# Dublin Core™ Metadata Initiative

# The Integrated Workfolw
# of Entity Management and Service
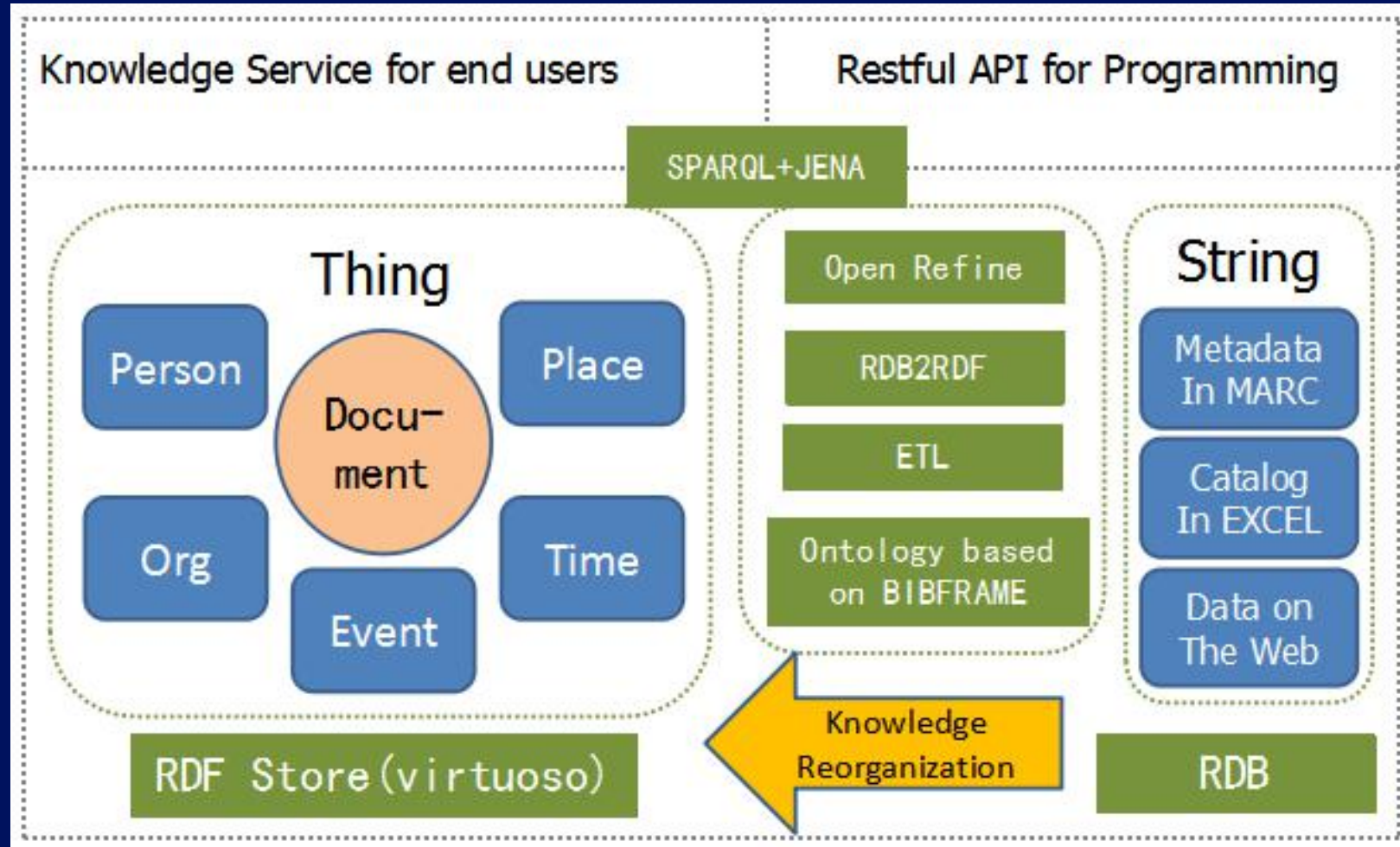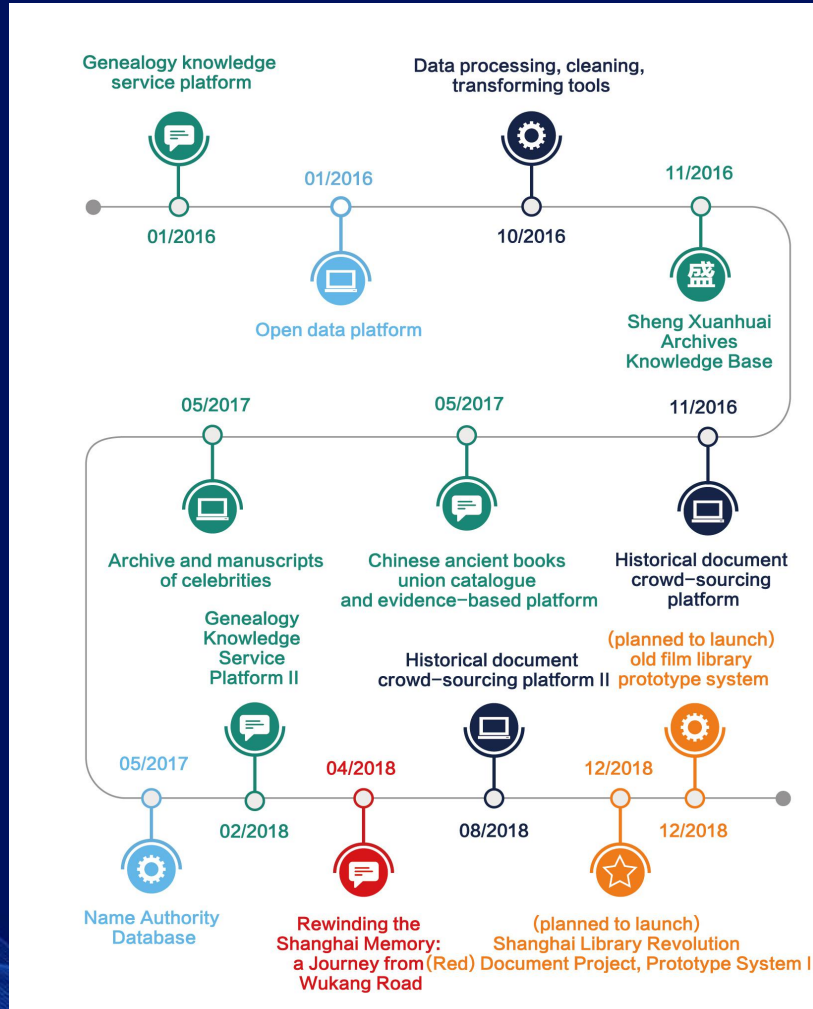# for Digital Humanities

