

GAKG: A Multimodal Geoscience Academic Knowledge Graph

Cheng Deng¹, Yuting Jia¹, Chong Zhang¹, Jingyao Tang¹, Hui Xu¹, Luoyi Fu¹, Weinan Zhang¹,
Haisong Zhang², Xinbing Wang¹, Chenghu Zhou³

¹Shanghai Jiao Tong University, ²Tencent AI Lab

³Institute of Geographical Science and Natural Resources Research, Chinese Academy of Sciences
{davendw,hnxxjyt,yiluofu}@sjtu.edu.cn,hansonzhang@tencent.com,zhouch@lreis.ac.cn

Presented by Cheng Deng

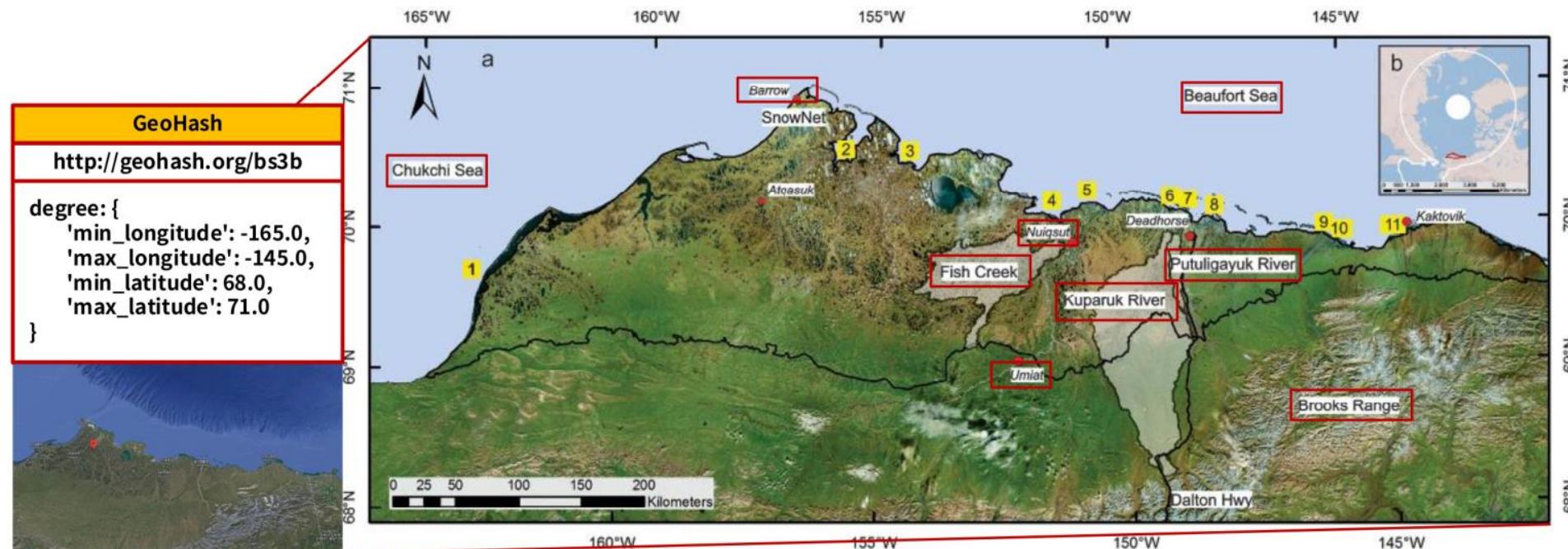
- ✉ <http://www.big-cheng.com>
- ⌚ <https://github.com/davendw49>
- ✉ davendw@sjtu.edu.cn



上海交通大学
SHANGHAI JIAO TONG UNIVERSITY

CIKM 2021

30th ACM INTERNATIONAL CONFERENCE ON INFORMATION AND KNOWLEDGE MANAGEMENT

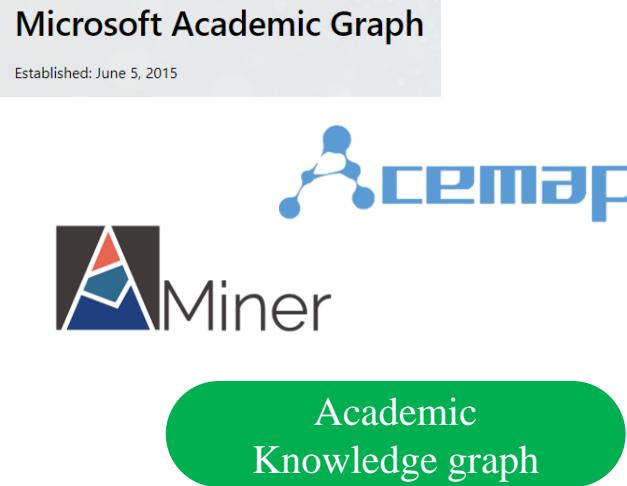
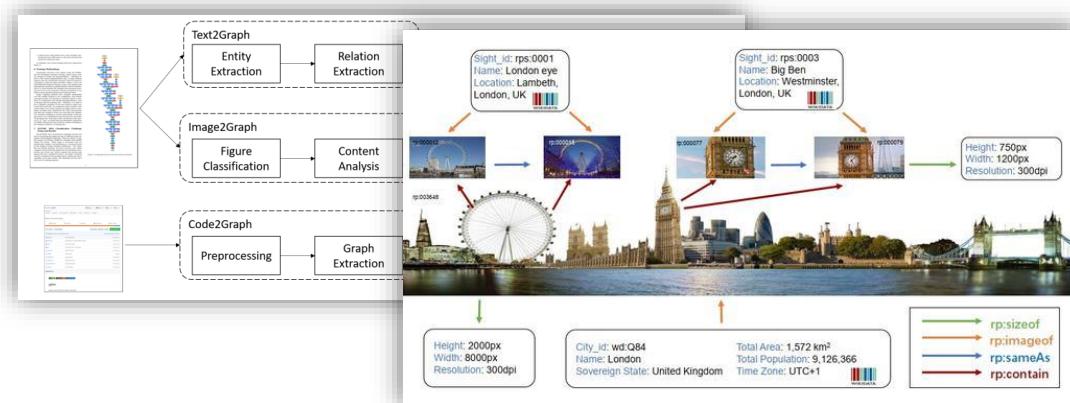


A large quantity of literatures

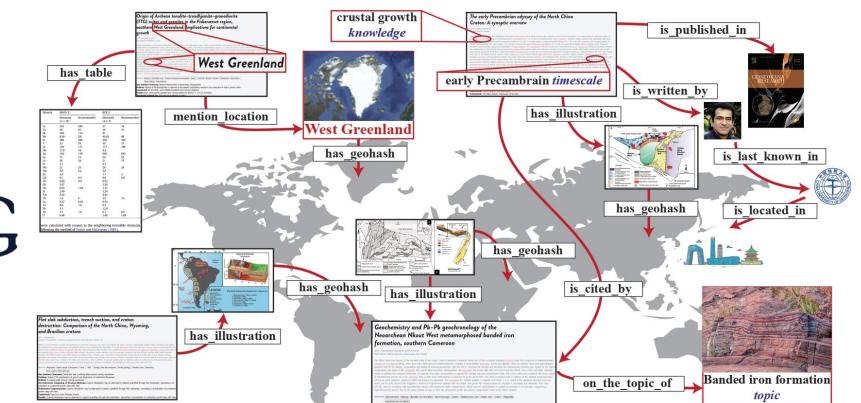
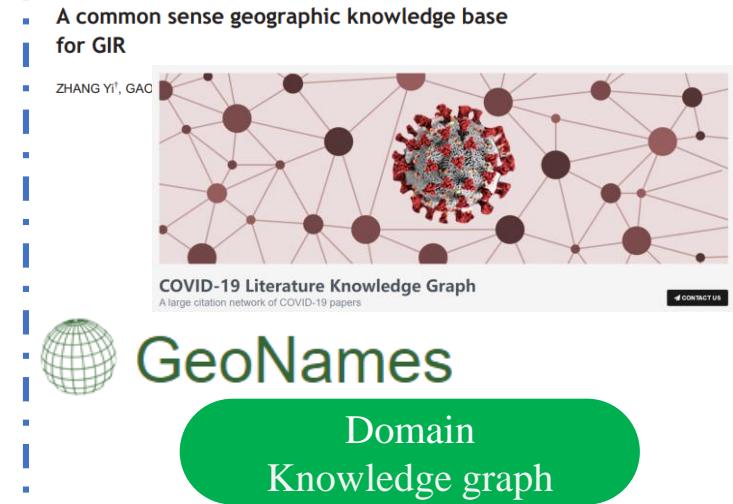
Papers, a major mean to disseminate knowledge

