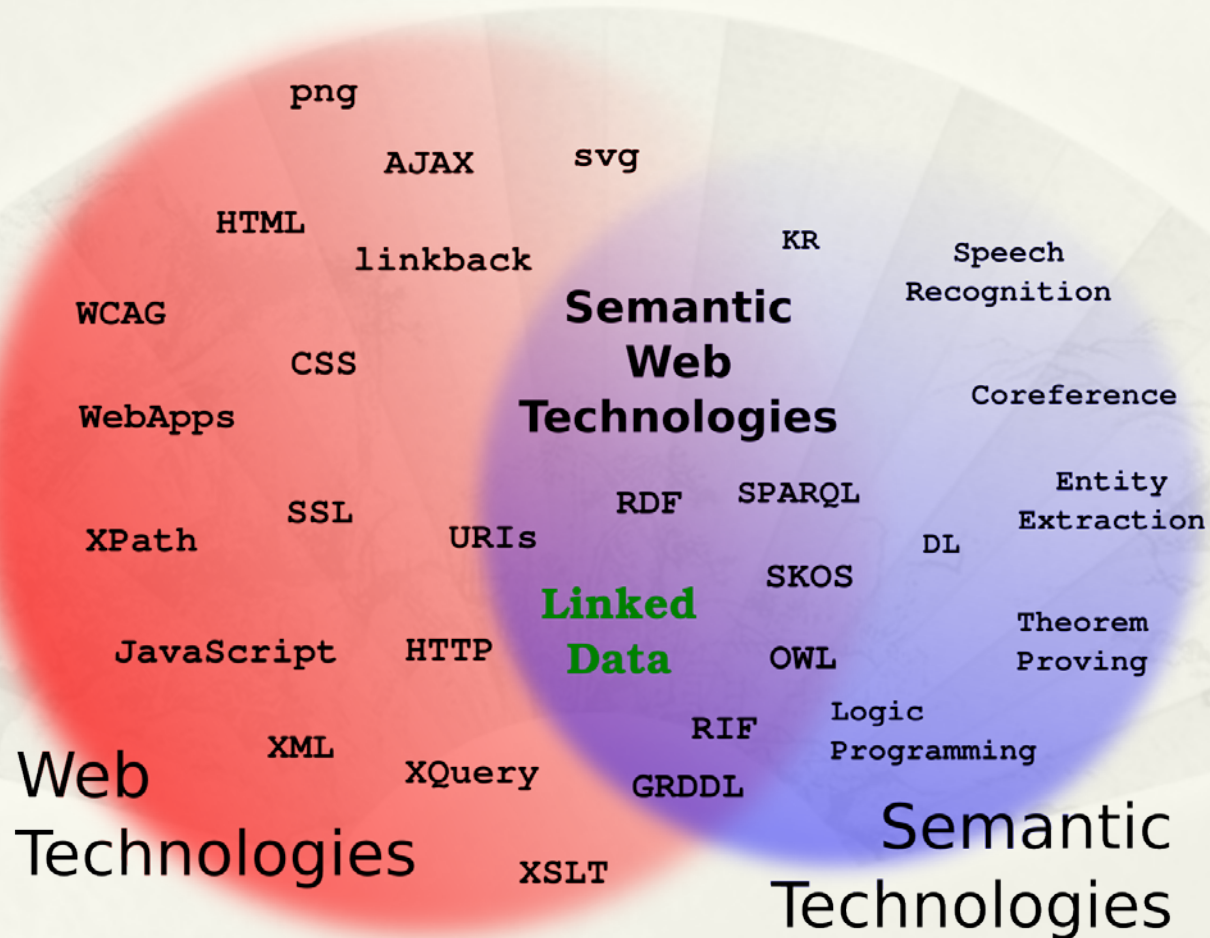
➤ LOD

- ✓ 上海图书馆数字人文系统
- ✓ DBPedia
- ✓ BBC
- ✓ LOC
- ✓ VIAF - Virtual International Authority File
- ✓ GeoNames
- ✓

