



Transforming Unstructured **Historical and Geographic Data** into **Spatio-Temporal Knowledge Graphs**

Basel Shbita


PhD Thesis Defense

05/01/2024

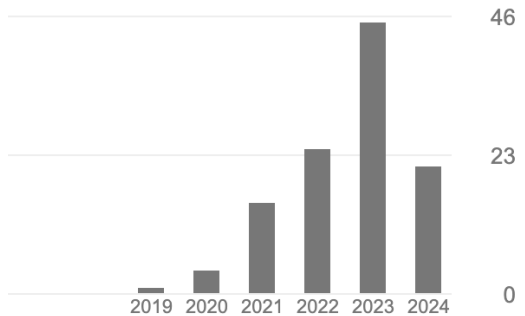
Committee Members:

Craig A. Knoblock, Cyrus Shahabi, John P. Wilson, Jay Pujara, and Yao-Yi Chiang

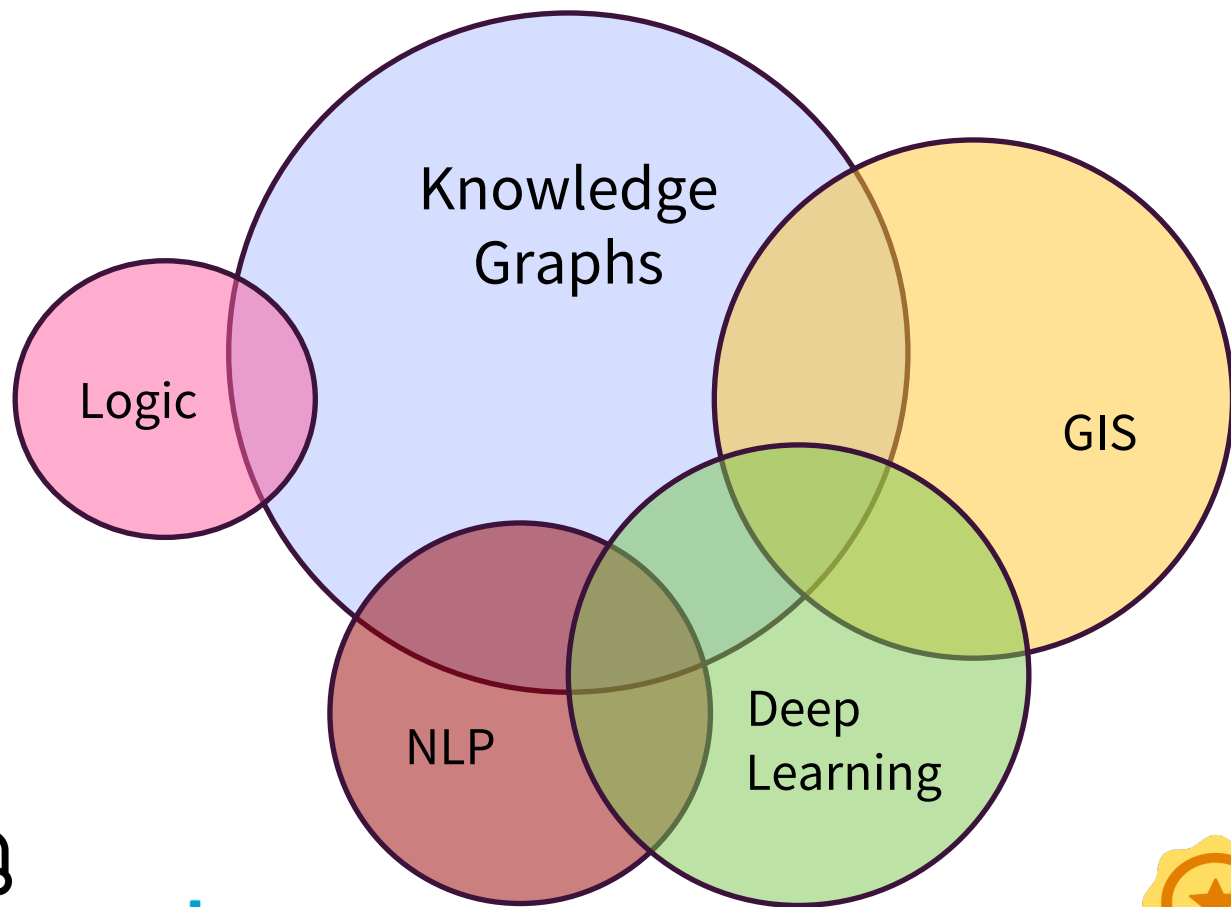
Agenda

- Basel's PhD Journey 
- Intro
- Thesis Overview
- Approach:
 - Building Spatio-Temporal KGs from Digitized Maps
 - Embedding Geo-Entities for Semantic Typing
 - From Digitized Reports to Spatio-Temporal KGs
- Conclusions & Future Directions

	All	Since 2019
Citations	110	110
h-index	6	6
i10-index	4	4



My PhD Journey



I take **different kinds of data** sources then instill them with **semantics**



ISWC 2022, 23, 24

WWW 2023

KDD 2024



GE Research 2021

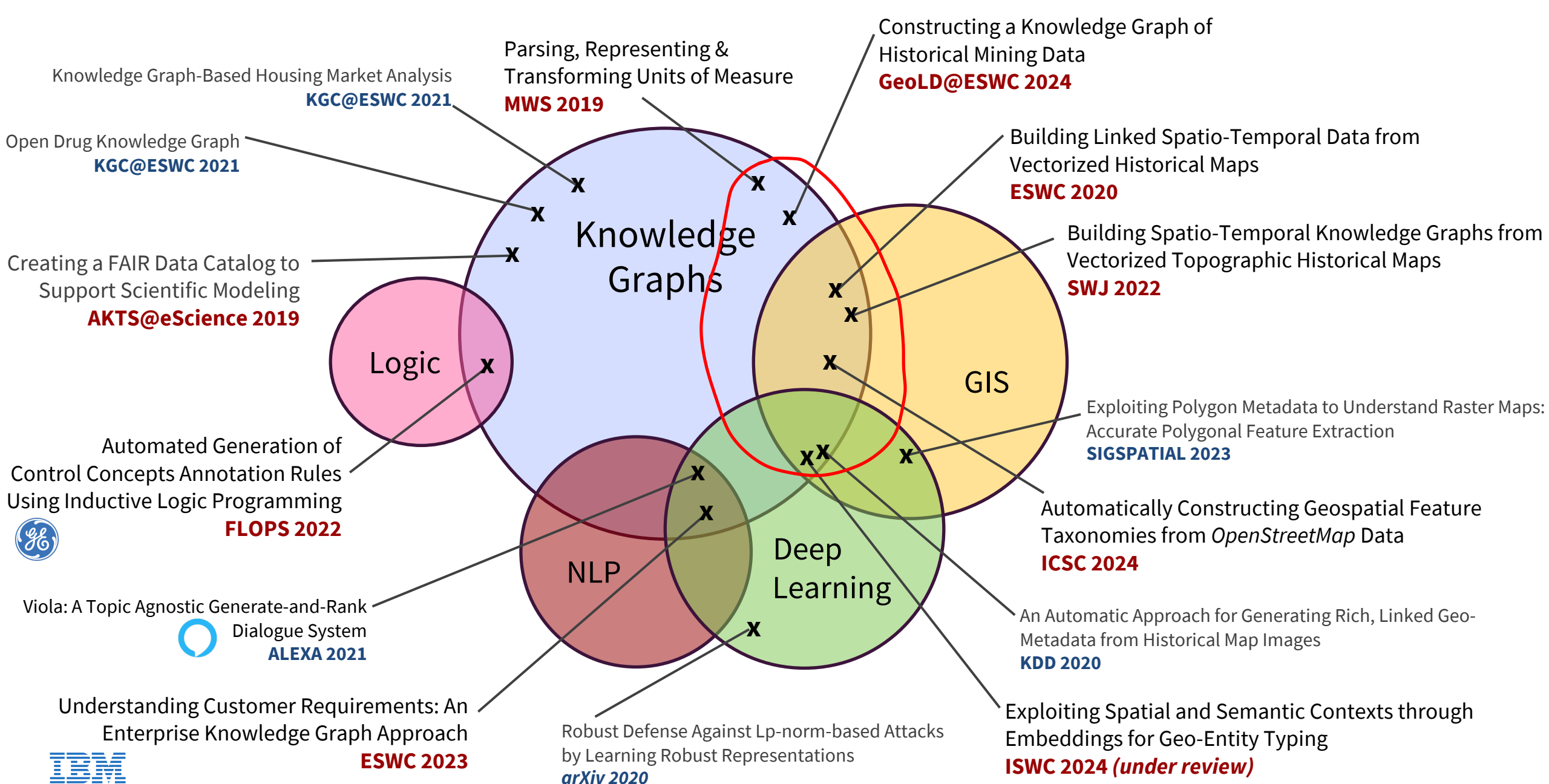
IBM Research 2022




semifinalists 2021



University Outstanding Teaching Assistant Award 2022
Highest Achievement

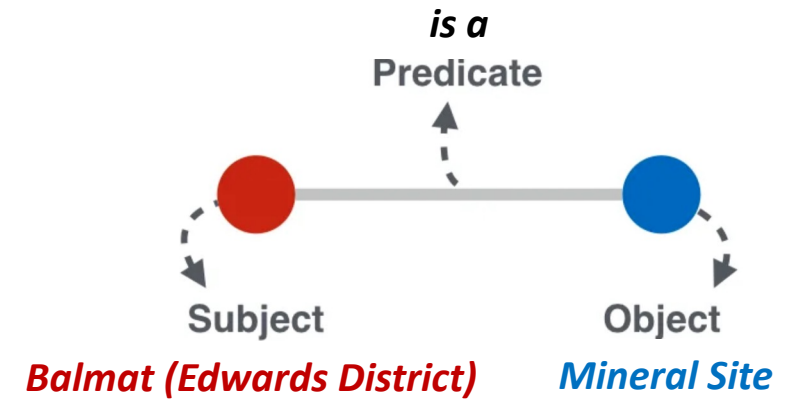


Agenda

- Basel's PhD Journey
- Intro 
- Thesis Overview
- Approach:
 - Building Spatio-Temporal KGs from Digitized Maps
 - Embedding Geo-Entities for Semantic Typing
 - From Digitized Reports to Spatio-Temporal KGs
- Conclusions & Future Directions