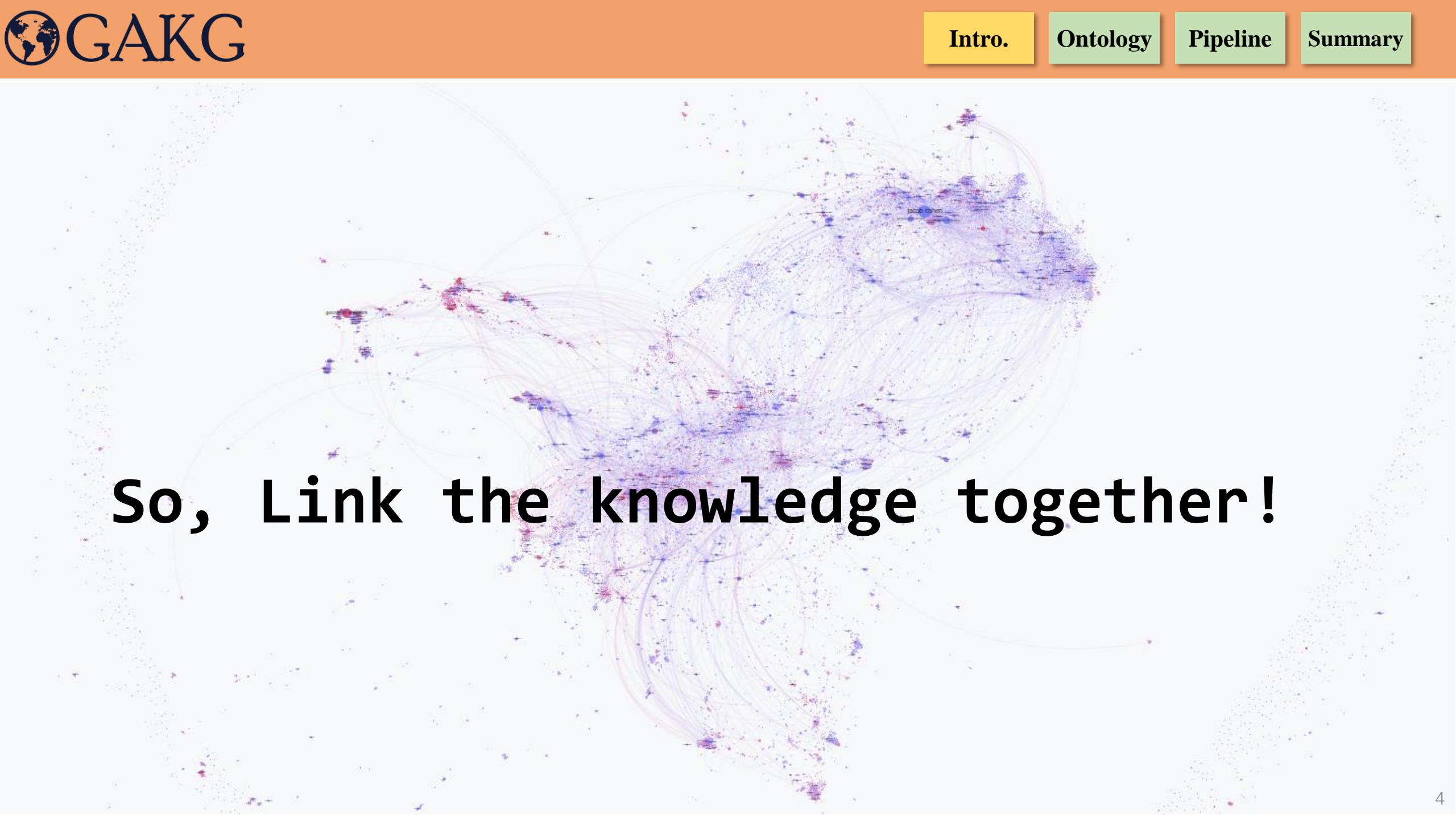
geoscience papers contain abundant multimodal data