➤ Non-LOD

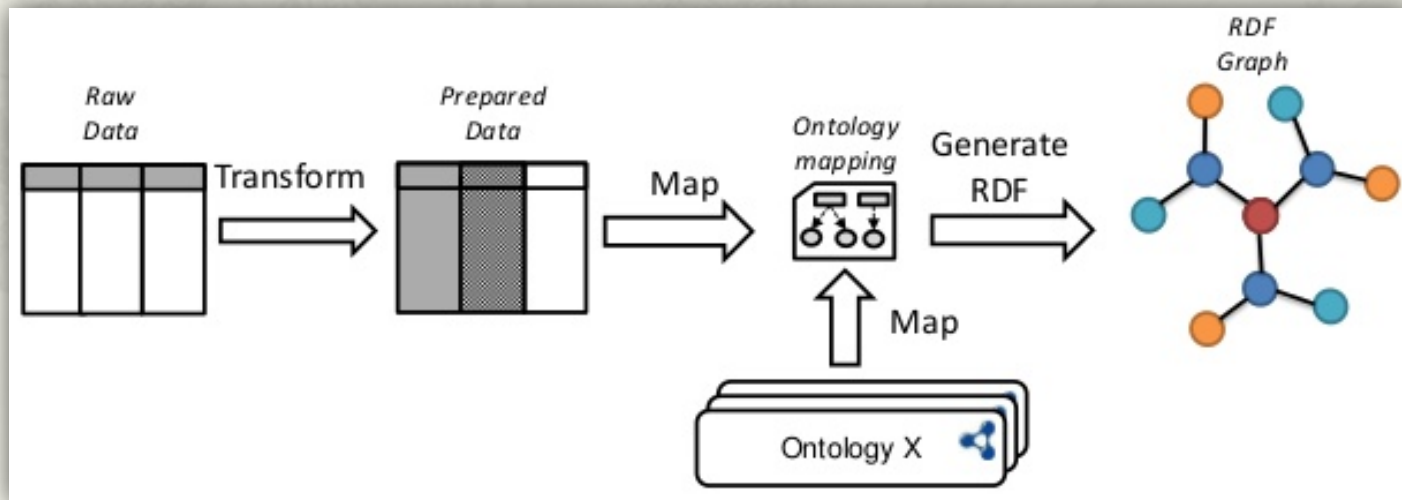
- ✓ CBDB
- ✓ 百科资源
- ✓ 图片
- ✓

相关技术



RDF数据生成

- ① 准备数据
- ② 转换/清洗数据
- ③ 准备映射文件（数据字段到本体的映射）
- ④ 生成RDF数据



Link to other Data Sets

* Linking Approaches

* Common Key Matching

- * Matching based on common keys
 - * e.g. ISBN, DOI, Wikipedia Article ID, Musicbrainz IDs
- * Matching locations based on geographic coordinates

* Label Matching

- * Comparing labels using string similarity measures
 - * e.g. object/page with title/label “The Shining (film)” on DBpedia/Wikipedia is the same as movie object with the “The Shining” on LinkedMDB
- * Comparing labels using semantic similarity measures
 - * e.g. “UofT” is the same “University of Toronto”, or a drug named “Tylenol” is the same another drug “Acetaminophen” (scientific name of brand name Tylenol)

* Graph/Ontology Matching

- * Compare labels, schema elements (e.g., types), and related objects (e.g., matching papers if they have the same set of authors)

Link to other Data Sets

* Linking Frameworks

Silk | Limes | LinQuer

The screenshot displays the Silk Workbench interface, which is used for creating and managing semantic links between data sets. The interface includes a menu bar with options like Start, Workspace, Editor, Generate Links, Learn, Reference Links, Status, and About. Below the menu bar, there are buttons for Undo, Redo, and Help, and a status bar showing Precision: 1,00 | Recall: 1,00 | F-measure: 1,00 with a green checkmark.