KGs

- Knowledge Graph (KG)?
 - **Graphs** are natural way to **encode** data
 - Using **semantic concepts & relationships**
 - Semantic Network = **Knowledge Graph**
- Why use KG?
 - Combine **expressivity, interoperability, & standardization** in the **Semantic Web** stack
- Semantic Web?
 - Extension of WWW, enabling the Web of Data (aka “Linked Data”)
 - Encoding of **semantics** with the data
 - Linked Open Data **principles** // FAIR



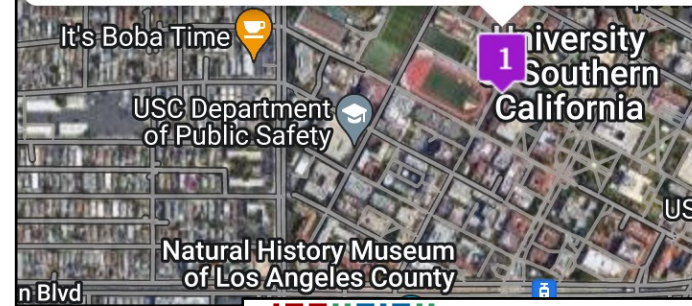
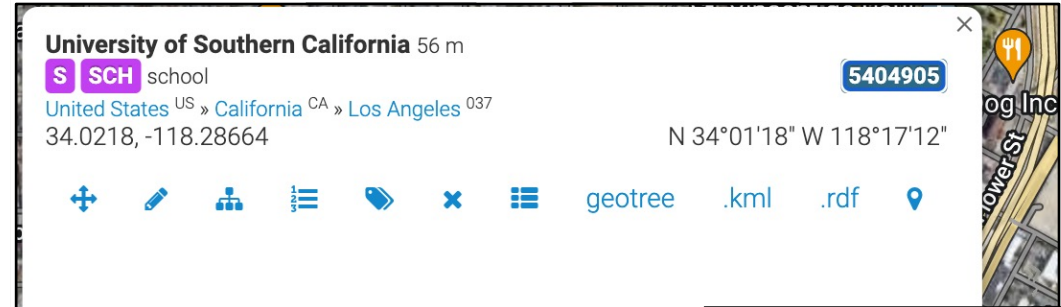
Geo & Spatio-Temporal KGs

- Spatio-Temporal KGs


- Contextual (**what**)
- **Spatial** (**where**)
- **Temporal** (**when**)

- Geo-**semantics**

- Representation, annotation, & reasoning
- Modeling & ontology development
- Integration & interoperability



DBpedia	
dbo:foundingDate	• 1880-10-06 (xsd:date)
dbo:mascot	• Tommy Trojan (unofficial) • Traveler
dbo:motto	• Palmam qui meruit ferat • "Let whoever earns the palm bear it"

 Wikidata

University of Southern California (Q4614)

private university in Los Angeles, California, United States
USC | University of Southern CA

Way: University of Southern California (161235655)

Version #25

Intro: Spatio-Temporal KGs

- So, what's so **special** about them?
 - Spatial analysis
 - Temporal analysis
 - Spatio-temporal aggregations
 - Geographic QA
 - Environmental & social science
 - Urban planning
 - Transportation
 - etc...

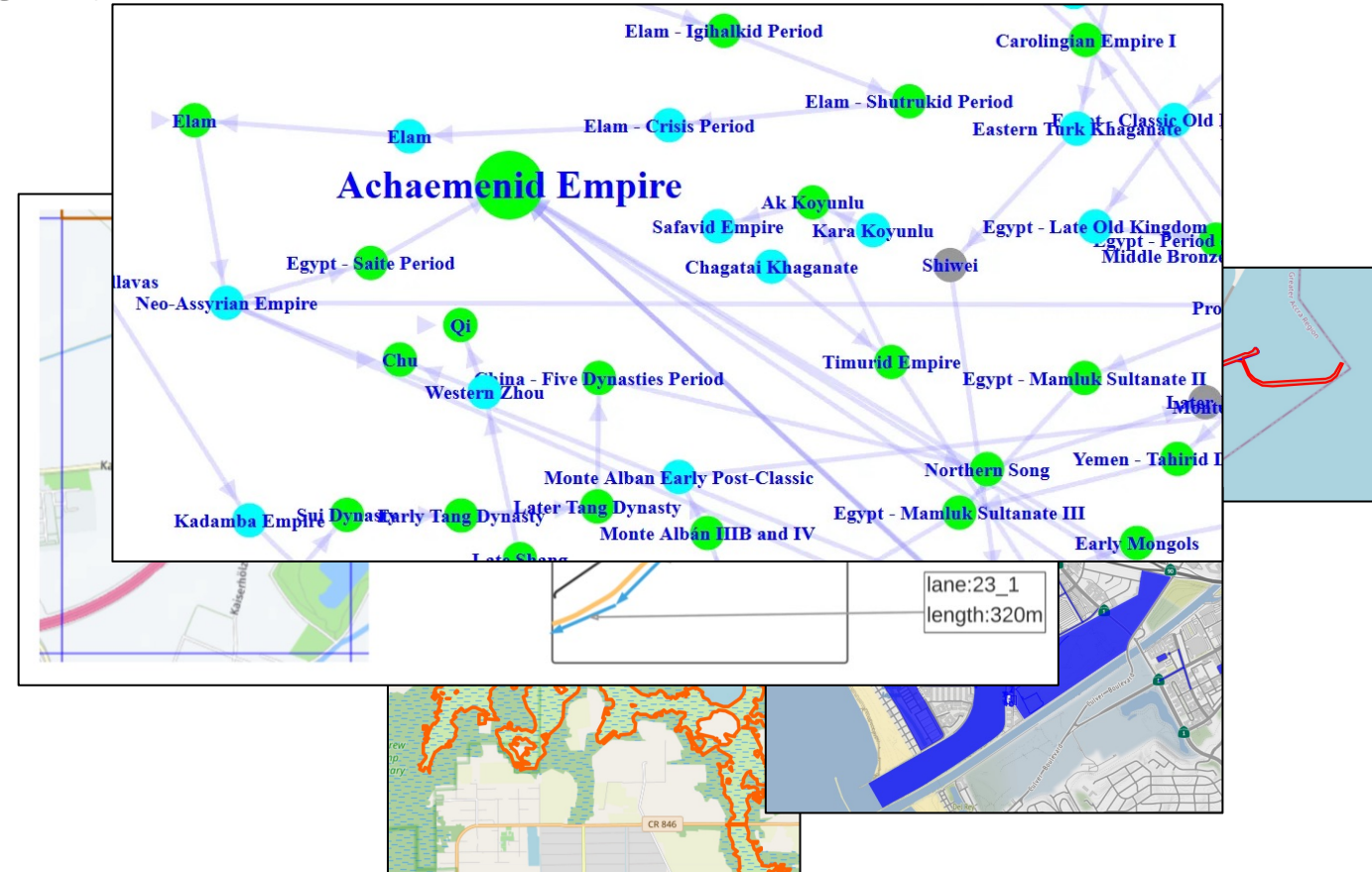



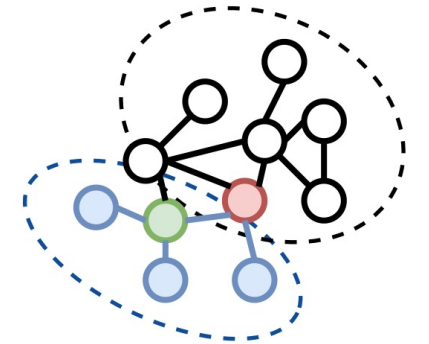
figure from Szwoch, G. (2019). Combining road network data from openstreetmap with an authoritative database. *Journal of Transportation Engineering, Part A: Systems*, 145(2), 04018085.
figure from <https://terminusdb.com/blog/human-history-knowledge-graph/>

Agenda

- Basel's PhD Journey
- Intro
- Thesis Overview 
- Approach:
 - Building Spatio-Temporal KGs from Digitized Maps
 - Embedding Geo-Entities for Semantic Typing
 - From Digitized Reports to Spatio-Temporal KGs
- Conclusions & Future Directions

Research Problem

- How can we **transform & link unstructured digitized & historical** geo-data into **structured, semantic, & queryable spatio-temporal KGs**?
- Objectives:
 - **Automated KG construction** from various historical geo-data sources
 - Contextual **geo-entity** recognition (**ER**) & **semantic typing**
 - **Semantic enrichment** by **linking** (**EL**) to additional sources on the web
 - Adherence to **Semantic Web principles**
 - shared, **accessible**, visualized, standardized **across-domains**, & **scalable** for easy use by downstream tasks for easy analysis & expressive **integration**



Thesis statement


This thesis provides **tools & techniques** for the automated **understanding & transformation** of **unstructured geospatial & historical data** from heterogeneous sources into a standardized representation of expressive & interoperable **spatio-temporal knowledge graphs**. I also present methodologies that both **integrate** the data **with other sources on the web** and **leverage web knowledge** for enhanced data analysis.