Papers thus have time and spatial characteristics

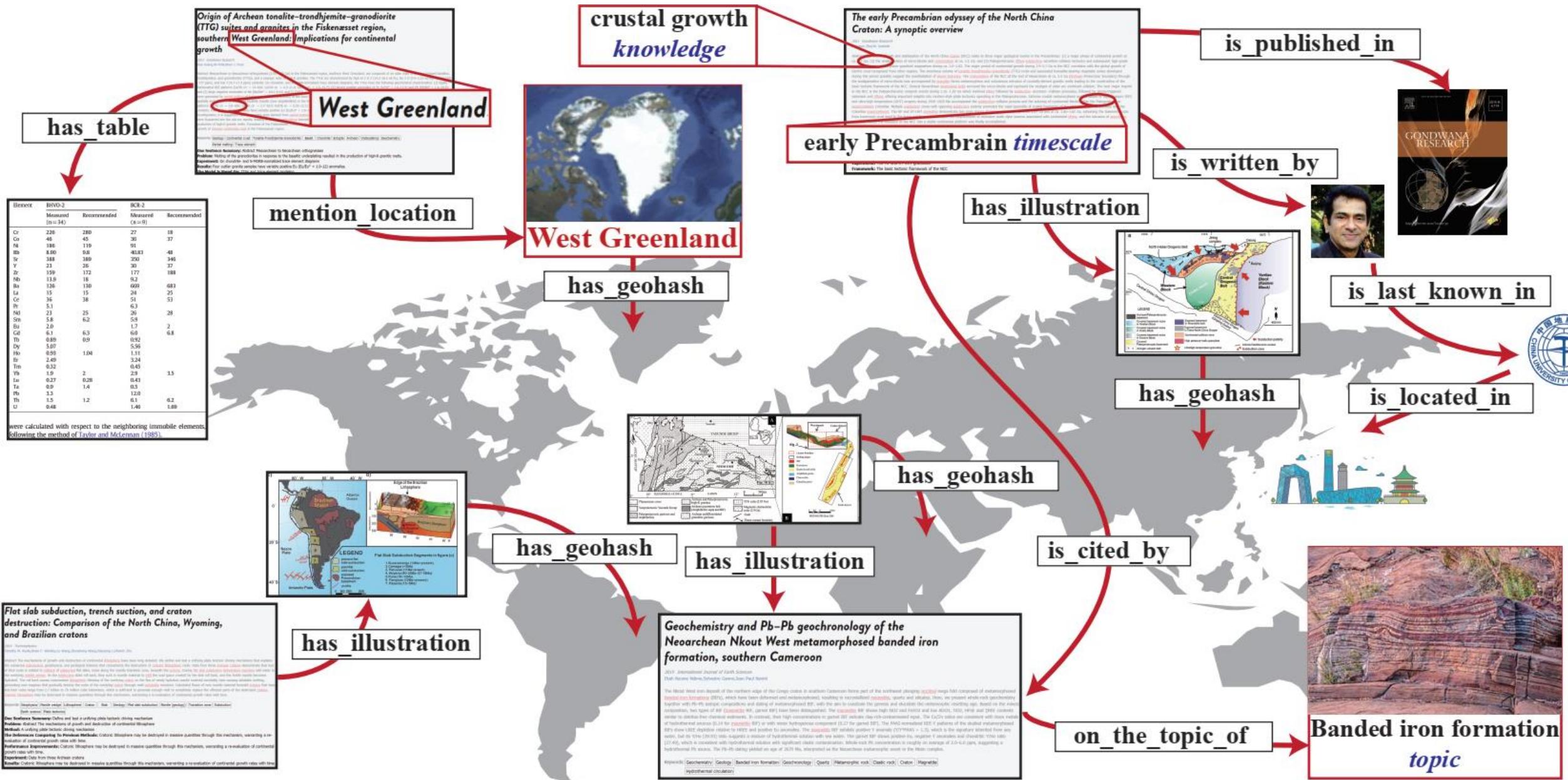


Multimodal Geoscience Academic Knowledge Graph





So, Link the knowledge together!



GAKG's schema-graph consists of **11** concepts connected by **19** relations. Five of them (*has_concluded*, *has_designed*, *is_located_in*, *has_developed* and *earn_in_the_way_of*) have a upper class relation *acer:mention_knowledge*. Since **GAKG** is the union of academic concepts and their relations, we manage GAKG as linked open data (LOD), we provide *#sameAs* axioms linking to the entities in other datasets, **271 thousands** in total.

The Graph base namespace (Graph IRI) is <https://www.acekg.cn>, all the concepts and relations shared.

To our knowledge, GAKG is currently the **largest** and **most comprehensive** geoscience academic knowledge graph, consisting more than **68 million triples, including 8 millions concepts and 41 million links**.

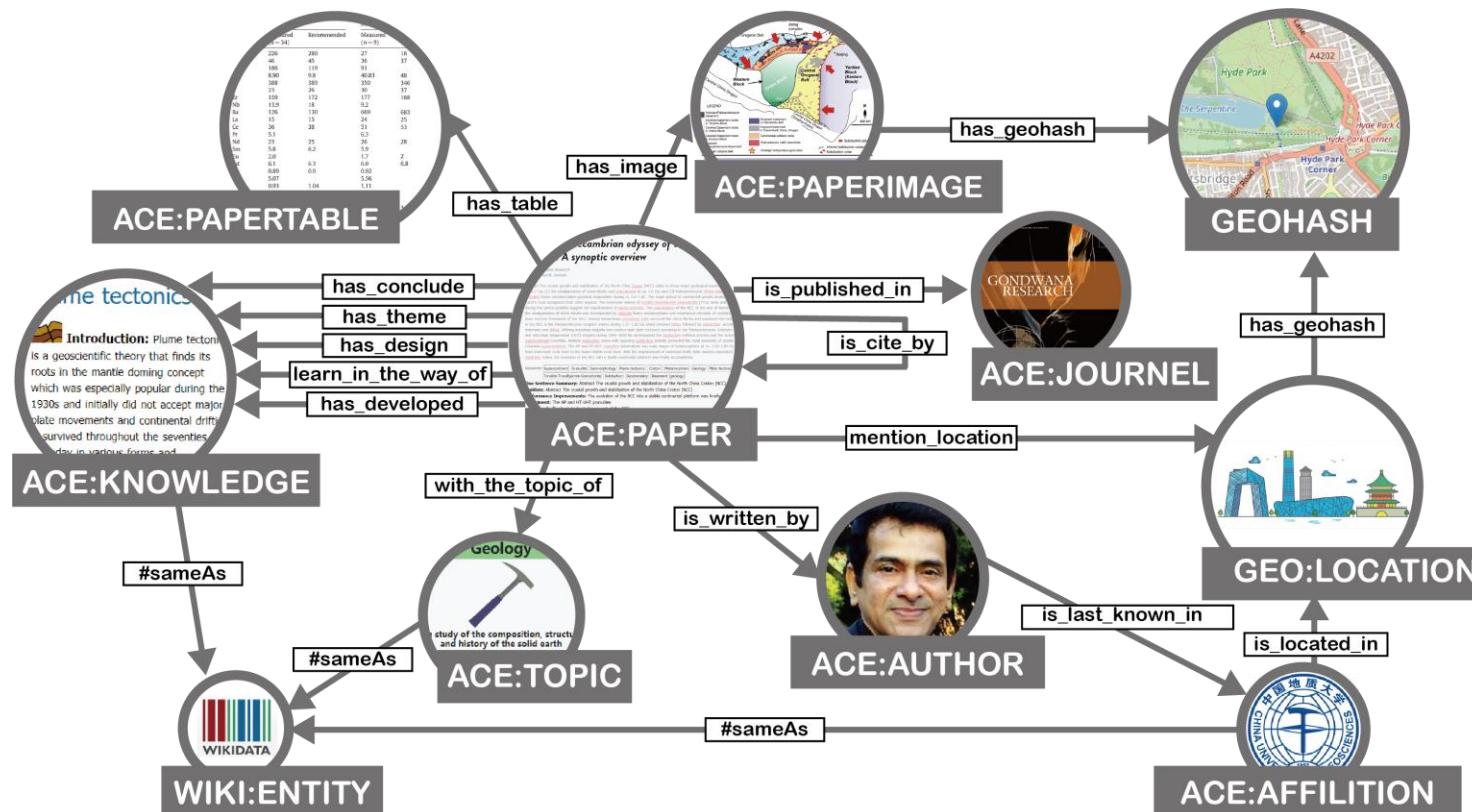


Table 1: Statistics of GAKG Concepts (Up to May 30, 2021).

Concept	Count	Concept	Count
paper	1,122,094	knowledge	62,576
author	908,933	illustration	3,562,816
affiliation	27,175	papertable	760,054
topic	765,184	location	784,279
journal	194	geohash	996,731
timescale	1,701	Total	8,991,737

Table 2: Statistics of GAKG Relations (Up to May 30, 2021).

Relation	Count	Relation	Count
is_cited_by	17,704,495	mention_knowledge	704,899
on_the_topic_of	10,401,972	mention_location	759,260
is_written_by	3,547,077	has_geohash	1,021,870
is_published_in	1,122,094	mention_timescale	1,120,398
is_last_known_in	662,850	in_the_period_of	189
is_located_in	25,019	before	155
has_illustration	3,562,816	#sameAs	271,156
has_table	760,054	Total	41,664,304

Snorql for GAKG

Snorql for GAKG is an extension of [Snorql](#). Use [basic endpoint](#) for application.

- Query examples are provided below (or right-hand side of) the text area. See also [GAKG RDF dataset](#) for the general description.
- Example for paper: [The early Precambrian odyssey of the North China Craton: A synoptic overview](#)
- Example for affiliation: [Utrecht University](#)
- Example for author: [M. Santosh](#)
- Example for timescale: [miocene](#)
- Example for location: [kitaké](#)
- Example for illustration: [Illustration/42989303](#)

SPARQL query:

```
PREFIX xsd: <http://www.w3.org/2001/XMLSchema#>
PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
PREFIX foaf: <http://xmlns.com/foaf/0.1/>
PREFIX skos: <http://www.w3.org/2004/02/skos/core#>
PREFIX iiswc: <http://annotation.semanticweb.org/iswc/iswc.daml#>
PREFIX acep: <https://www.acekg.cn/property/>
PREFIX ace: <https://www.acekg.cn/relation/>
PREFIX acec: <https://www.acekg.cn/concept/>
PREFIX acepgeo: <https://www.acekg.cn/property/geo/>
PREFIX geor: <https://www.acekg.cn/relation/geo/>
PREFIX geo: <https://www.acekg.cn/concept/geo/>
PREFIX owl: <http://www.w3.org/2002/07/owl#>

SELECT DISTINCT ?s ?t ?img WHERE {
  ?s rdf:type ace:paper;
  ?s acep:year 2000 .
  ?s acec:has_illustration ?img .
  ?s acep:title ?t
}LIMIT 10
```

Results: [Browse](#) [Run Query](#)

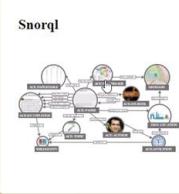
AUTHOR: AFFILIATION: YEAR: TITLE: CONCEPT: EasySPARQL

[Snorql](#) for GAKG v1.0. Use [basic SPARQL endpoint](#) for your application.