Cuijuan Xia   Shanghai Library   Oct. 14,2021

DCMI 2021 VIRTUAL
Online, 2021 October 4-15

# 01 Background

Current Situation and Problems

# From Digital Library systems to semantic knowledgebases

# Process

- Transform all metadata records of different resources in different formats into RDF triples. Take authors and contributors, publishers and organizations, events and places as entities but not as strings.

- Give every entity cool URI as global identifier and locator. And then link them together after Name Entity Recognization(NER) and entity disambiguation.

- Enrich more semantic data about the entities by extracting structured data from the content of digital resource objects or open datasets on the Web like wiki data.

- Provide authority control and knowledge linking service among resources of different kinds of collections in a web-scale. Provide digital Humanities services for researchers, and open data APIs for the third party developers.

End Users

Developers

Crowd Sourcing | Search

HTTP URI Link | Restful API | Sparql Endpoint

Wukang Road Tour | Red Tour | Old Movie Tour

**Shanghai Memory mobile Web Apps**

Old Photos | Old Movies | Red Revolution Books

Modern Books, Magazines & Newspapers

Manuscripts & Archives | Genealogy | Acient Books

**Special collection knowledgebases**

URI Link

Rest APIs

SH Architecture Knowledgebase | SH Geo Names Dataset

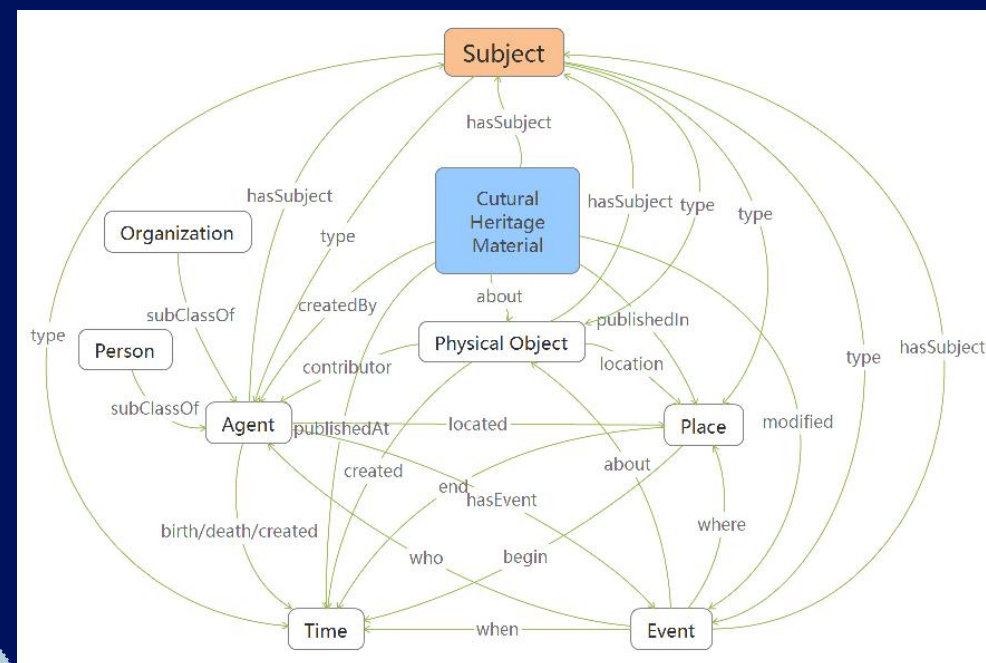GLAM Organization Dataset | Immovable cultural relics Dataset

Name Authority Database | His & Cul Event Knowledgebase

CN Chronology | **Basic knowledgebases** | CN Geo Names

Subject

hasSubject

Organization

hasSubject | type | type

type

Cutural Heritage Material

hasSubject type

createdBy | about | publishedIn

subClassOf

type | Person | Physical Object | location | type | hasSubject

contributor

subClassOf | Agent | publishedAt | located | Place | modified

created | about

end | hasEvent

birth/death/created | where

who | begin | Event

Time | when

Linked entities and data of different
knowledgebases with URIs
semanticly by the One Ontology
Abstract Model.

# Difficulties and Problems

- Data Cleaning and NER.

  - Have to correct the Inconsistency of metadata records in digital library systems and deal with the entity disambiguation manually.

- Decentralized process of entities extraction, creation, modification, publishing, and management.

# Solutions

- **Integrated Workflow**

    - Bridge the gaps among the process of entities management(including extraction, creation, modification, publishing, interlinking) and Service(federal search, visualizaiton).

- **New Technology application**

    - Machine Learning to make the OCR and NER work more efficiently.

- **One platform**

    - support the **Integrated Workflow and new technology application.**

    - support the services(across knowledgebases and visulizations for SNS and spacial-temporal analysis) based on entities and the relations among entites.