Building Spatio-Temporal KGs from Digitized Maps

Embedding Geo-Entities for Semantic Typing

From Digitized Reports to Spatio-Temporal KGs

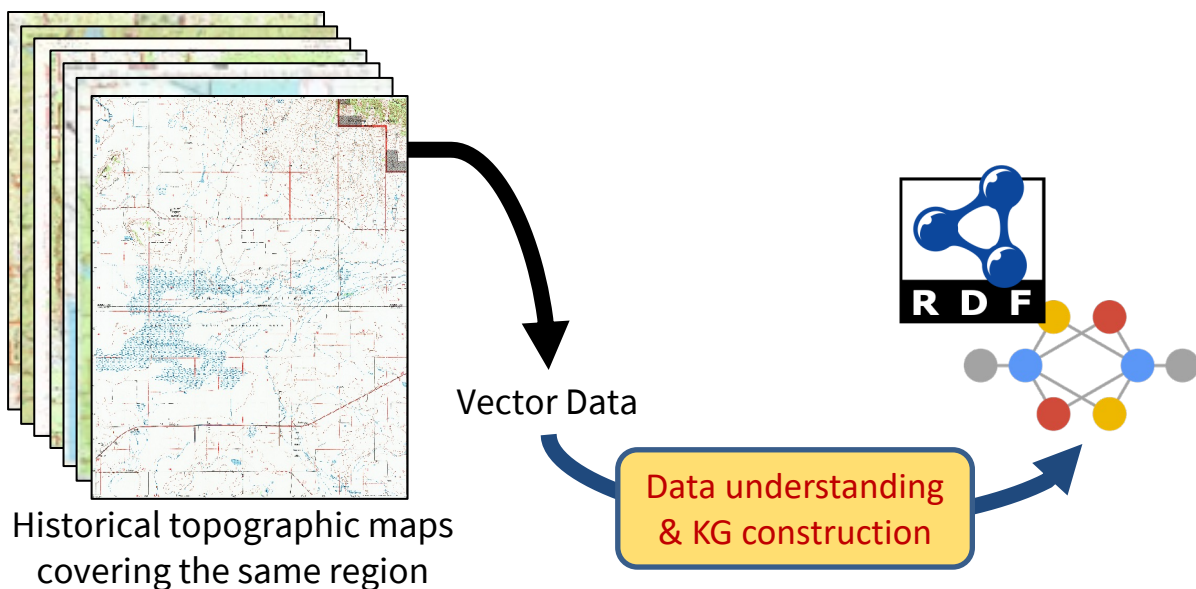
Agenda

- Basel's PhD Journey
- Intro
- Thesis Overview
- Approach:
 - Building Spatio-Temporal KGs from Digitized Maps 
 - Embedding Geo-Entities for Semantic Typing
 - From Digitized Reports to Spatio-Temporal KGs
- Conclusions & Future Directions

Building Spatio-Temporal KGs from Digitized Maps

- Goals

- Automatically **integrate**, **represent**, **relate** & **interlink** geospatial data from overlapping digitized resources
 - Line & Polygon features (e.g. **railroads**, **wetlands**)
 - Geo-semantic **representation** following LD & SW principles
- User-assisted (**known type**) geo-entity linking (**EL**)



Building Spatio-Temporal KGs from Digitized Maps

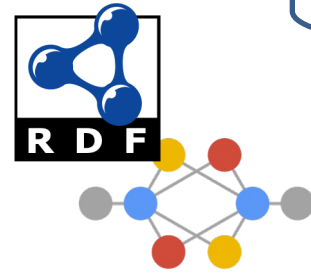
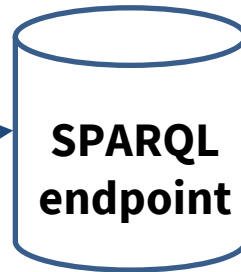


Way: Long Bell Lumber Company Railroad (177559134)
Version #3
Tags

name	Long Bell Lumber Company Railroad
railway	abandoned
tiger:cfcc	B11
tiger:county	Siskiyou, CA
tiger:name_base	Long Bell Lumber Company RR

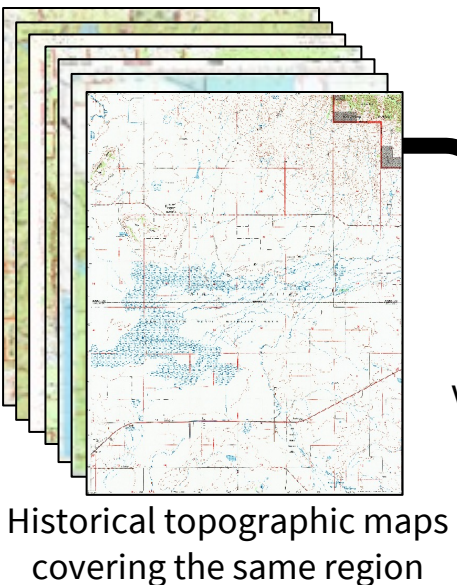
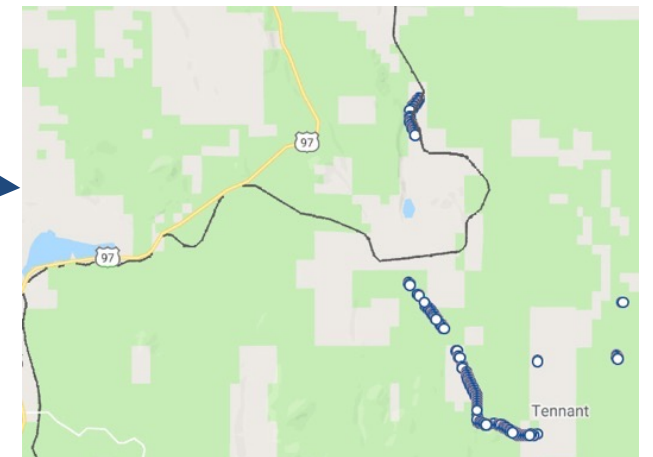
“can you show me a subset of instances tagged as **abandoned**?”

“**railroad** segments that are **present** in 1962 but are **not present** in 2001”



Data understanding & KG construction

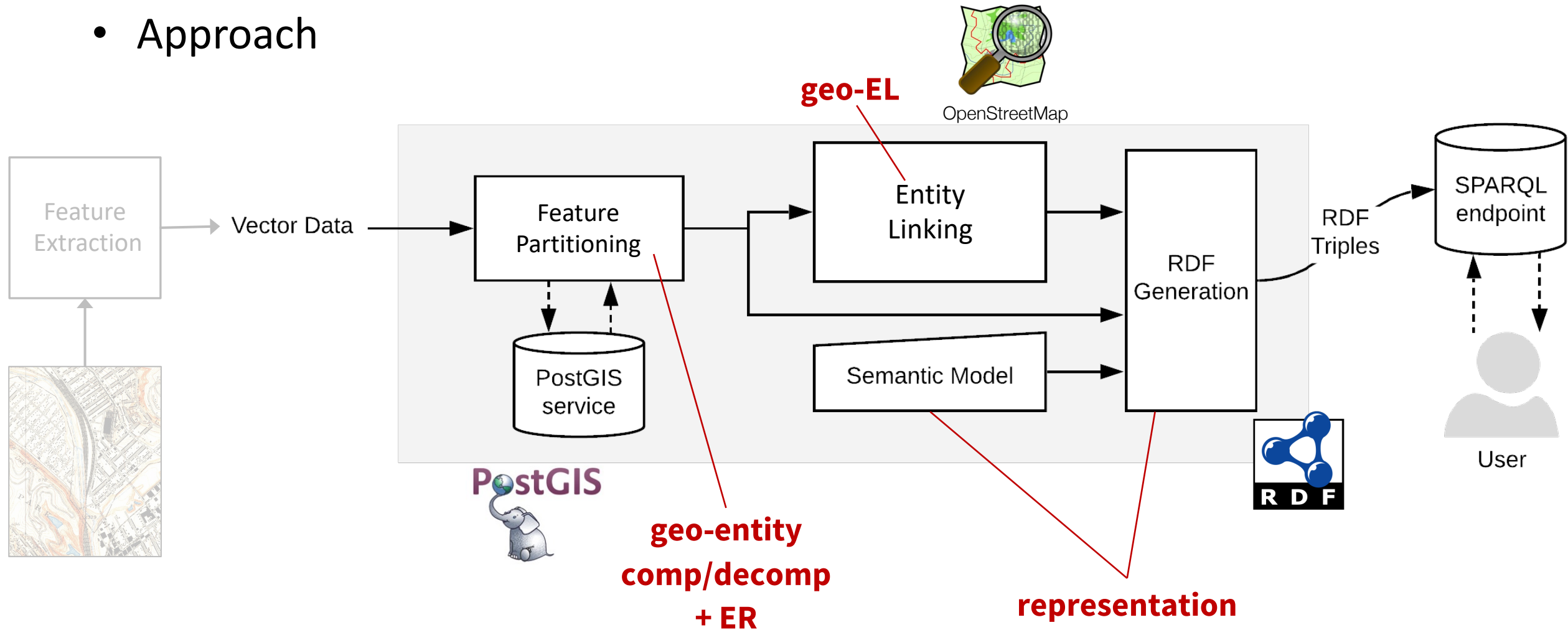
Vector Data



Historical topographic maps covering the same region

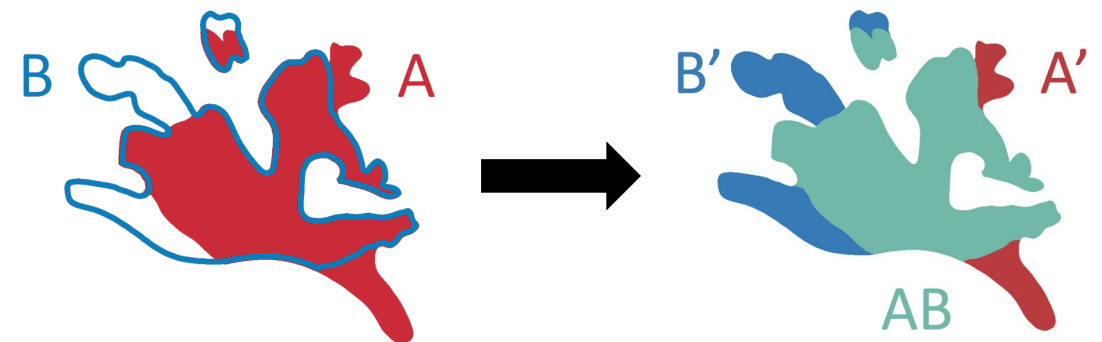
Building Spatio-Temporal KGs from Digitized Maps