The main workspace is divided into several sections:

- Property Paths:** This section is on the left and contains two columns: Source and Target. The Source column has a custom path with the URI `?a/<http://xmlns.com/foaf/0.1/n`. The Target column has a custom path with the URI `?b/rdfs:label`.
- Transformations:** This section is below the Property Paths and contains two transformations: Lower case and Tokenize.
- Comparators:** This section is below the Transformations and contains three comparators: Equality, Jaccard, and another one partially visible.
- Aggregators:** This section is at the bottom left and contains two aggregators: Average and Maximum.

The main workspace shows a workflow diagram with the following steps:

- Path (Source):** A purple box with name `unnamed_1` and URI `?a/<http://xmlns.com/foaf/0.1/n`.
- Path (Target):** A red box with name `unnamed_2` and URI `?b/rdfs:label`.
- Lower case (Transfc):** A green box with name `unnamed_8` (receiving input from the Source path) and another green box with name `unnamed_9` (receiving input from the Target path).
- Levenshtein distance:** An orange box with name `unnamed_6`, threshold `0.0`, weight `1`, `minChar: 0`, and `maxChar: z`. It receives input from both `unnamed_8` and `unnamed_9`.
- Date (Compare):** An orange box with name `unnamed_5`, threshold `400.0`, and weight `1`. It receives input from `unnamed_4` and `unnamed_3`.
- Minimum (Aggregate):** A blue box with name `unnamed_7`, weight `1`. It receives input from `unnamed_6` and `unnamed_5`.

At the bottom of the interface, there are two dropdown menus: Link Limit: unlimited and Link Type: owl:sameAs.

关联属性

常用的链接属性 (Predicates)

- * owl:sameAs
- * foaf:homepage
- * foaf:topic
- * foaf:based_near
- * foaf:maker/foaf:made
- * foaf:depiction

- * foaf:page
- * foaf:primaryTopic
- * rdfs:seeAlso

SPARQL Queries

HOW TO DO?

1. Issue follow-up queries to different endpoints
2. Querying a central collection of datasets
3. Build store with copies of relevant datasets
4. Use query federation system

Follow-up Queries

Idea: issue follow-up queries over other datasets based on results from previous queries

Substituting placeholders in query templates

✓ Advantage

Queried data is up-to-date

✗ Drawbacks

Requires the existence of a SPARQL endpoint for each dataset

Requires program logic

Very inefficient

Query a Collection

Idea: Use an existing SPARQL endpoint that provides access to a set of copies of relevant datasets

Example:

SPARQL endpoint over a majority of datasets from the LOD cloud

✓ Advantage

No need for specific program logic

✗ Drawbacks

Queried data might be out of date

Not all relevant datasets in the collection

Build Store

Idea: Build your own store with copies of relevant datasets and query it

✓ Advantages

- No need for specific program logic

- Can include all datasets

- Independent of the existence, availability, and efficiency of SPARQL endpoints

✗ Drawbacks

- Requires effort to set up and to operate the store

- Ideally, data sources provide RDF dumps; if not?

- How to keep the copies in sync with the originals?

- Queried data might be out of date

Federate Query

Idea: Querying a mediator which distributes sub-queries to relevant sources and integrates the results

DARQ (Distributed ARQ) <http://darq.sourceforge.net/>

Semantic Web Integrator and Query Engine(SemWIQ) <http://semwiq.sourceforge.net/>

✓ Advantages

- No need for specific program logic

- Queried data is up to date

✗ Drawbacks

- Requires the existence of a SPARQL endpoint for each dataset

SinoPedia入口: 嗖啪

上海图书馆 SPARQL | Namespace



嗖啪

数字人文关联数据知识库

嗖啪一下

尝试: 闻一多 | 李政道 | 威廉·劳伦斯·布拉格 | 上海图书馆
<http://data.nobelprize.org/resource/laureate/21>