# 02 Semantic Content Architecture

The Semantic Content Architecture of Entity Management

# Ontology Abstract Model of All Entities
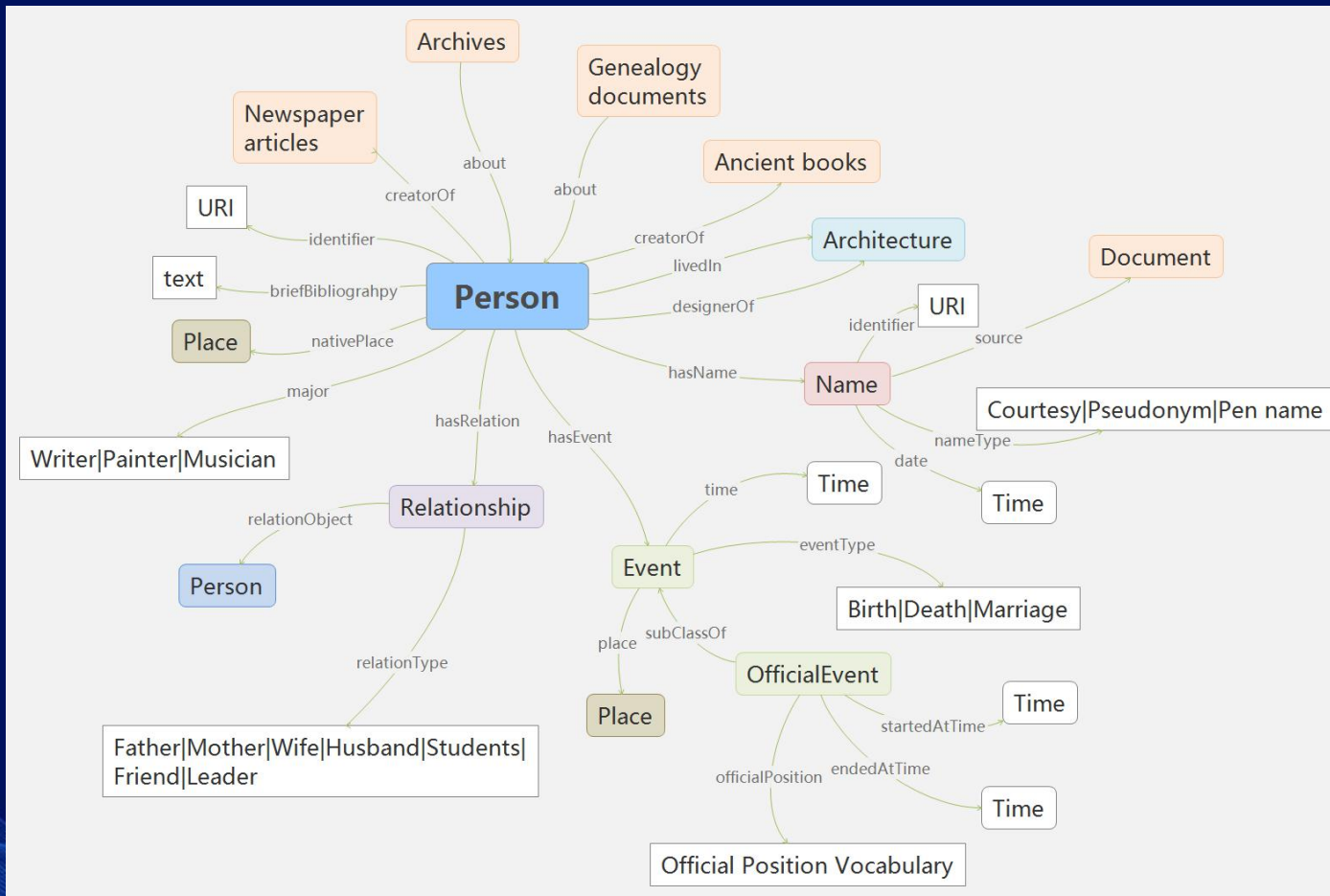
http://data.library.sh.cn/ont/ontology/search



one unified ontology for knowledge modeling, RDF sepecification for consistent knowledge representation