- Approach



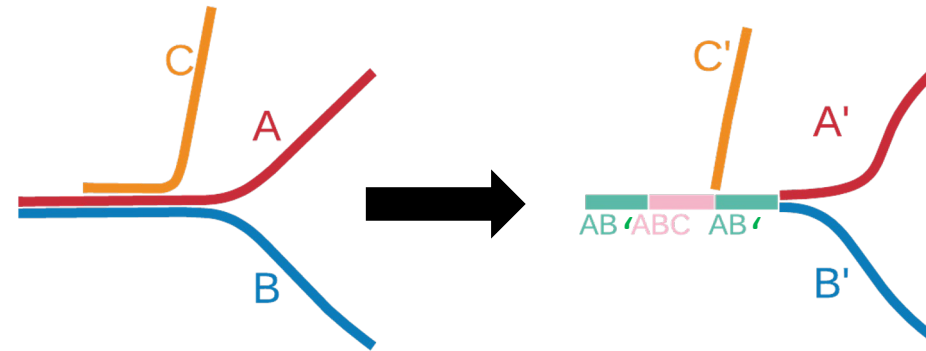
Building Spatio-Temporal KGs from Digitized Maps

- Approach
 - Step 1. Feature Partitioning
 - Generate **building block** geometries (i.e. geo entities) to **represent** the topographic **features** from different map sheets
 - Represent **common & distinct parts** (changes) of the geo-features
 - Allow **incremental** additions over time
 - **Create a DAG** of building-block **geometries** (nodes) and their **relations** (edges)



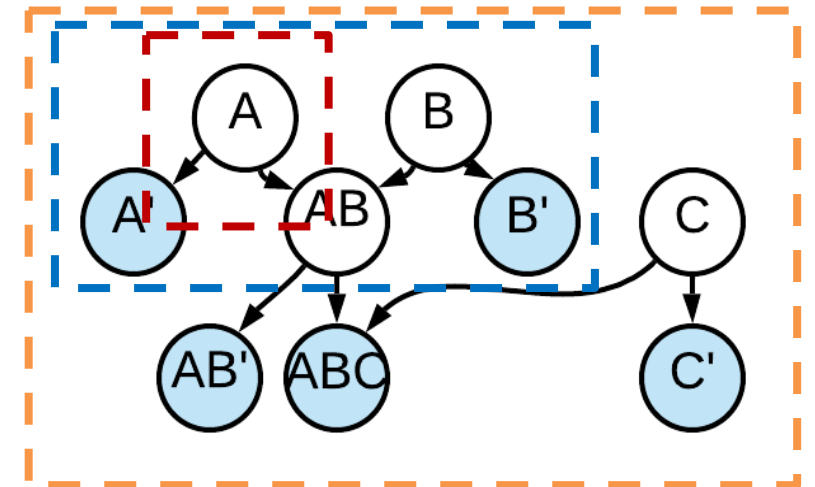
Building Spatio-Temporal KGs from Digitized Maps

- Approach
 - Step 1. Feature Partitioning



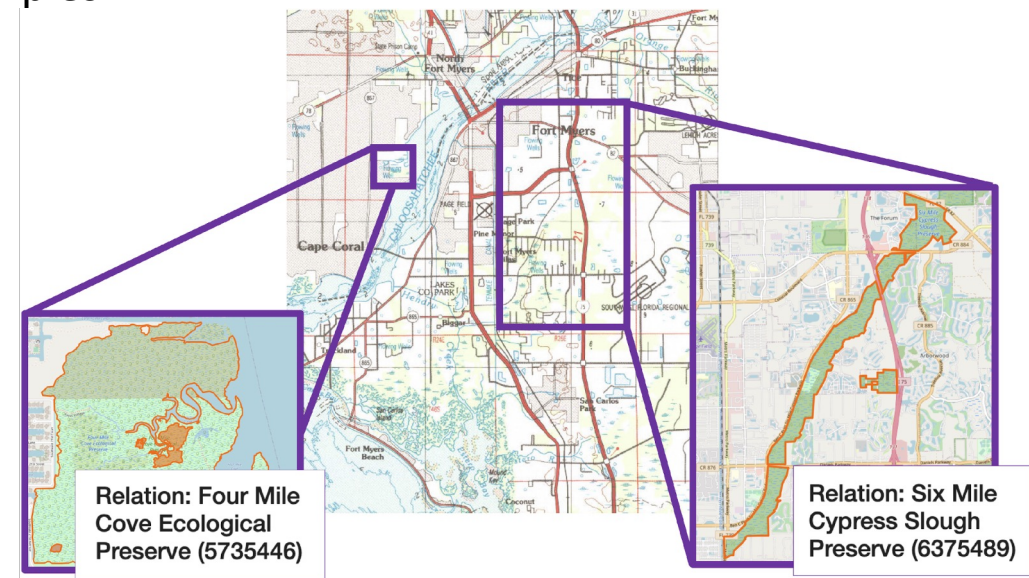
```

foreach  $i \in \mathcal{M}$  do
  foreach  $k \in \mathcal{L}$  do
     $\mathcal{F}_\alpha = \mathcal{F}_i \cap \mathcal{F}_k;$ 
     $\mathcal{F}_\gamma = \mathcal{F}_k \setminus \mathcal{F}_\alpha;$ 
  end
   $\mathcal{F}_\delta = \mathcal{F}_i \setminus (\bigcup_{j \in \mathcal{L}} \mathcal{F}_j);$ 
end
    
```



Building Spatio-Temporal KGs from Digitized Maps

- Approach
 - Step 2. Geo-entity Linking
 - Link the generated entities to **LoD** (OSM, geoNames, LinkedGeoData, Wikidata)
 - **Sampling**
 - **Reverse geocoding** for initial filtering (**geo feature type is known**)
 - Determine confidence by **frequency** of (random) samples

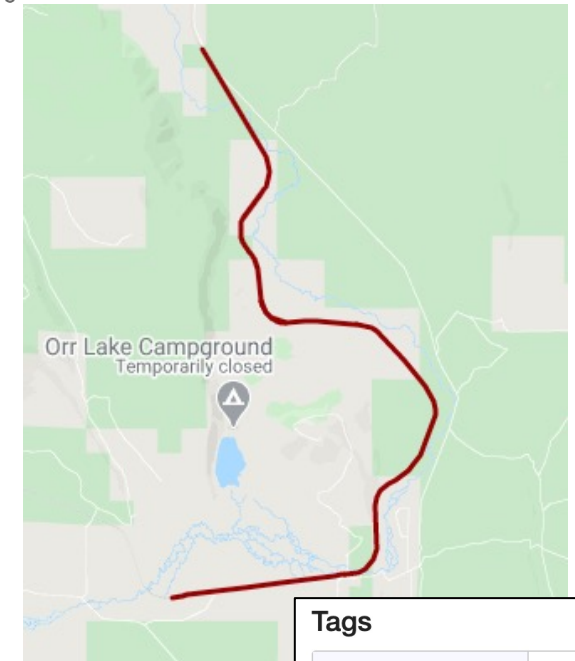
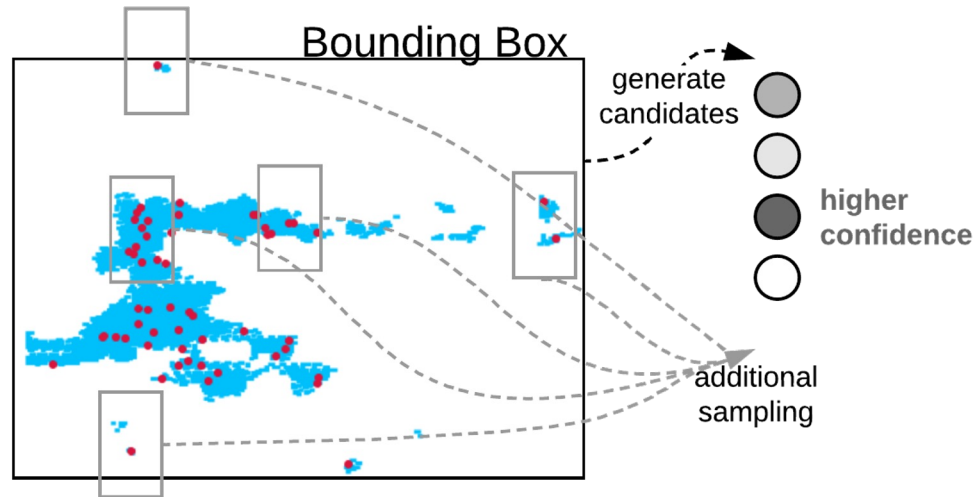
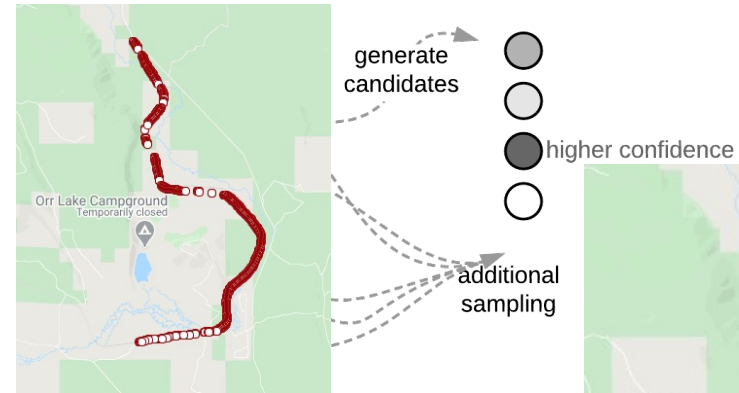