Home SPARQL-Endpoint Svirql Downloads HTL Github About Us

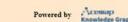
 **GAKG**
Multimodal Geoscience Academic Knowledge Graph

Q Carbonate Rock [SEARCH](#)

Snorql 

Map Search 

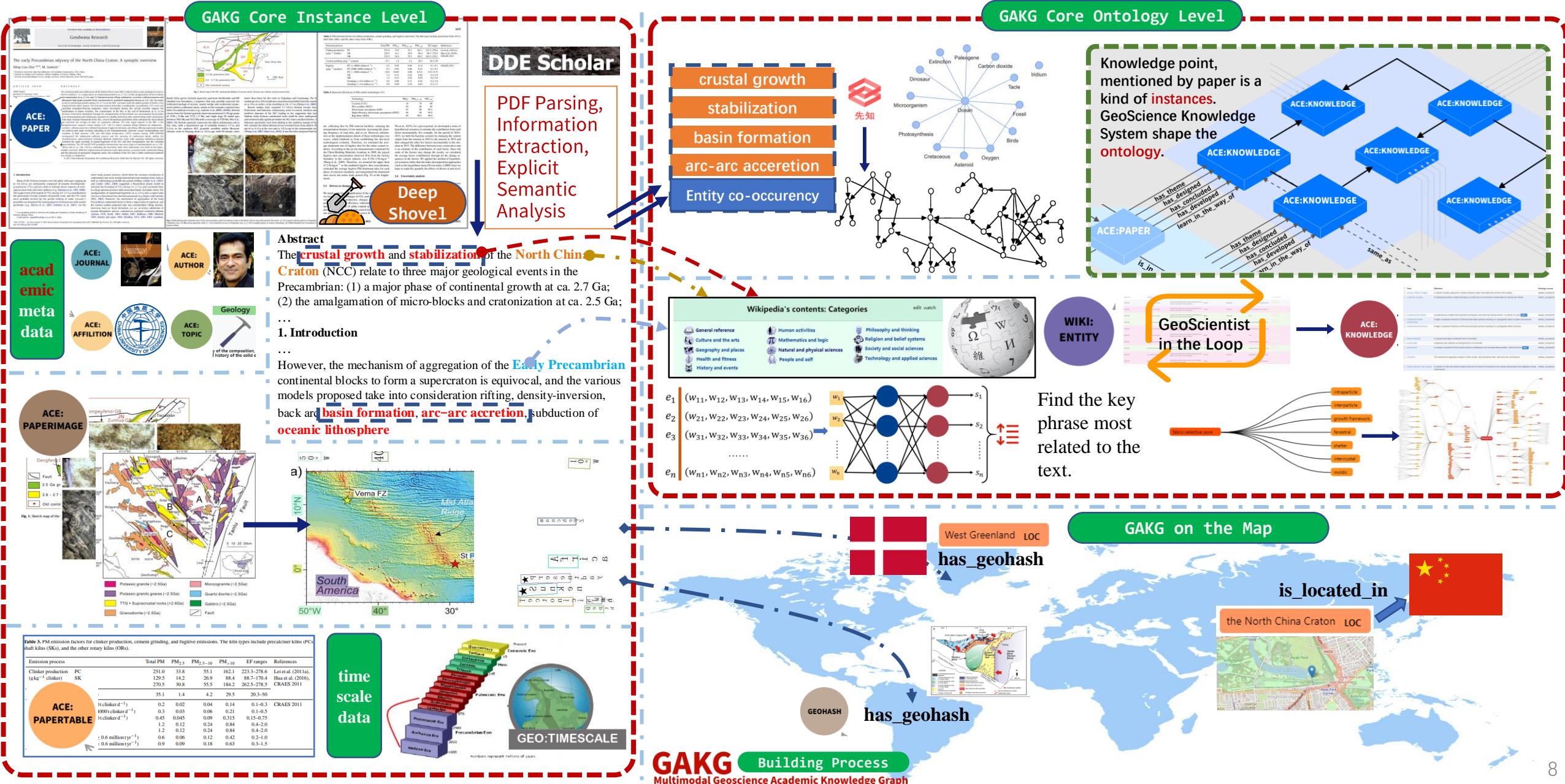
Svirql 

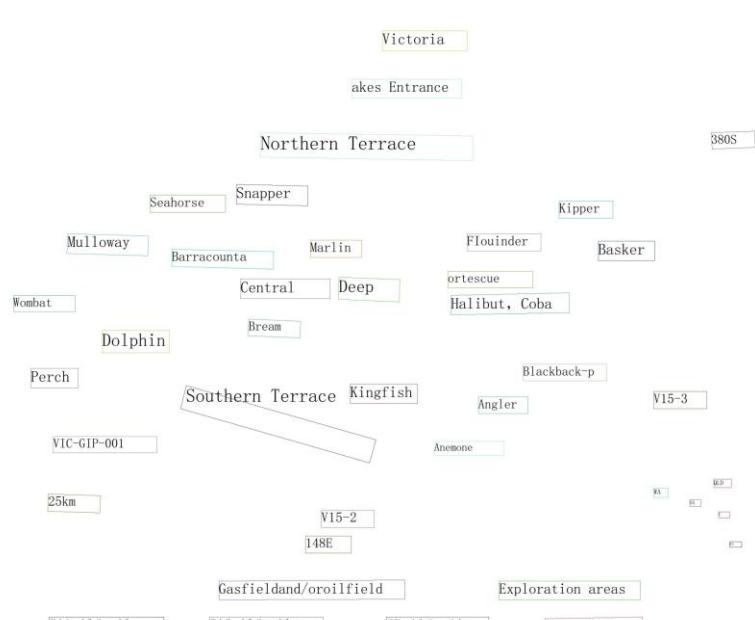
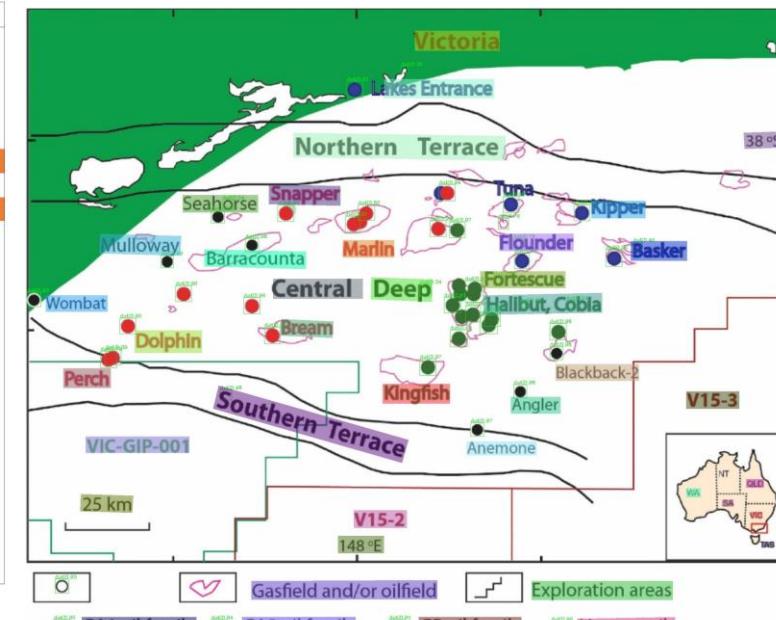
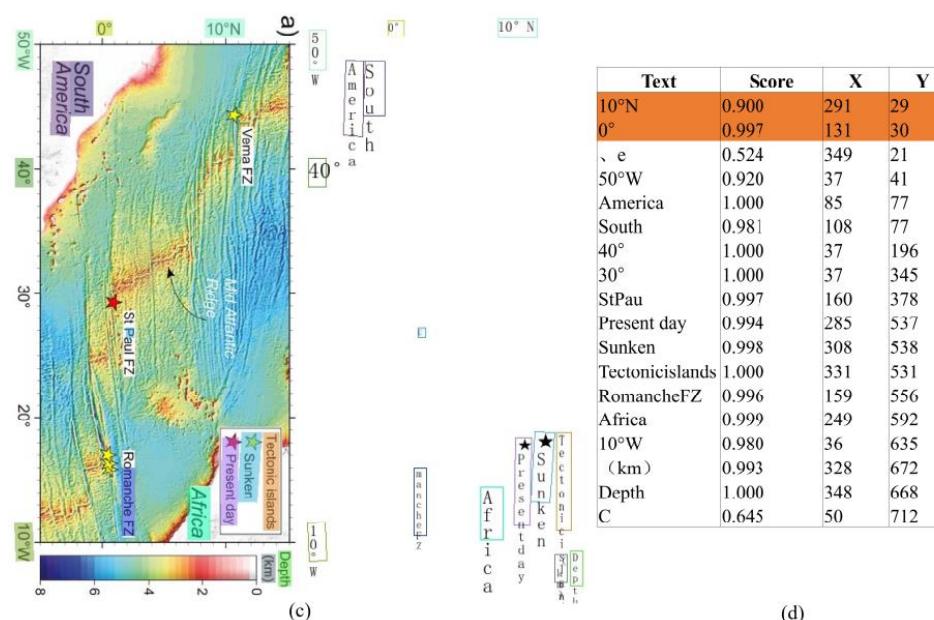