已接入站点

- π index评价库 (<http://bm.pi-index.com:8181/sparql>)
- 诺贝尔奖知识库 (<http://data.nobelprize.org/sparql>)
- 第一次世界大战知识库 (<http://ldf.fi/ww1lod/sparql>)
- Getty艺术知识库 (<http://vocab.getty.edu/sparql>)
- 地理信息知识库 (<http://geo.linkeddata.es/sparql>)



扫一扫，直接在手机上打开

嗖啪检索平台

嗖啪 Sopa

鲁迅

嗖啪一下

检索结果 (8) | 仅限题名

鲁迅美术学院 (RDF/XML | JSON-LD | NTriples) 机构
http://sinopedia.library.sh.cn/entity/organization/df9282ad585d42bcb3610c44dc5ffd5b
鲁迅美术学院，为中国辽宁省的一所艺术类高校。现有沈阳、大连两个校区。鲁迅艺术学院始建于延安。1945年，迁校至东北。1958年，改为鲁迅美术学院。鲁美设有11个系，分别为中国画系、版画系、油画系、雕塑系、摄影系、视觉传达艺术设计系、环境艺术设计系、染织服装艺术设计系、工业设计系、美术史论系、文化传播与管理系。
同等词条 [DBPedia]

鲁迅博物馆 (RDF/XML | JSON-LD | NTriples) 机构
http://sinopedia.library.sh.cn/entity/organization/qemavdomvzbx1cafqemavdomvzbx1caf
同等词条 [上海图书馆]

鲁迅 (RDF/XML | JSON-LD | NTriples) 人物
http://sinopedia.library.sh.cn/entity/person/dd411250a2c24fa8a2bdecc1bc819ab0
周树人（1881年9月25日—1936年10月19日），原名樟寿，字豫才，豫山、豫亭，以笔名鲁迅闻名于世，浙江绍兴人，为中国的现代著名作家，新文化运动的领导人、支持者，中国现代文学的奠基人和开山巨匠，在西方世界享有盛誉的中国现代文学家、思想家。鲁迅的主要成就包括杂文、短中篇小说、文学、思想和社会评论、学术著作、自然科学著作、古代典籍校勘与研究、散文、现代散文诗、旧体诗、外国文学与学术翻译作品和木刻版画的研究，对于五四运动以后的中国社会思想文化发展产生一定的影响，蜚声世界文坛，尤其在韩国、日本思想文化领域有极其重要的地位和影响，被誉为“二十世纪东亚文化地图上占最大领土”。

RDF/XML

JSON-LD

NTriples

内容协商

- <rdf:RDF>

- <rdf:Description rdf:about="http://sinopedia.library.sh.cn/entity/person/dd411250a2c24fa8a2bdecc1bc819ab0">

<owl:sameAs rdf:resource="http://dbpedia.org/resource/Lu_Xun"/>

<foaf:depiction rdf:resource="http://commons.wikimedia.org/wiki/Special:FilePath/LuXun1930.jpg"/>

<owl:sameAs rdf:resource="http://data.library.sh.cn/entity/person/x5xiy94zid7u351t"/>

- <rdf:comment xml:lang="zh">

原名周树人，字豫才。浙江绍兴人。现代文学家、思想家。1898年就读于江南水师学堂。1902年赴日本留学，原学医，后转志于文学。

1905年参加革命派与改良派的大论战，力主反清革命。1907年加入光复会。1909年回国后，曾在杭州、绍兴等地任教。辛亥革命后任职于中华民国政府教育部社会教育司，曾参与筹建京师图书馆和历史博物馆，同时在北京大学和北京女子师范大学等校兼课。后投身于五四新文化运动。1918年1月参加《新青年》杂志编委会。同年5月发表中国现代文学史上第一篇白话文小说《狂人日记》，奠定了新文学运动的基石。此后，又陆续创作发表了《呐喊》、《坟》、《热风》、《彷徨》等小说专集，并开始接触马克思主义。1925年发起并领导语丝社、未名社，主编《莽原》杂志。1926年因支持北京爱国学生运动而遭军阀政府通缉，同年8月南下，执教于厦门大学。1927年1月转赴广州，任中山大学文学系主任兼教务主任，并与中国共产党组织建立联系。四一五反革命政变发生后愤而辞职。同年10月迁居上海，同时彻底放弃进化论思想，接受马克思主义观点。1928年创办《奔流》杂志，并编辑