Building Spatio-Temporal KGs from Digitized Maps

- Approach
 - Step 2. Geo-entity Linking

```

 $B_s$  = bounding box wrapping  $s$ ;
 $\mathcal{L}$  = reverse-geocoding( $B_s, T$ );
for 1... $N$  do
     $e$  = randomly sample a Point in segment  $s$ ;
     $E$  = reverse-geocoding( $e, T$ );
     $\mathcal{L}$ .add( $E$ );
end
    
```



Way: Black Butte Subdivision (322131253)

Version #6

Tags	
name	Black Butte Subdivision
operator	Union Pacific Railroad;Amtrak
owner	Union Pacific Railroad
railway	rail

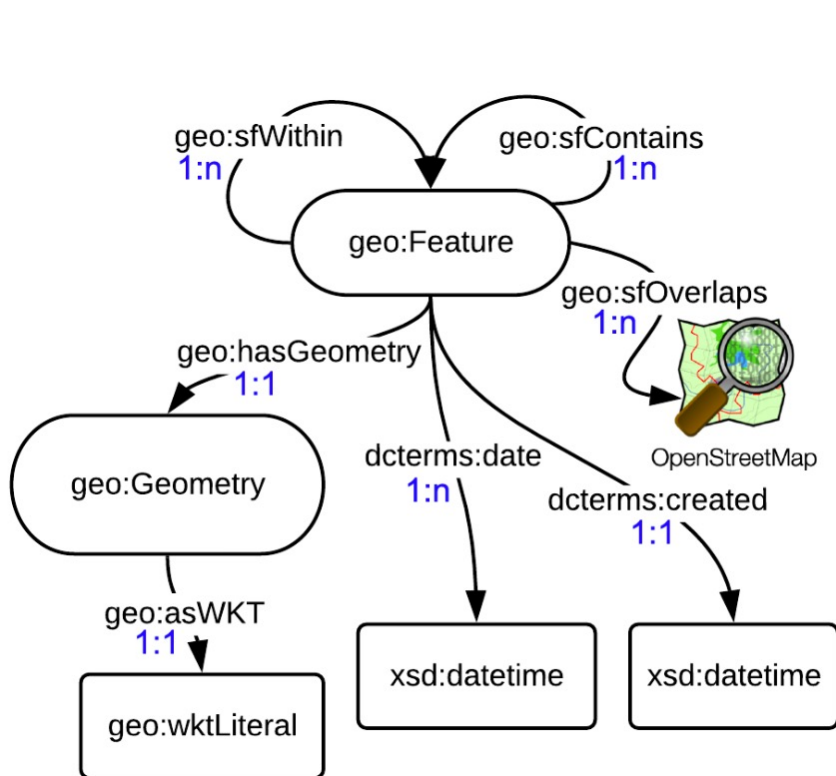
Building Spatio-Temporal KGs from Digitized Maps

- Approach
 - Step 3. Semantic Modeling (Data Representation)
 - Transform & materialize the data (construct KG)
 - Follows **linked data principles**
 - Provide a **useful semantic representation** supporting downstream tasks
 - Construct a **meaningful semantic model**
 - **Hierarchically-driven**
 - Follows W3C & OGC **standards (GeoSPARQL)**



Building Spatio-Temporal KGs from Digitized Maps

- Approach
 - Step 3. Semantic Modeling (Data Representation)



The top part shows an RDF graph with **geo:Feature** as the root node. It has outgoing edges for **uri**, **geo:asWKT**, **dcterms:date**, **geo:sfWithin**, **geo:sfContains**, **dcterms:created**, and **geo:sfOverlaps**.

gid	wkt	year	predecessor_id	successor_id	time_generated	OSM_uri
http://linkedmaps.isi.edu/13	MULTIPOLYGON((... 41.123592071033... 41.123592362641... 41.120256913283...))	1961	http://linkedmaps.isi.edu/5	http://linkedmaps.isi.edu/12	2021-02-10T15:23:23	https://www.openstreetmap.org/way/86302769

```
SELECT ?f ?wkt WHERE {
  ?f a geo:Feature ;
  geo:hasGeometry [ geo:asWKT ?wkt ] ;
  dcterms:date 1962^^xsd:gYear .
  FILTER NOT EXISTS { ?f geo:sfContains__:_ }
  MINUS { ?f dcterms:date 2001^^xsd:gYear . }
```

Fig. 16. Query feature geometries present in 1962 but not in 2001

Building Spatio-Temporal KGs from Digitized Maps

- Evaluation

- Feature Partitioning
- Geo EL
- RDF Query Performance

Data: 2 x railroad, 3 x wetland (3-7 sheets)

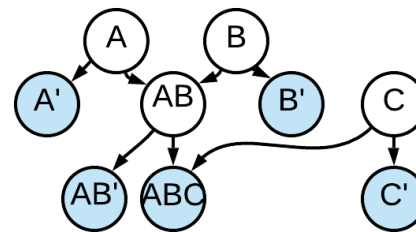
Geo-entity linking results; Area is in square kilometers

		Area	Precision	Recall	F_1
Railroads	CA-baseline		0.193	1.000	0.323
	CA	420.39	0.800	0.750	0.774
	CO-baseline		0.455	1.000	0.625
	CO	132.01	0.833	1.000	0.909
Wetlands	CA-baseline		0.556	1.000	0.714
	CA	224.05	1.000	1.000	1.000
	FL-baseline		0.263	1.000	0.417
	FL	27493.98	0.758	0.272	0.400
	TX*	16.62	-	-	-

Year	# vecs	Runtime (s)	# nodes
1954	2382	<1	1
1962	2322	36	5
1988	11134	1047	11
1984	11868	581	24
1950	11076	1332	43
2001	497	145	57
1958	1860	222	85

Year	# vecs	Runtime (s)	# nodes
1961	12	<1	1
1993	17	<1	5
1990	27	6	11
2018	9	6	24

Year	# vecs	Runtime (s)	# nodes
1987	184	<1	1
1956	531	180	5
2020	5322	1139	13



Query time ranges **10-50 [ms]**
 No significant change with respect to

- # of map editions we process
- Complexity of the query we compose

Building Spatio-Temporal KGs from Digitized Maps

- Related Work

- Transforming geospatial vector data into RDF (Kyzirakos 2014, Usery 2012)
 - Do not address **geo-entity linking** or **entity semantics**
- Contextualizing geospatial data (Vaisman 2019, Smeros 2016)
 - Do not address Linking **unlabeled** geo entities
- Geospatial change analysis (Perez 2015, Kauppinen 2014)
 - Do not address **incremental process** of change over time

Building Spatio-Temporal KGs from Digitized Maps

- Takeaways

- Paradigm for **geospatial data integration** on the web

- **Unsupervised**

- **Hierarchy-driven & incremental** semantic model for simple & efficient **querying**

- Follows LD & SW principles

- Does not require re-generation of data

- URIs are preserved

- **Linked to Web**

- Fuels further discovery & enrichment

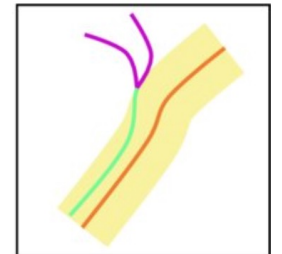
- Still, many challenges exist

- **Complexity** of changes in original topographic maps

- Quality & **level of detail**

- Crowdsourcing: **availability, granularity** (e.g., mud vs. wetland)


- **User-informed** knowledge (target schema)



<https://linked-maps.isi.edu>

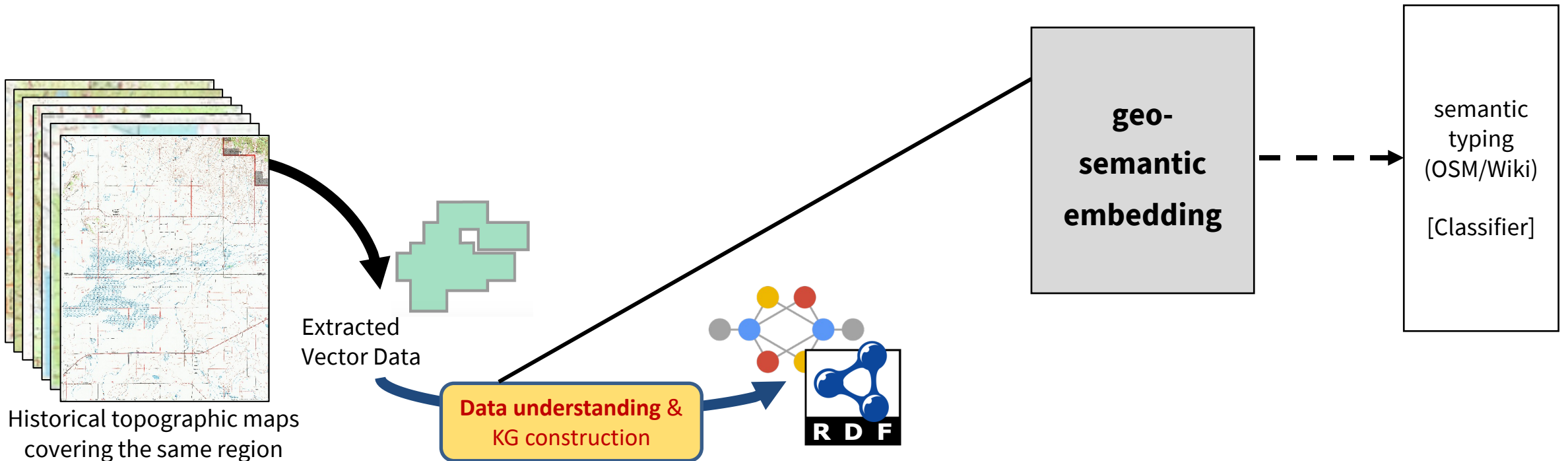
<https://github.com/usc-isi-i2/linked-maps>