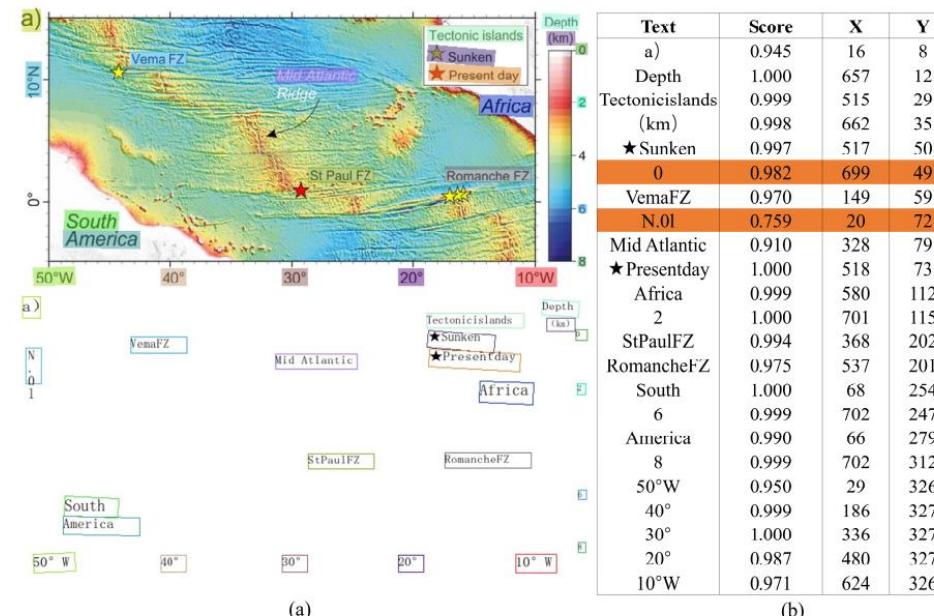
Powered by 
© Copyright 2015-2021 Acemap, Inc. Shanghai Jiao Tong University.

GAKG SPARQL: <https://www.acekg.cn/sparql>

GAKG Snorql: <https://snorql.acemap.cn/>

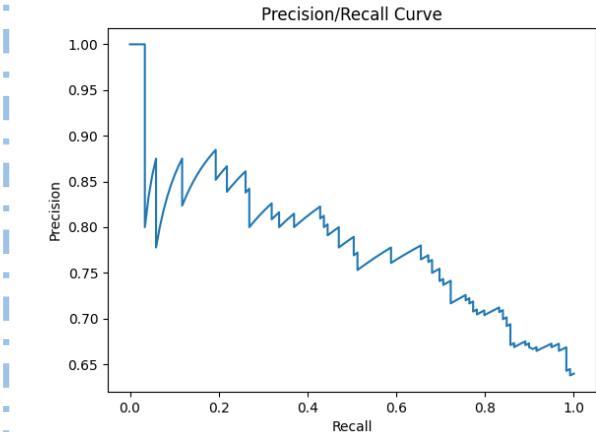
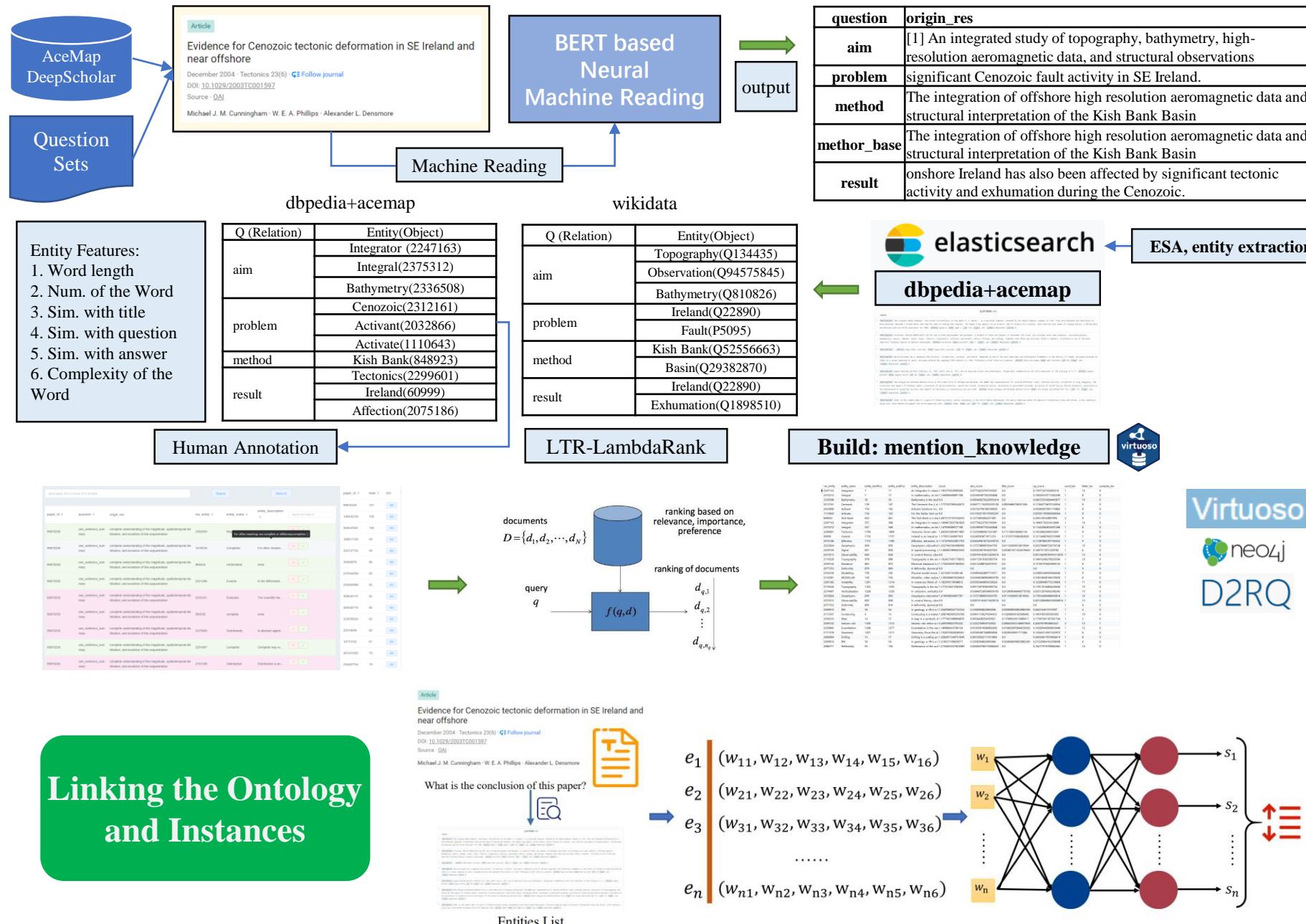
- Ref: SNORQL: <https://github.com/kurtjx/SNORQL>
- Ref: Virtuoso: <https://virtuoso.openlinksw.com/>