夏翠娟. 文化记忆资源的知识融通： 从异构资源元数据应用纲要到一体化本体设计.图书情报知识,2021(1)
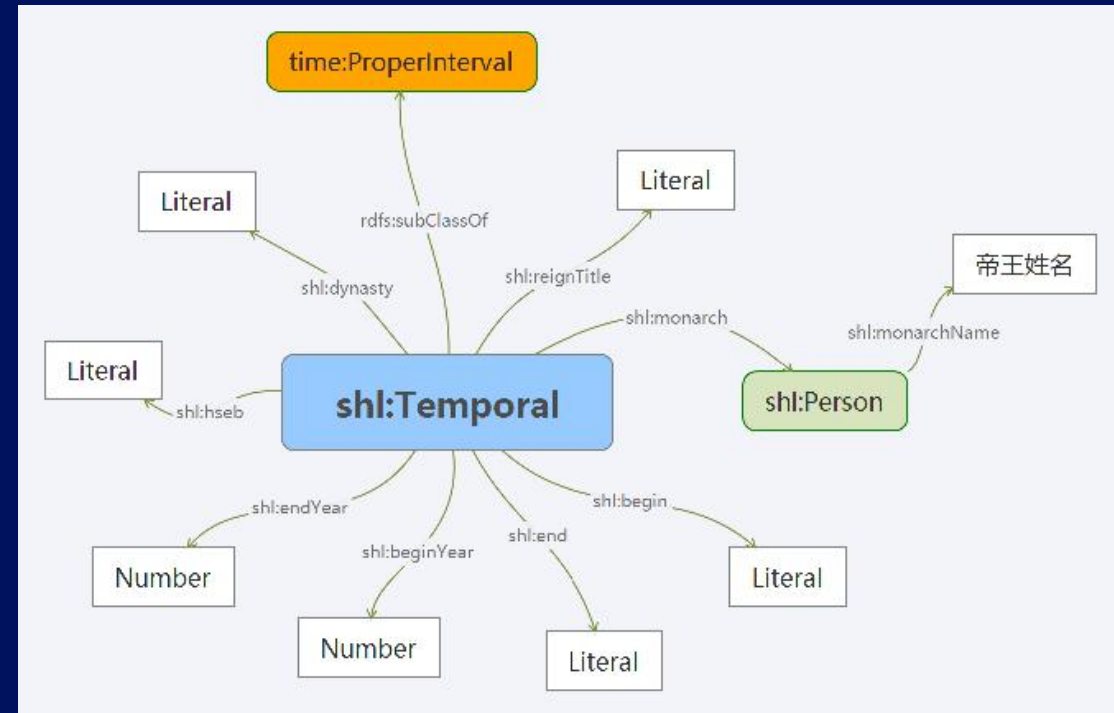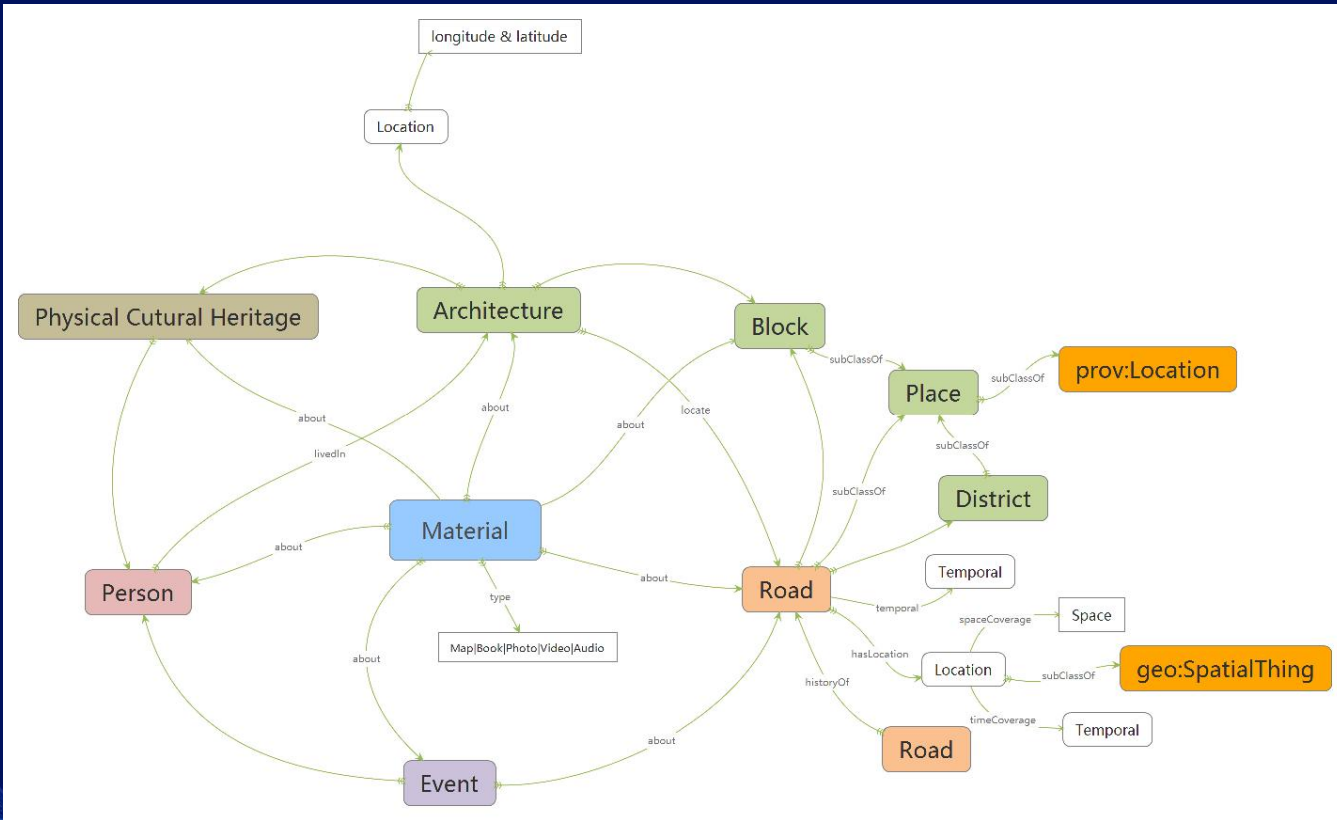
# Ontology Application Profile of Organization



金家琴.夏翠娟.数字人文数据基础设施建设中机构本体的构建：研究和应用[J].图书馆论坛,2020,40(04):30-39.

# Ontology Application Profile of Person



Xia, Cuijuan, and Liu, Wei. "Name Authority Control in Digital Humanities: Building a Name Authority Database of Shanghai Library." International Journal of Libraryship,3.1(2018):21.