Agenda

- Basel's PhD Journey
- Intro
- Thesis Overview
- Approach:
 - Building Spatio-Temporal KGs from Digitized Maps
 - Embedding Geo-Entities for Semantic Typing 
 - From Digitized Reports to Spatio-Temporal KGs
- Conclusions & Future Directions

Embedding Geo-Entities for Semantic Typing

- Goals:
 - Embed geospatial data into **high-dimensional vector space**
 - Preserve its **semantic meaning** & relationships between entities
 - Develop techniques for **semantic typing/labeling** of geospatial entities



Embedding Geo-Entities for Semantic Typing



waterway



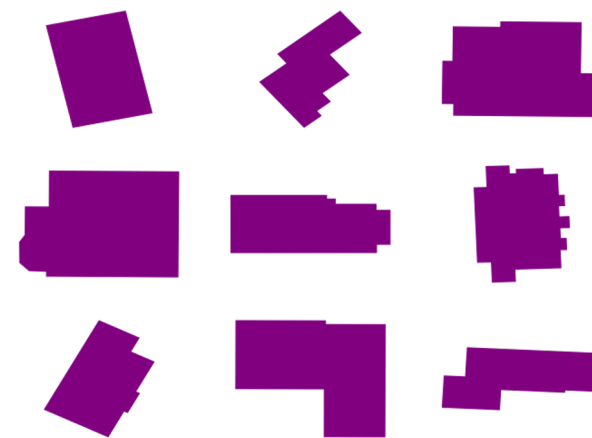
- waterway
 - canal [5748]
 - dam [2425]
 - ditch [37754]
 - drain_waterway [7044]
 - river_waterway [16925]
 - stream [799887]



water



- water
 - basin [1870]
 - lake [3214]
 - pond [6835]
 - reservoir [4176]
 - river_water [1570]



building



- building
 - apartments [31601]
 - commercial [3875]
 - house [118030]
 - industrial [3223]
 - residential [19763]
 - retail_building [4109]
 - school [2400]
 - warehouse [1098]

*“Everything is related to everything else.
But near things are more related than distant things.”*

-Waldo R. Tobler

Embedding Geo-Entities for Semantic Typing

- Approach:

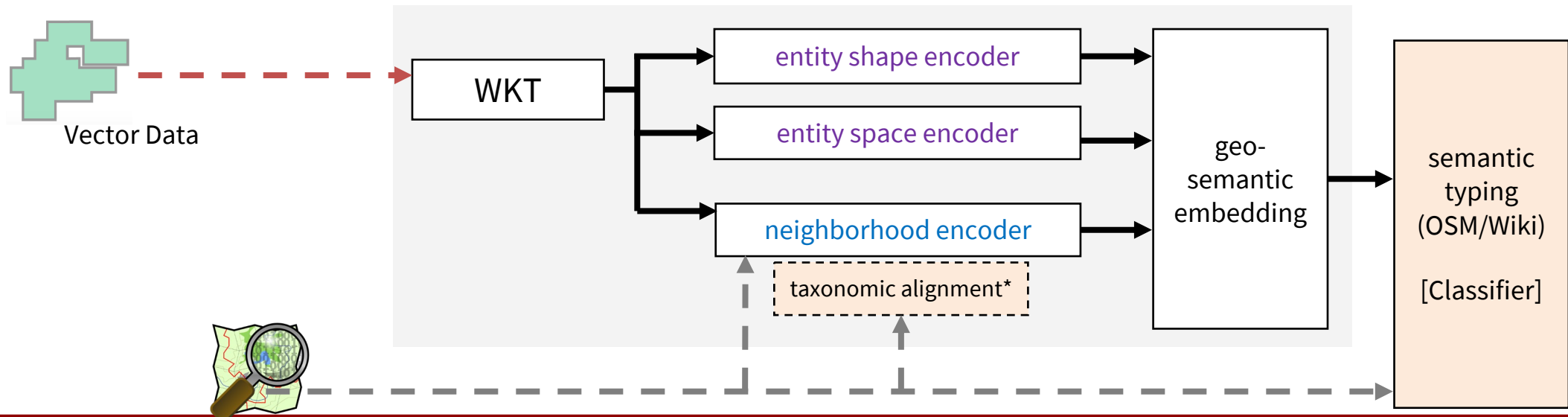
- Method to **embed**:

- Geometric attributes (**shape**)
 - Spatial attributes (**area, length**)
 - Neighborhood context (**nearby geo-entities**)

geospatial

semantic

to generate a **representation** that can learn & infer properties about geo-entities



Embedding Geo-Entities for Sem

- Data



OpenStreetMap

Relation: 10052899

Version #1

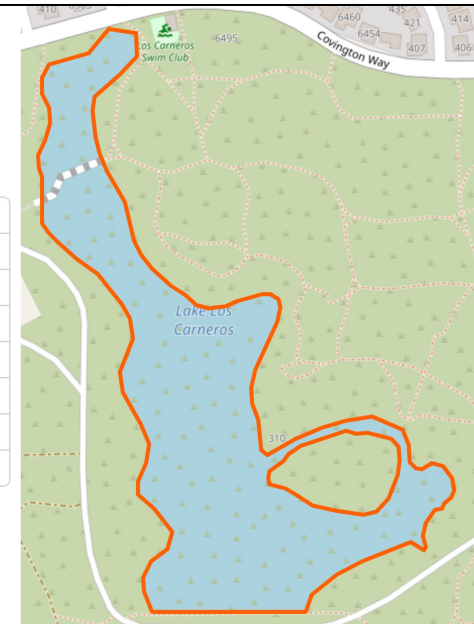
Changeset #74650407

Tags

ele	22
gnis:county_id	083
gnis:created	06/13/2000
gnis:feature_id	1871851
gnis:state_id	06
type	multipolygon
natural	water
water	reservoir

Members

- ▼ 2 members
- Way 23145279 as outer
- Way 726021752 as inner



- Nodes** - dots used to mark locations
- Ways** - connected line of nodes
- Relation** - used to create more complex shapes



CA OSM Snapshot

index		
0	node_tagged	1000170
1	node_untagged	133590505
2	way_untagged	4029438
3	way_tagged	9017322
4	relation_untagged	84613
5	relation_tagged	68462

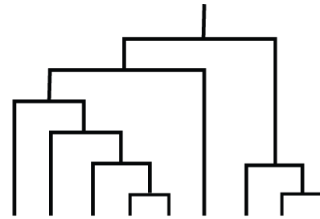
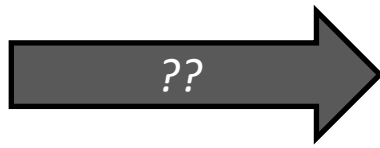
shp_type		
0	Polygon	4876318
1	LineString	4176529
2	Point	1000170
3	MultiPolygon	12694
4	GeometryCollection	1364

Embedding Geo-Entities for Semantic Typing

- but, OSM data is
 - **Inconsistent** across regions
 - Varying-**granularity**
 - **Noisy**



OpenStreetMap
data/dump



sidewalk (Q177749)

pedestrian path along the side of a road
pavement | footpath | footway | platform

Statements

subclass of thoroughfare

Statements

subclass of public space
line construction
axis of communication
geographical feature

OSMonto - An Ontology of OpenStreetMap Tags

Mihai Codescu*, Gregor Horsinka*, Oliver Kutz**
Till Mossakowski***, Rafaela Rau*

* DFKI GmbH Bremen, Germany
** Research Center on Spatial Cognition,
SFB/TR 8, University of Bremen, Germany

OUTDATED

Embedding Geo-Entities for Semantic Typing

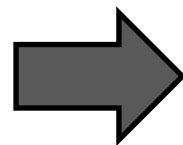
- **Auxiliary** (Appendix A):
 - Generate a lightweight **taxonomy** from OSM **tag data**



```

<way id="232250107" visible="true" vers
2019-05-06T23:22:23Z" user="Enock4seth"
<nd ref="5058536215"/>
<nd ref="1797433673"/>
<nd ref="4992821222"/>
<tag k="highway" v="tertiary"/>
<tag k="name" v="Nana Kana Street"/>
</way>
<way id="244376453" visible="true" vers
2015-04-02T14:55:17Z" user="sidneys" ui
<nd ref="2517024878"/>
<nd ref="2517024879"/>
<nd ref="2517024880"/>
<nd ref="2517024881"/>
<nd ref="2517024878"/>
<tag k="building" v="industrial"/>
</way>
<way id="244376454" visible="true" vers
2015-04-02T13:43:25Z" user="sidneys" ui
<nd ref="2517024878"/>

```



construct base terminology

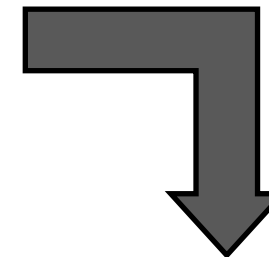
frequent non-informative
infrequent informative

```

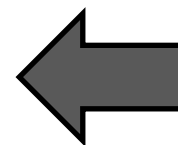
{'apartments',
'building',
'driveway',
'highway',
'house',
'residential',
'service'}