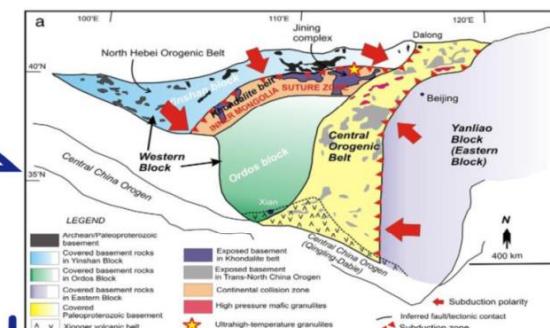
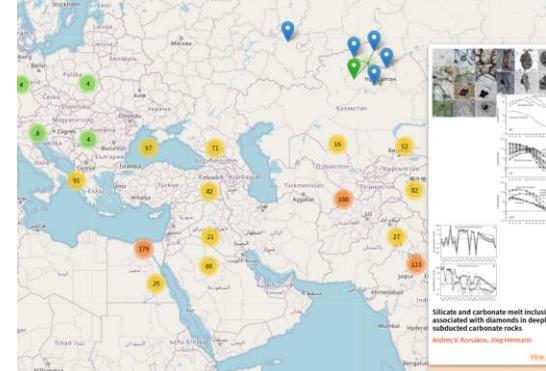
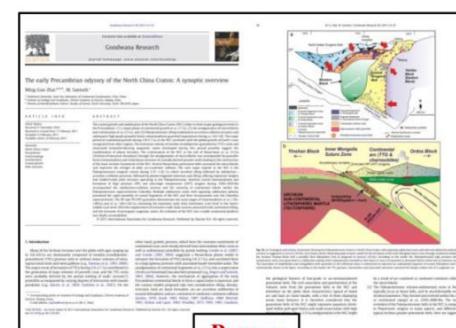
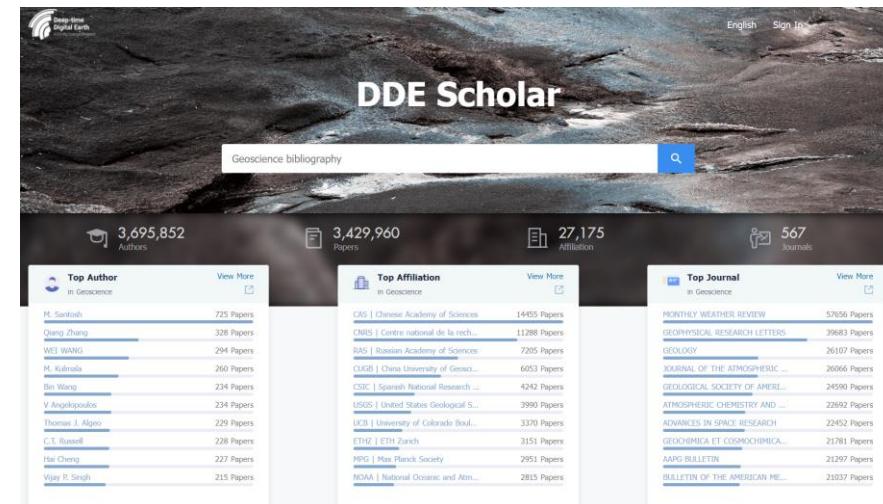
An Example of Illustrations Coordinates Extraction.

- (a) is the result of OCR,
- (b) highlights the recognition error,
- (c) is the result of OCR after using a rule-based method adjusting image,
- and (d) highlights the corrected coordinates.

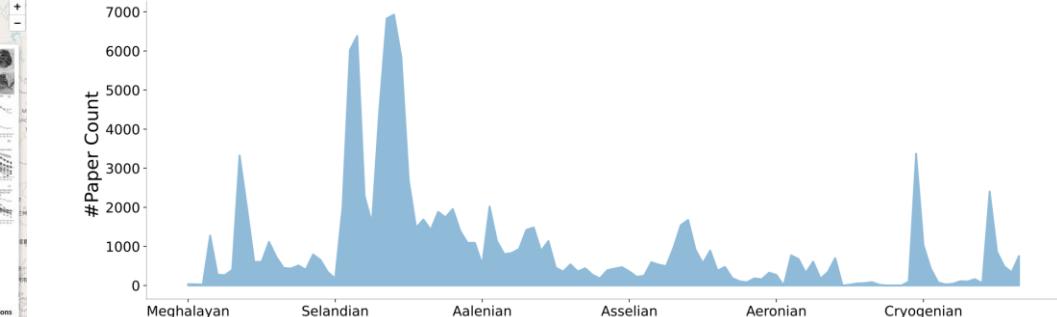


Models	precision
ESA without Supervision	0.63
ESA + Learn2rank (after 1 st loop HITL)	0.79 (^{↑25%})

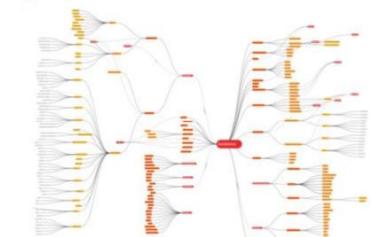
Models	precision
Autophrase without Supervision	0.59
Autophrase + Learn2rank (after 1 st loop HITL)	0.78 (^{↑32%})



Extract the Knowledge Entities, Locations and Timescales



Geographic Information Extraction



Knowledge Entity Extraction



Time scale Information Extraction

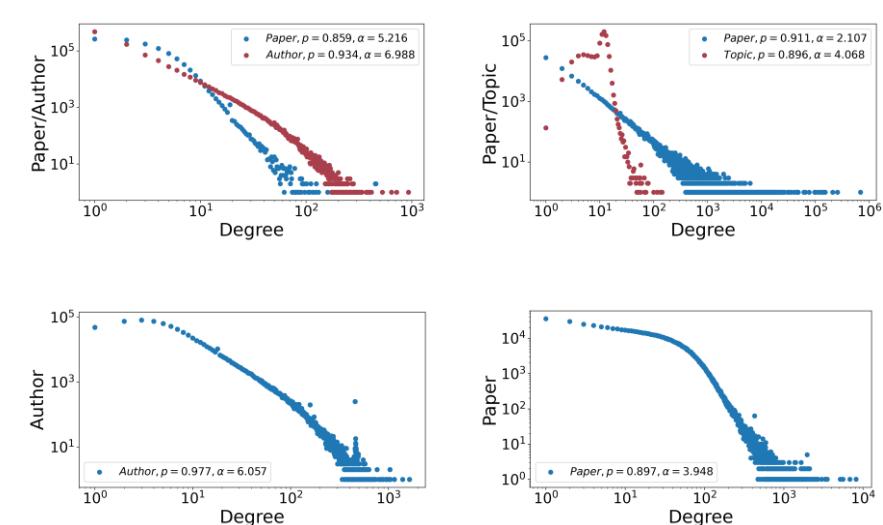


Figure 8: Papers Distribution along with Geologic Era.