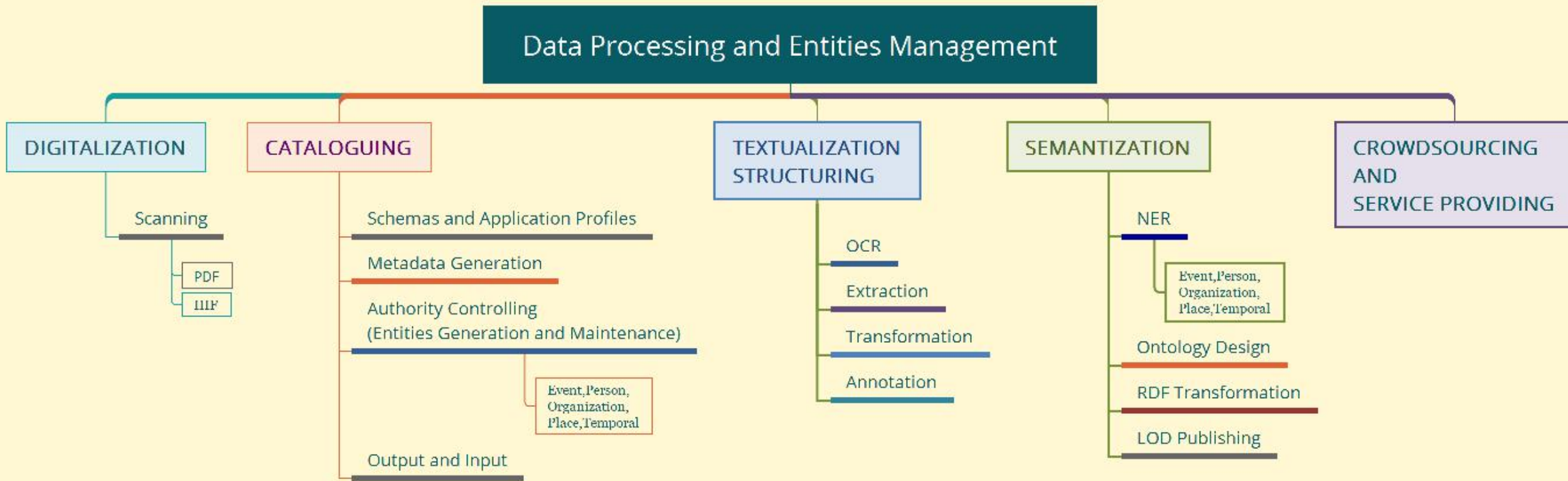# Ontology Application Profile of Place and Temporal



Xia Cuijuan,Wang Lihua,Liu Wei. Shanghai memory as a digital humanities platform to rebuild the history of the city. Digital Scholarship in the Humanities, 2021(3)  https://academic.oup.com/dsh/advance-article/doi/10.1093/llc/fqab023/6178577?guestAccessKey=2a5451a0-5080-4b49-9e5f-bf0fae4f2e10

# Ontology Application Profile of Event

Xia Cuijuan,Wang Lihua,Liu Wei. Shanghai memory as a digital humanities platform to rebuild the history of the city. Digital Scholarship in the Humanities, 2021(3)  https://academic.oup.com/dsh/advance-article/doi/10.1093/llc/fqab023/6178577?guestAccessKey=2a5451a0-5080-4b49-9e5f-bf0fae4f2e10

# 03 Integrated Workflow

Integrated workflow to bridge the gaps among the different stages during the data processing

# Integrated Workflow for
# Data Processing and Entity Management

# Genealogy Union Cataloging System as an Example



Work

Instance

Item

Service Platform

Complete Cat

Abstract only

Holdings

Authority Control
(Entities Maintenance)

Entities Kowledgebases

Check duplicate

create-edit-verify
Workflow

Cataloguing System

Statistics

folio
future of libraries is open
Holdings Manage System

Data synchronization

Import

Export

# Authority Control (Entities Maintenance) during cataloging



Geographic Names Kowledgebase

Chinese History Chronology Kowledgebase

https://jplb.library.sh.cn

# 04 Services Support

SNS and spacial-temporal analysis

# Name Authority Database
# as a Data Linking Hub



http://names.library.sh.cn

# SNS  Social Network analysis



http://dhc.library.sh.c

# HGIS  Spatial-temporal analysis



## Spatial-temporal
## Data Infrastructure

http://dhc.library.sh.c

# Event Knowledgebase
# About the cultural history of Shanghai

http://scc.library.sh.cn  (go live soon)