```

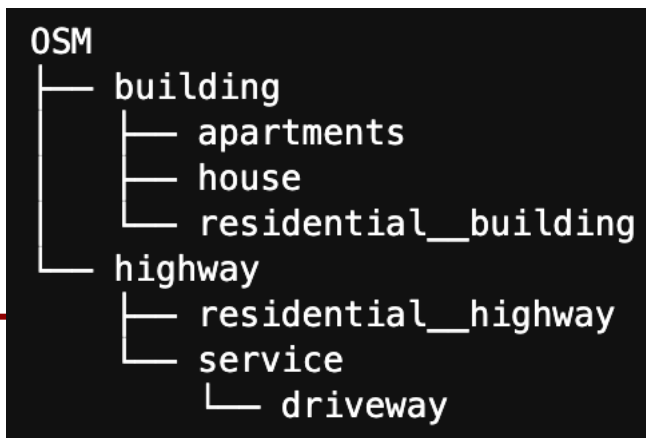
count parent-child relations
path frequency



parent	child	counter
building	house	15
highway	service	14
building	residential	33
highway	residential	22
building	apartments	2
service	driveway	5



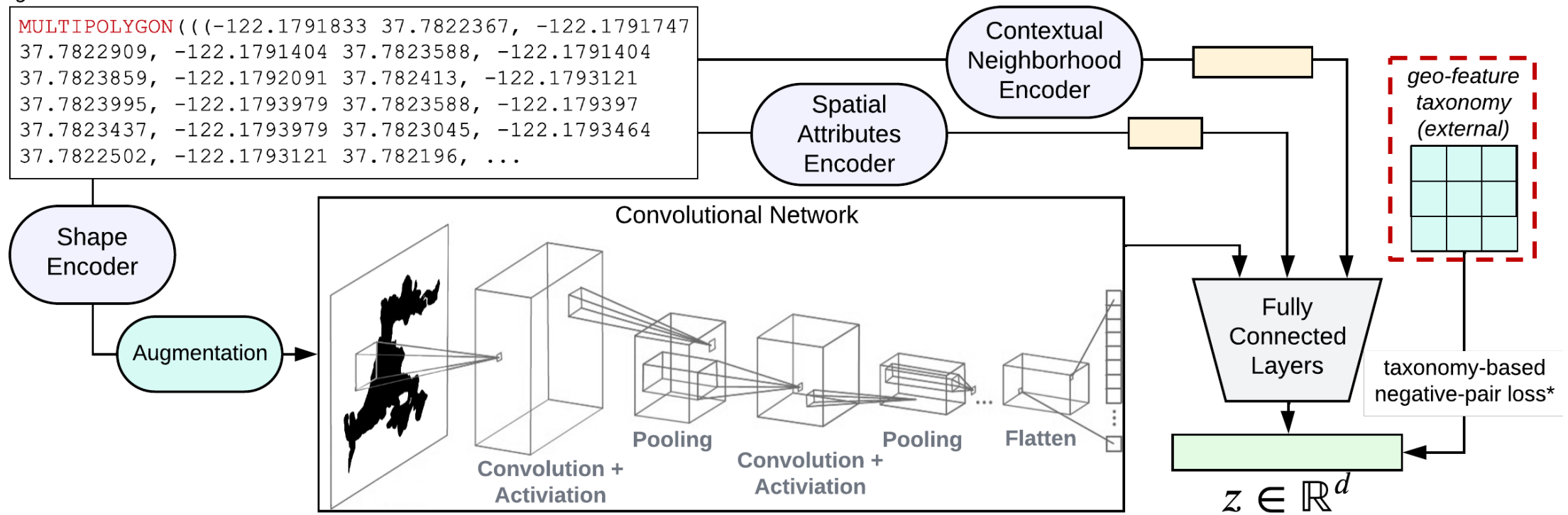
build taxonomy
conflict resolution



Embedding Geo-Entities for Semantic Typing

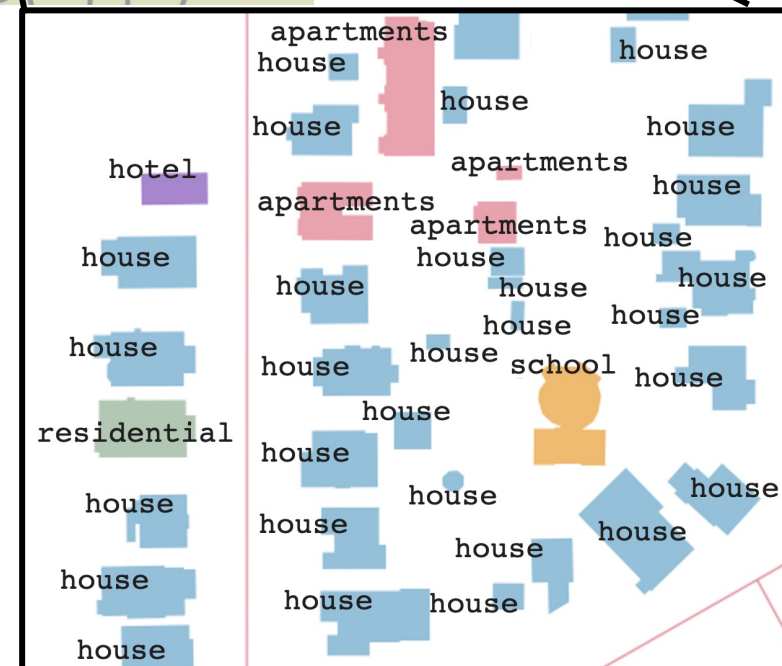
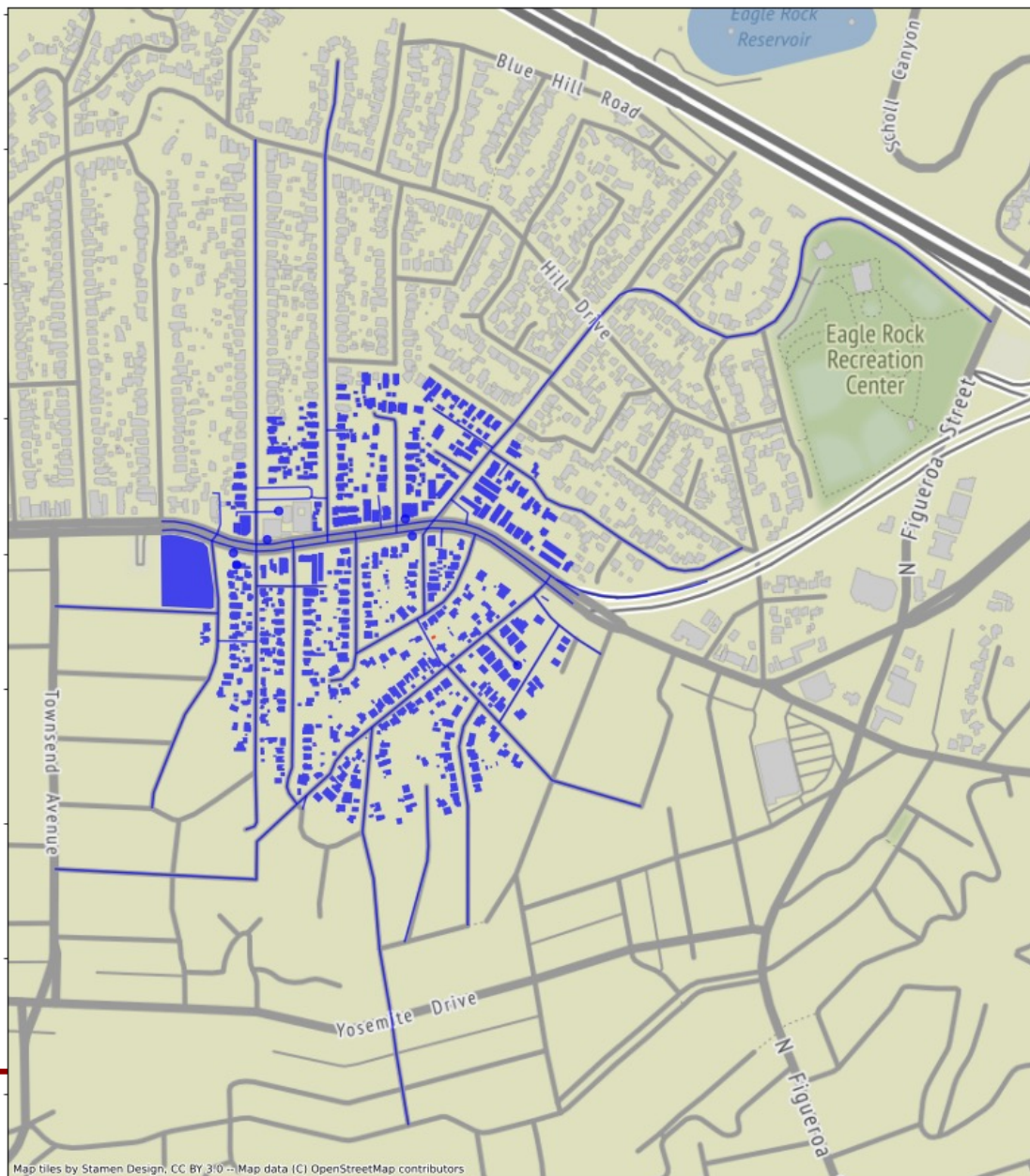
geo-referenced vector data

```
MULTIPOLYGON (((-122.1791833 37.7822367, -122.1791747
37.7822909, -122.1791404 37.7823588, -122.1791404
37.7823859, -122.1792091 37.782413, -122.1793121
37.7823995, -122.1793979 37.7823588, -122.179397
37.7823437, -122.1793979 37.7823045, -122.1793464
37.7822502, -122.1793121 37.782196, ...
```