<http://sinopedia.library.sh.cn/entity/person/dd411250a2c24fa8a2bdecc1bc819ab0> <http://www.w3.org/2002/07/owl#sameAs> <http://dbpedia.org/resource/Lu_Xun> .

<http://sinopedia.library.sh.cn/entity/person/dd411250a2c24fa8a2bdecc1bc819ab0> <http://xmlns.com/foaf/0.1/depiction> <http://commons.wikimedia.org/wiki/Special:FilePath/LuXun1930.jpg> .

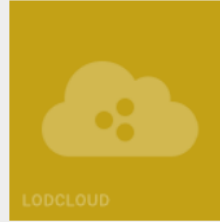
<http://sinopedia.library.sh.cn/entity/person/dd411250a2c24fa8a2bdecc1bc819ab0> <http://www.w3.org/2002/07/owl#sameAs> <http://data.library.sh.cn/entity/person/x5xiy94zid7u351t> .

<http://sinopedia.library.sh.cn/entity/person/dd411250a2c24fa8a2bdecc1bc819ab0> <http://www.w3.org/2000/01/rdf-schema#comment> "原名周树人，字豫才。浙江绍兴人。现代文学家、思想家。1898年就读于江南水师学堂。1902年赴日本留学，原学医，后转志于文学。1905年参加革命派与改良派的大论战，力主反清革命。1907年加入光复会。1909年回国后，曾在杭州、绍兴等地任教。辛亥革命后任职于中华民国政府教育部社会教育司，曾参与筹建京师图书馆和历史博物馆，同时在北京大学和北京女子师范大学等校兼课。后投身于五四新文化运动。1918年1月参加《新青年》杂志编委会。同年5月发表中国现代文学史上第一篇白话文小说《狂人日记》，奠定了新文学运动的基石。此后，又陆续创作发表了《呐喊》、《坟》、《热风》、《彷徨》等小说专集，并开始接触马克思主义。1925年发起并领导语丝社、未名社，主编《莽原》杂志。1926年因支持北京爱国学生运动而遭军阀政府通缉，同年8月南下，执教于厦门大学。1927年1月转赴广州，任中山大学文学系主任兼教务主任，并与中国共产党组织建立联系。四一五反革命政变发生后愤而辞职。同年10月迁居上海，同时彻底放弃进化论思想，接受马克思主义观点。1928年创办《奔流》杂志，并编辑出版《马克思主义文艺论丛》。1930年起，先后参加和参与发起中国自由运动大同盟、中国左翼作家联盟、中国民权保障同盟等，曾任“左联”常务委员，民权保障同盟上海分会执行委员等，同时继续进行小说、杂文及其他文学创作，并与国民党

RDF/XML

NTriples

链接数据



rdfs:label	鲁迅
foaf:gender	男 @zh
dbpedia-owl:name	鲁迅
dbpedia-owl:birthDate	1881-09-25
dbpedia-owl:deathDate	1936-10-19
rdf:type	foaf:Person
owl:sameAs	http://data.library.sh.cn/entity/person/x5xiy94zid7u351t

	dbpedia:Lu_Xun DBpedia
foaf:depiction	http://commons.wikimedia.org/wiki/Special:FilePath/LuXun1930.jpg