Benchmark-CD

Table 5: Statistics of Community Detection Benchmarks.

Benchmarks	Number of communities	Nodes	Edges	Nodes in Largest WCC	Edges in Largest WCC	Nodes in Largest SCC	Edges in Largest SCC	Average Cluster Coefficient	Triangles	Diameter
GPCN	194	842,121	16,034,510	838,219 (0.995)	16,031,892 (0.999)	0	0	0.0699	38,789,469	176
GACN	194	860,280	5,381,861	752,718 (0.875)	5,282,032 (0.972)	752,718 (0.875)	5,282,032 (0.972)	0.6897	43,502,542	15
Email-Eu-core	42	1,005	25,571	986 (0.981)	25,552 (0.999)	803 (0.799)	24,729 (0.967)	0.3994	105,461	7
CORA	7	2,708	5,429	2,485 (0.918)	2,604 (0.493)	13 (0.005)	14 (0.003)	0.1314	1,630	15
DBLP (Collaboration Network)	2,547	317,080	1,049,866	317,080 (1.000)	1,049,866 (1.000)	317,080 (1.000)	1,049,866 (1.000)	0.6324	2,224,385	21
Amazon (Product Network)	5,000	334,863	925,872	334,863 (1.000)	925,872 (1.000)	334,863 (1.000)	925,872 (1.000)	0.3967	667,129	44

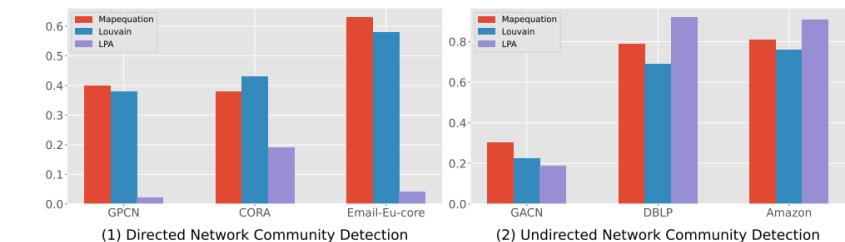


Figure 9: Community Detection Evaluation Results.

Benchmark-KE

Table 6: Statistics of KRL Benchmarks.

Benchmark	relation	entity	triple
FB15K	1,345	14,951	483,142
WN18	18	40,943	141,442
GA16K	10	16,363	151,662

Table 7: Results of Link Prediction Task.

Models	FB15K		WN18		GA16K	
	MR	hit@10	MR	hit@10	MR	hit@10
RESCAL	683	0.441	1,163	0.528	4,300	0.001
TransE	125	0.471	251	0.892	280	0.320
TransH	84	0.585	303	0.867	337	0.325
RotatE	40	0.884	309	0.959	214	0.366
SimpIE	74	0.876	412	0.947	311	0.260

Contribution

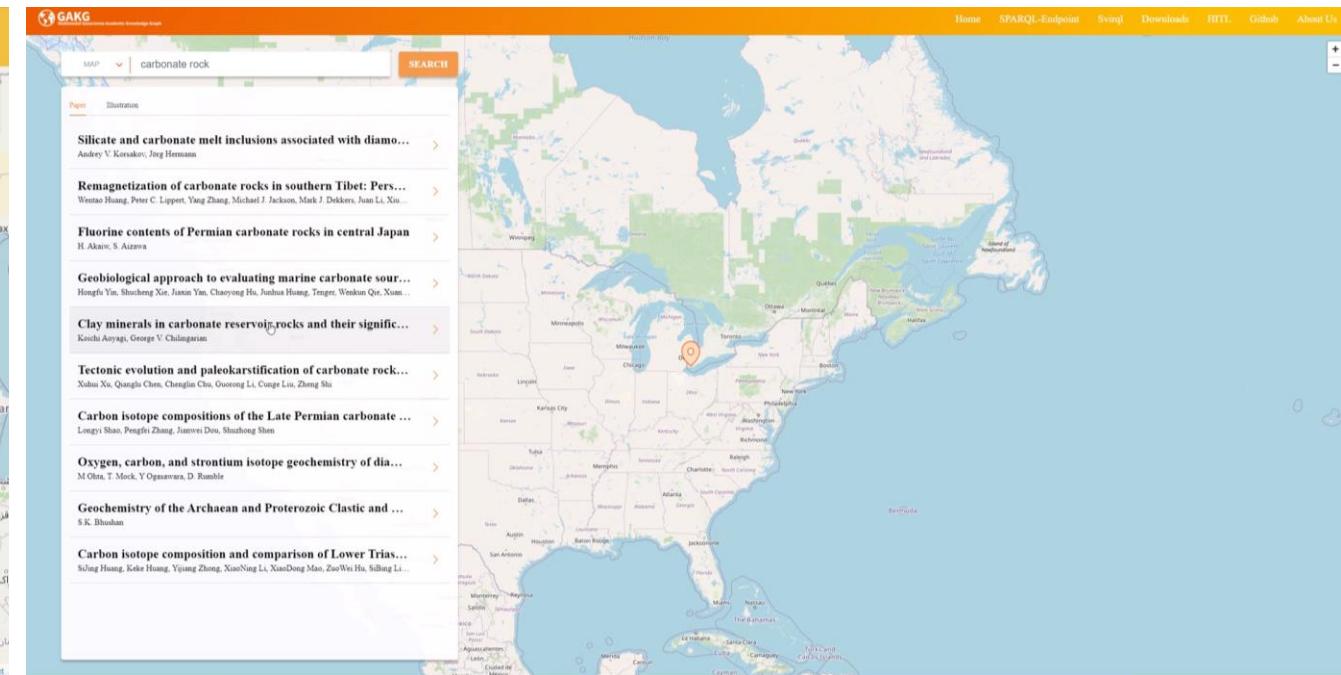
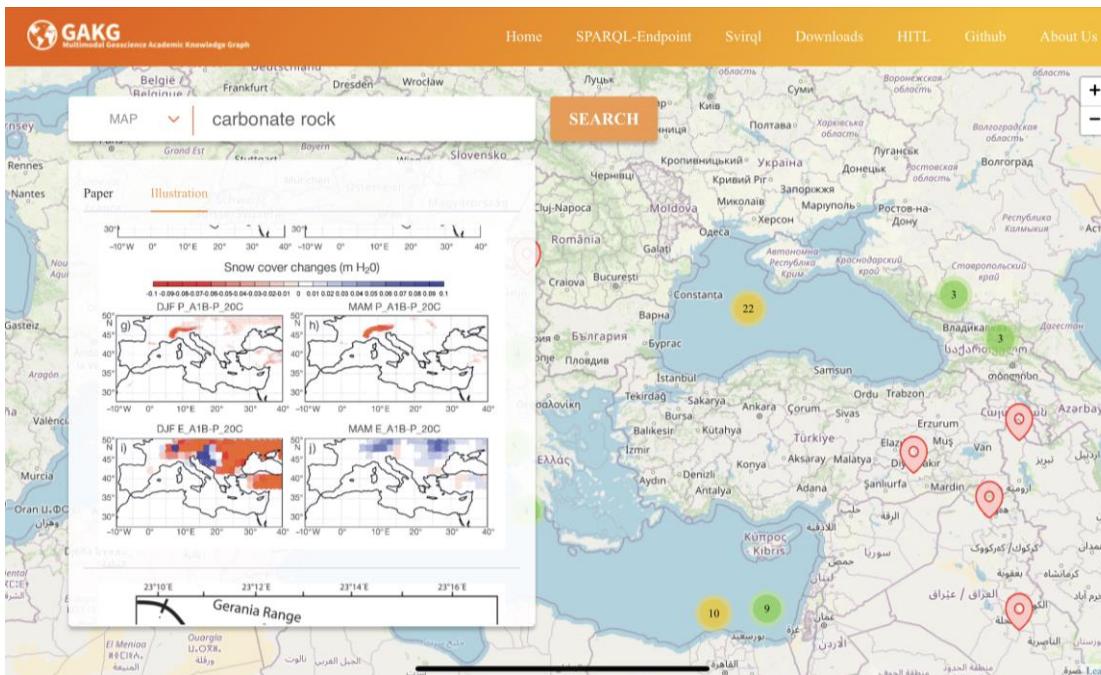
- We propose a multimodal **GeoScience Academic Knowledge Graph (GAKG)** framework by fusing papers' images, text, and bibliometric data.
- With a world map, all the illustrations, text, and geologic time scale extracted from the selected geoscience papers can be highly connected to the geographical information
- We put forward a Human-In-the-Loop knowledge extraction pipeline to extract paper's knowledge entities and mapping them to a crowd-sourcing knowledge taxonomy.