文化年谱

| 1825-1863 | 1864-1902 | 1903-1941 | 1942-1980 | 1981- |

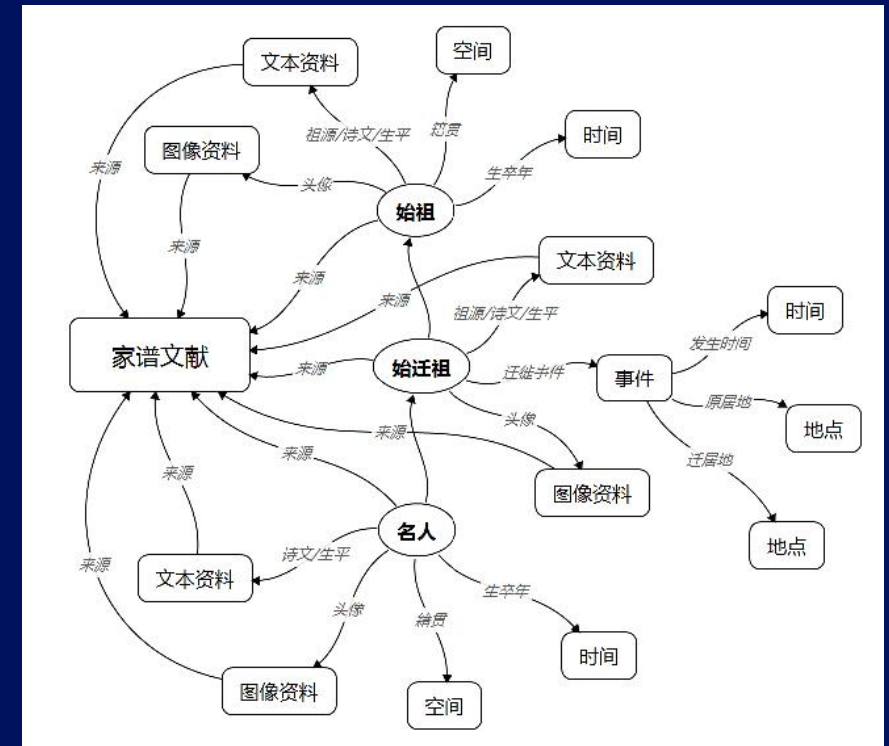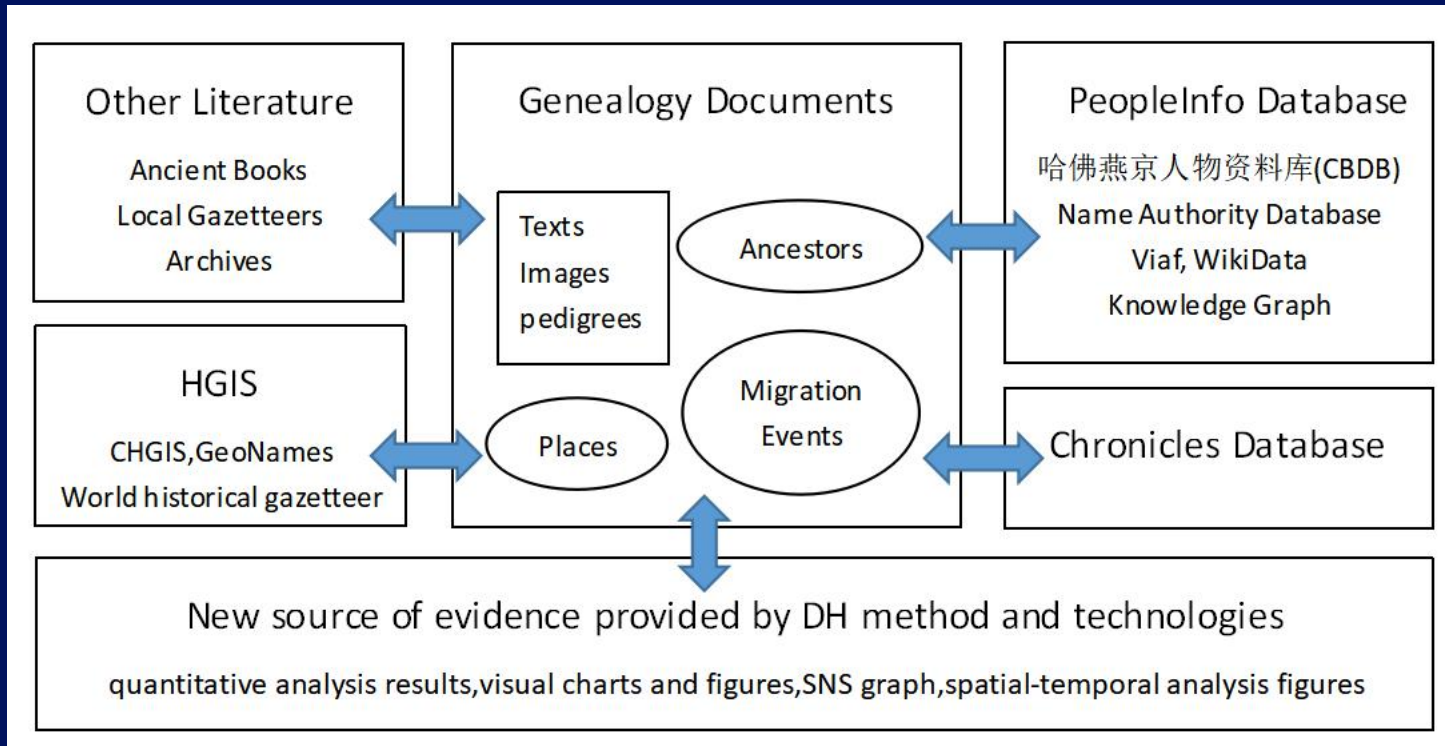| 3241 美术 | 2223 出版 | 1564 文学 | 1010 戏剧 |
| 970 音乐 | 807 电影 | 789 教育 | 740 新闻 |
| 642 建筑 | 369 宗教 | 364 戏曲 | 288 舞蹈 |