上图人名规范档

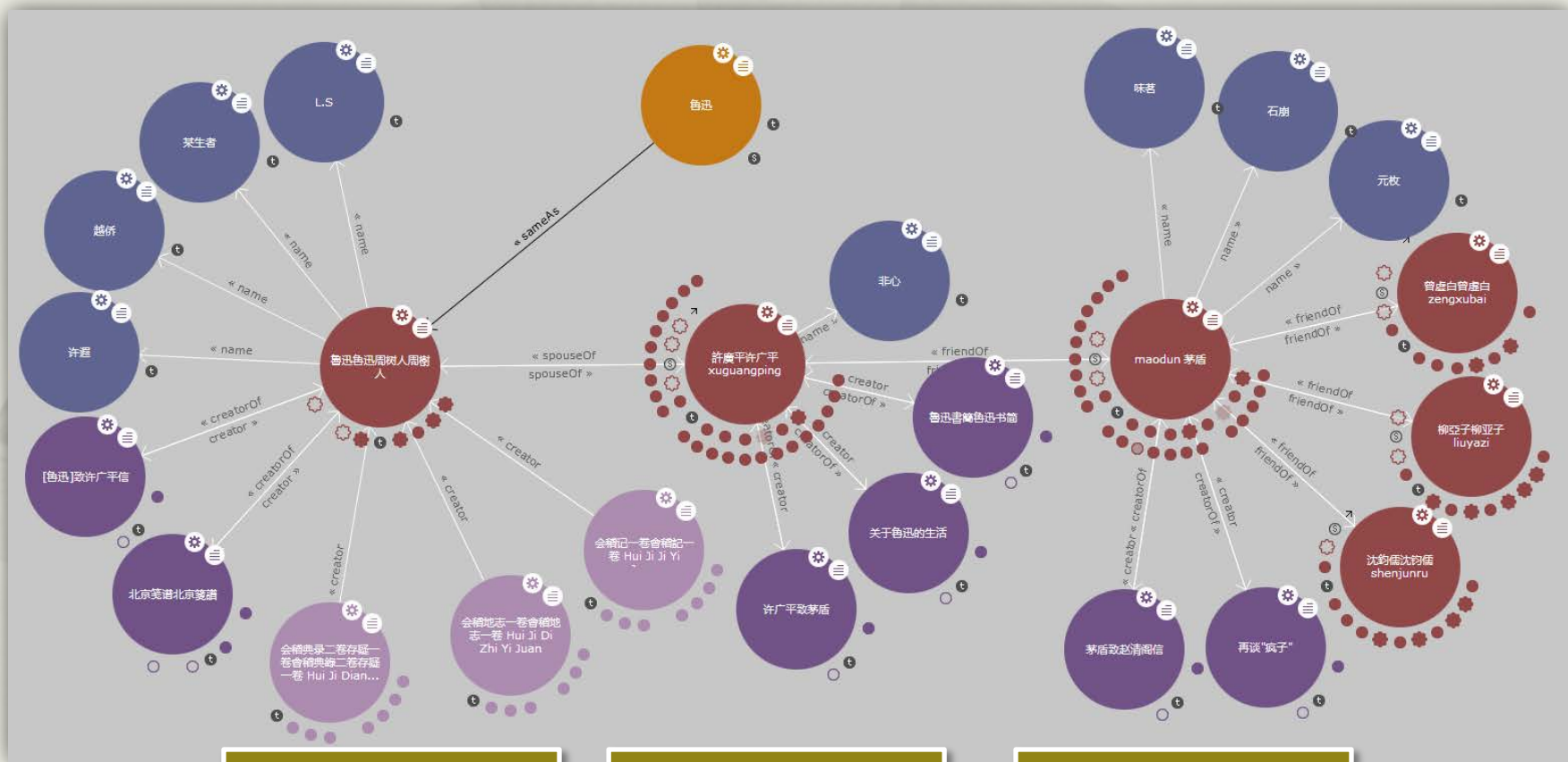
知识图谱：鲁迅

人名规范档

手稿

古籍

SinoPedia