Observation

- *Availability.*
 - *Dump files*
 - *SPARQL Endpoint*
 - *Snorql Query System*
 - *Papers on the Map*
 - *Source Code on Github*
- *Quality.*
 - *Data sources: Acemap, DBpedia, DDE*
 - *Verified partially by Geoscientists*
- *Limitations*
 - *Lack of annotation during the Human-In-The-Loop*

- **Geographic Information Retrieval.** We provide a knowledge based search engine on a geographical map for the literature of geoscience.
- If the researcher drags the window, the distribution of the papers will change accordingly.
- **LINK:** <https://gakg.acemap.info/>



- **Geoscience KBQA.** Based on GAKG, research can know more information about the relation between papers. These template-based queries can be applied in scientific research and academic communication. **These questions are also generally inextricable by existing Q&A systems and search engines.**

- Show the pictures of the Luca Pozzoli's papers if the paper is talking about Chemistry GO

< BACK



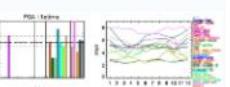
Using SEVIRI fire observations to drive smoke plumes in the CMAQ air quality model: a case study over Antalya in 2008
<https://www.acemap.info/paper/1003002>



Using SEVIRI fire observations to drive smoke plumes in the CMAQ air quality model: a case study over Antalya in 2008
<https://www.acemap.info/paper/1003002>



The AeroCom evaluation and intercomparison of organic aerosol in global models
<https://www.acemap.info/paper/489199435>



The AeroCom evaluation and intercomparison of organic aerosol in global models
<https://www.acemap.info/paper/489199435>

- One-hop queries, such as returning papers targeting a particular topic,
- Two-hop queries, such as querying illustrations in a specific field,
- Three-hop queries, such as querying geographic locations that a certain affiliation often studies,
- Four-hop queries, such as querying the relationship between geographic locations and affiliations.

Home SPARQL-Endpoint Sparql Downloads HTML GitHub About Us

Sparql for RDF schema AUTHOR AFFILIATION YEAR TITLE CONCEPT GENERATE SPARQL SHOW ▾

About Papers' Images

- Show the pictures about Organic matter and Chemistry GO

- Show the pictures of the Luca Pozzoli's papers if the paper is talking about Chemistry GO

About Papers

- Show papers about Organic matter and Chemistry GO

- What did the paper Organic have concluded designed developed targeted at test in the way of GO

- Which affiliation is undergoing researching Organic matter GO

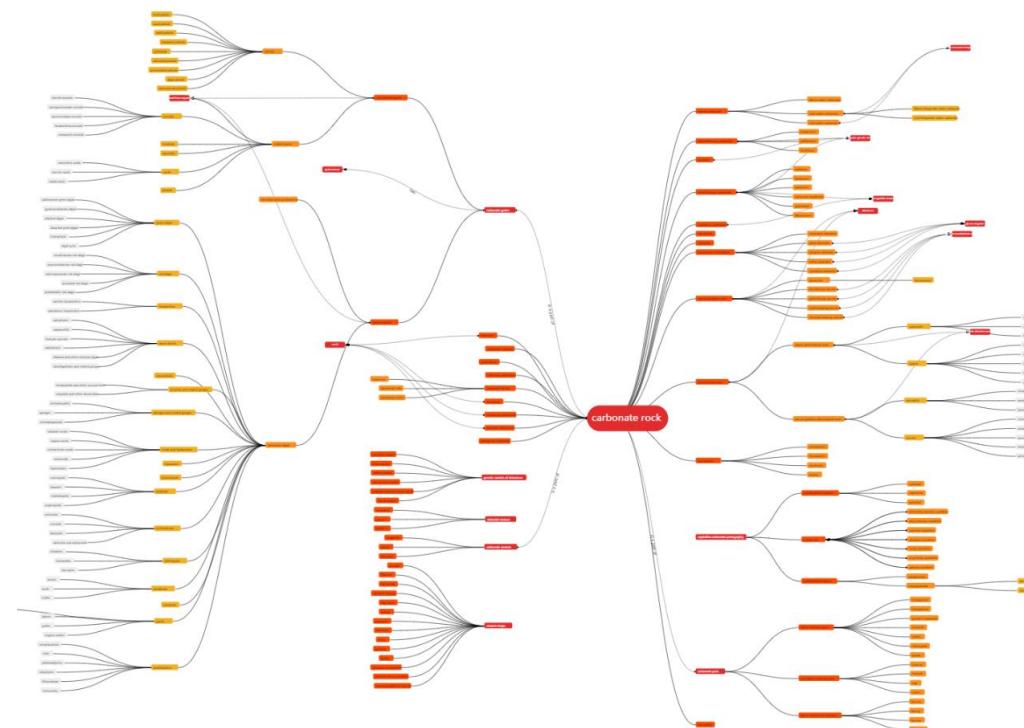
- Which author is undergoing researching Organic matter GO

- Show me the papers in the field of Chemistry published in 2020 GO

- Show me the papers written by Gongjun Tang published before 2020 GO

Future Works

- *Academic Knowledge System Construction with Human-In-the-Loop.*
- *Social Community Detection in Geoscience.*
- *Scientific Articles' Geographical Information Extraction.*



GAKG

<https://gakg.acemap.info/>

Multimodal Geoscience Academic Knowledge Graph

Github: <https://github.com/davendw49/gakg>

Thank You!

GAKG: A Multimodal Geoscience Academic Knowledge Graph

Cheng Deng¹, Yuting Jia¹, Chong Zhang¹, Jingyao Tang¹, Hui Xu¹, Luoyi Fu¹, Weinan Zhang¹,
Haisong Zhang², Xinbing Wang¹, Chenghu Zhou³

¹Shanghai Jiao Tong University, ²Tencent AI Lab

³Institute of Geographical Science and Natural Resources Research, Chinese Academy of Sciences
{davendw,hnxxyt,yiluofu}@sjtu.edu.cn, hansonzhang@tencent.com, zhouch@lreis.ac.cn



上海交通大学
SHANGHAI JIAO TONG UNIVERSITY

Presented by Cheng Deng

<http://www.big-cheng.com>
 <https://github.com/davendw49>
 davendw@sjtu.edu.cn