# Genealogy Documents as Evidence for Migration History Study



夏翠娟.文化记忆视域下家谱文献价值的再认识和内容的深开发[J].图书与情报,2019.10

data-drven multiple sources of evidence

# Cutural Tourisim

http://wkl.library.sh.cn



Navigate on the street map

open the architecture page and read or listen the brief description

Know more detail about the events happened or people lived in

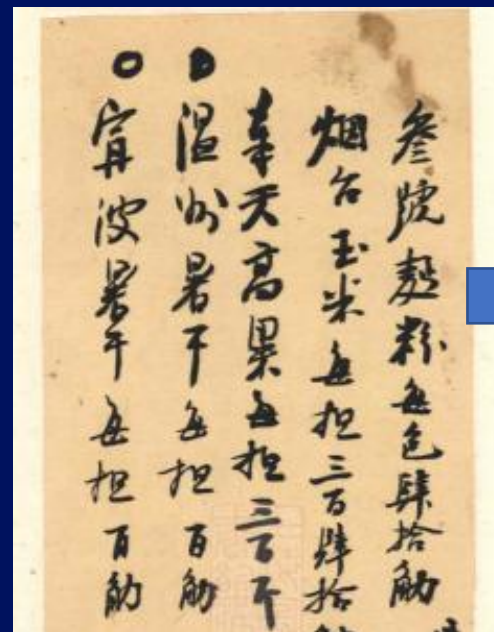Access the different kinds of resources related to the architecture or people

Xia Cuijuan,Wang Lihua,Liu Wei. Shanghai memory as a digital humanities platform to rebuild the history of the city. Digital Scholarship in the Humanities, 2021(3).
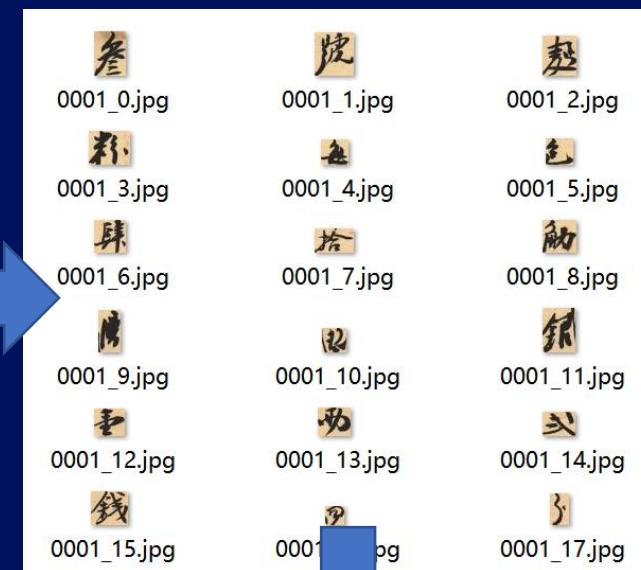
# Crowd-sourcing for UGC (User Generated Content)

- Transcription online

- Captcha



http://zb.library.sh.cn

贺晨芝, 张磊. 图书馆数字人文众包项目实践[J]. 图书馆论坛, 2020, 040(005):3-9.