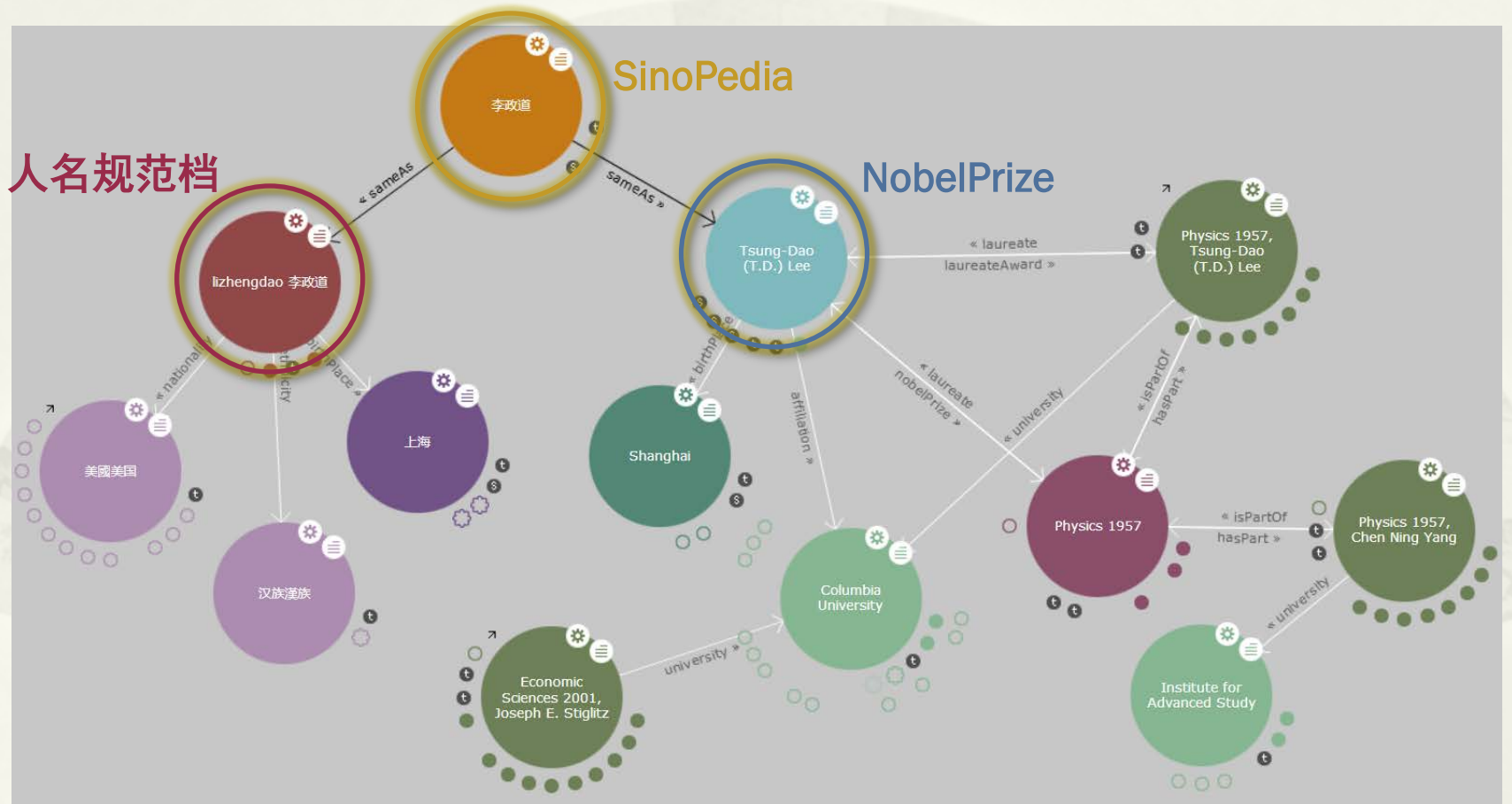


鲁迅

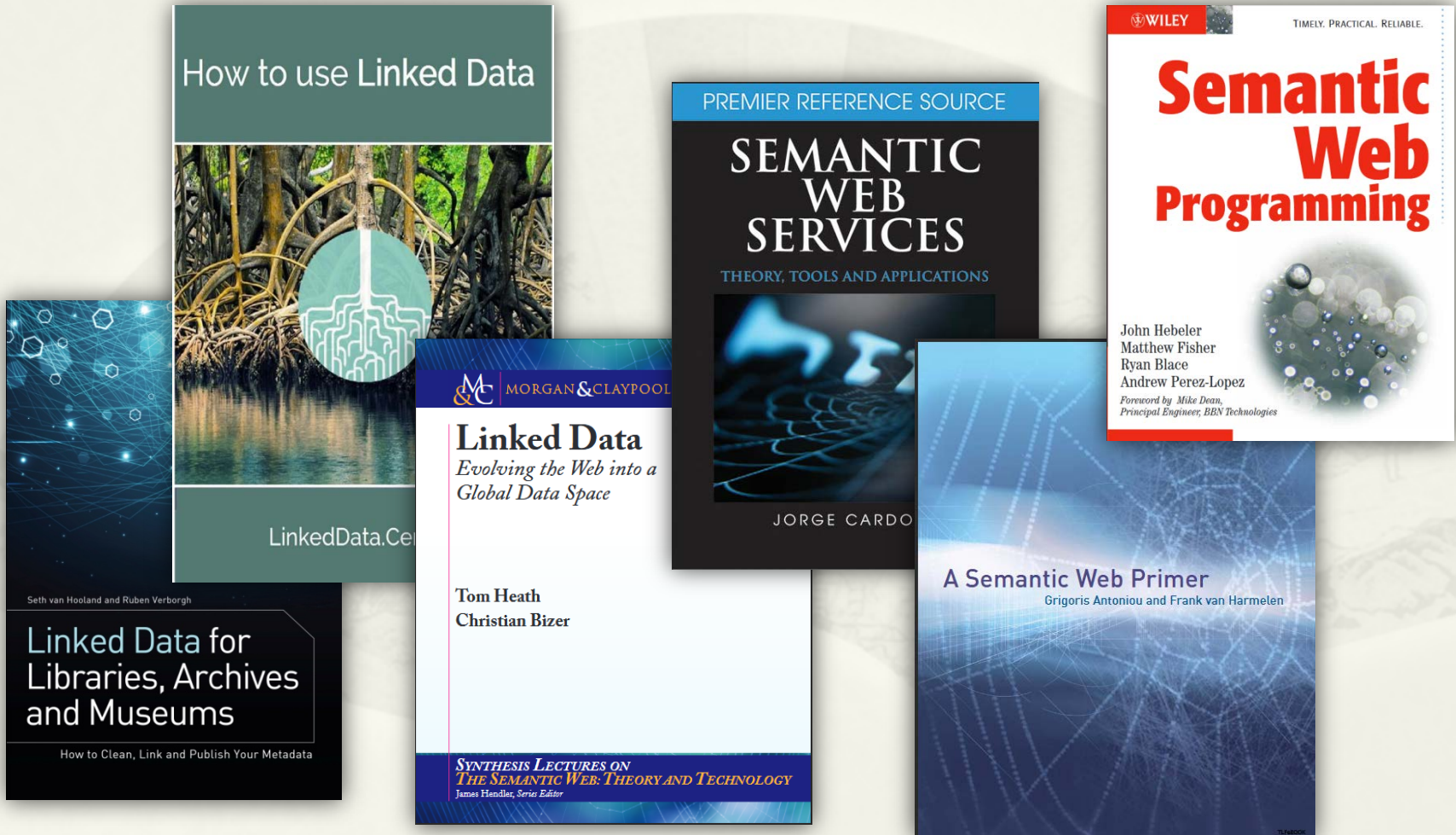
许广平

茅盾

知识图谱：李政道



好书推荐



嗖啪 <http://sinopedia.library.sh.cn/soopa>



SinoPedia
Share it now

共建 共享
共链 共赢



三人行（语义有你）
扫一扫二维码，加入该群。