Embedding Geo-Entities for Semantic Typing

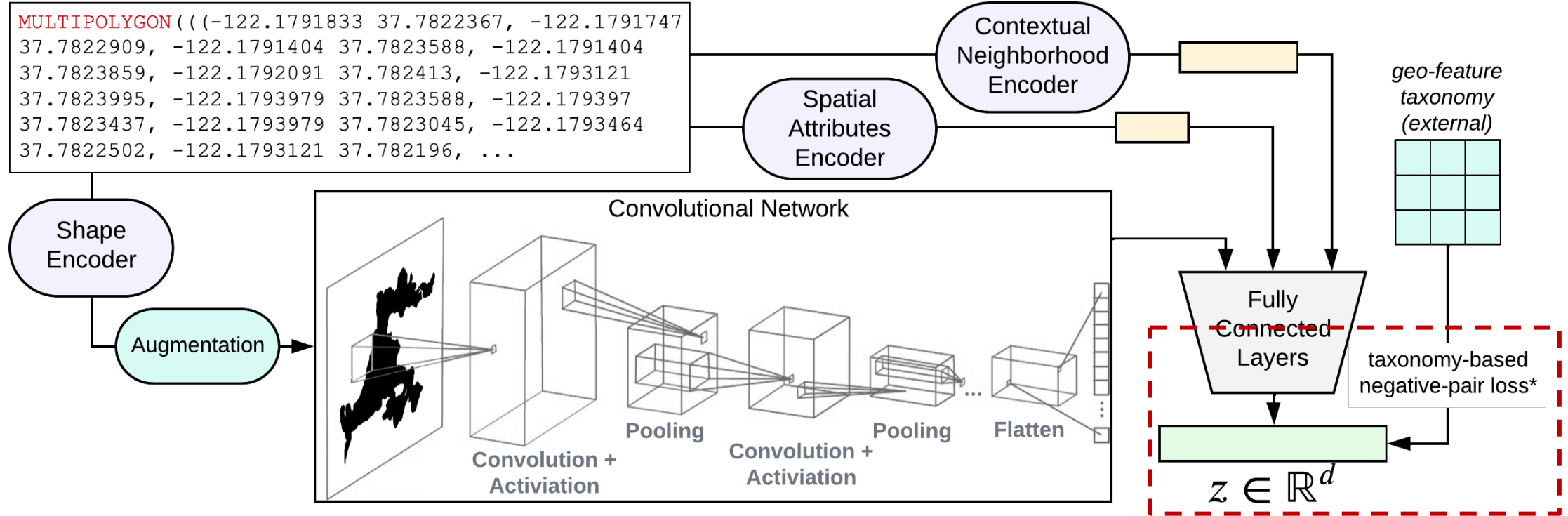
house	414
residential_building	57
primary	29
apartments_building	28
residential_highway	21
service	11
commercial_building	8
retail_building	6
hotel_tourism	6
platform	4
alley	3
hotel_building	2
school_building	2
tertiary	2
industrial_building	1
steps	1
warehouse	1
restaurant	1
turning_circle	1
motorway_link	1
place_of_worship	1
driveway	1
school_amenity	1



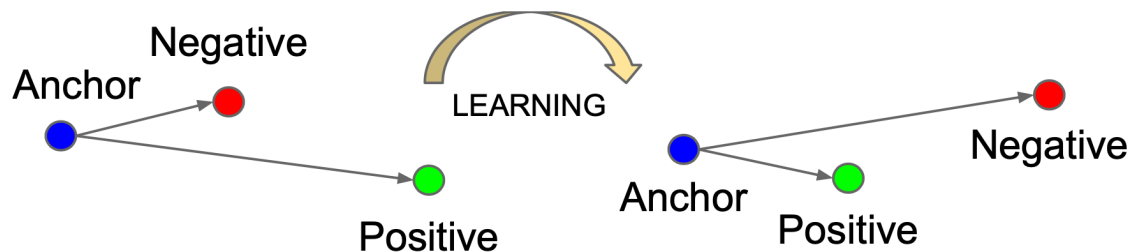
Embedding Geo-Entities for Semantic Typing

geo-referenced vector data

```
MULTIPOLYGON (((-122.1791833 37.7822367, -122.1791747
37.7822909, -122.1791404 37.7823588, -122.1791404
37.7823859, -122.1792091 37.782413, -122.1793121
37.7823995, -122.1793979 37.7823588, -122.179397
37.7823437, -122.1793979 37.7823045, -122.1793464
37.7822502, -122.1793121 37.782196, ...
```



Embedding Geo-Entities for Semantic Typing



$$L_q = -\log \frac{\exp(\text{sim}(e_q, e_+)/\tau)}{\sum_{i=0}^K \exp(\text{sim}(e_q, e_i) \cdot w_{q,i}/\tau)}$$

$$w_{i,j} = \frac{d_{tree} - d_{i,j}}{d_{tree}}$$

Amenity	Fast Food	0.000	0.667	0.667	0.667	0.667	1.000	1.000	1.000	1.000
	Parking	0.667	0.000	0.667	0.667	0.667	1.000	1.000	1.000	1.000
	Mosque	0.667	0.667	0.000	0.333	0.333	1.000	1.000	1.000	1.000
	Synagogue	0.667	0.667	0.333	0.000	0.333	1.000	1.000	1.000	1.000
	Church	0.667	0.667	0.333	0.333	0.000	1.000	1.000	1.000	1.000
Residential Building	Apartments	1.000	1.000	1.000	1.000	1.000	0.000	0.667	1.000	1.000
	House	1.000	1.000	1.000	1.000	1.000	0.667	0.000	1.000	1.000
Highway	Cycleway	1.000	1.000	1.000	1.000	1.000	1.000	1.000	0.000	0.667
	Footway	1.000	1.000	1.000	1.000	1.000	1.000	1.000	0.667	0.000
	Sidewalk	1.000	1.000	1.000	1.000	1.000	1.000	1.000	0.667	0.000
	Fast Food									
	Parking									
	Mosque									
	Synagogue									
	Church									
	Apartments									
	House									
	Cycleway									
	Sidewalk									

Embedding Geo-Entities for Semantic Typing

- Evaluation

- 8-fold SVC on embeddings

- Data: 2k+ instances → 11 WD classes

- 16k+ instances → 18 OSM tags

4 settings



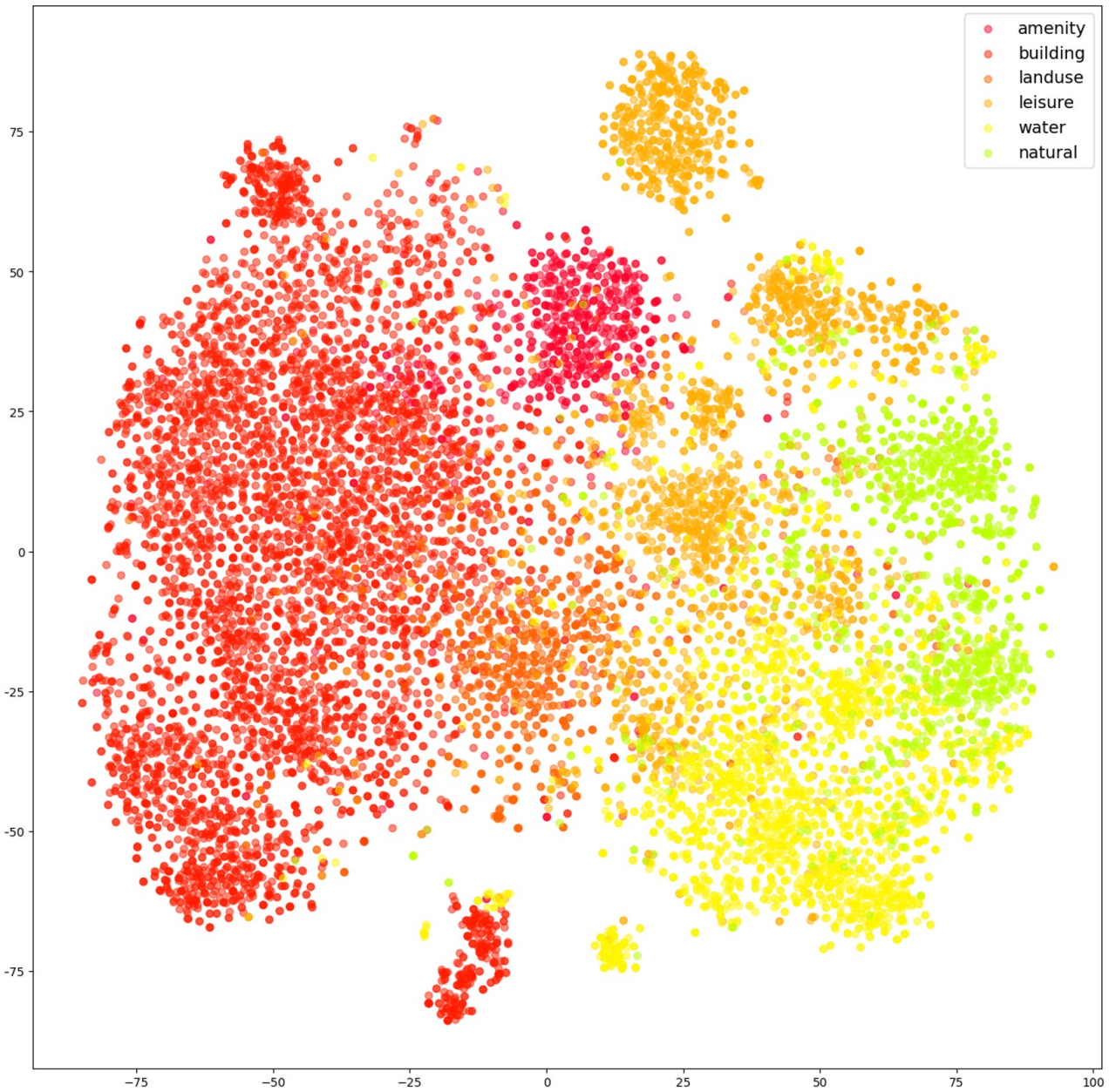
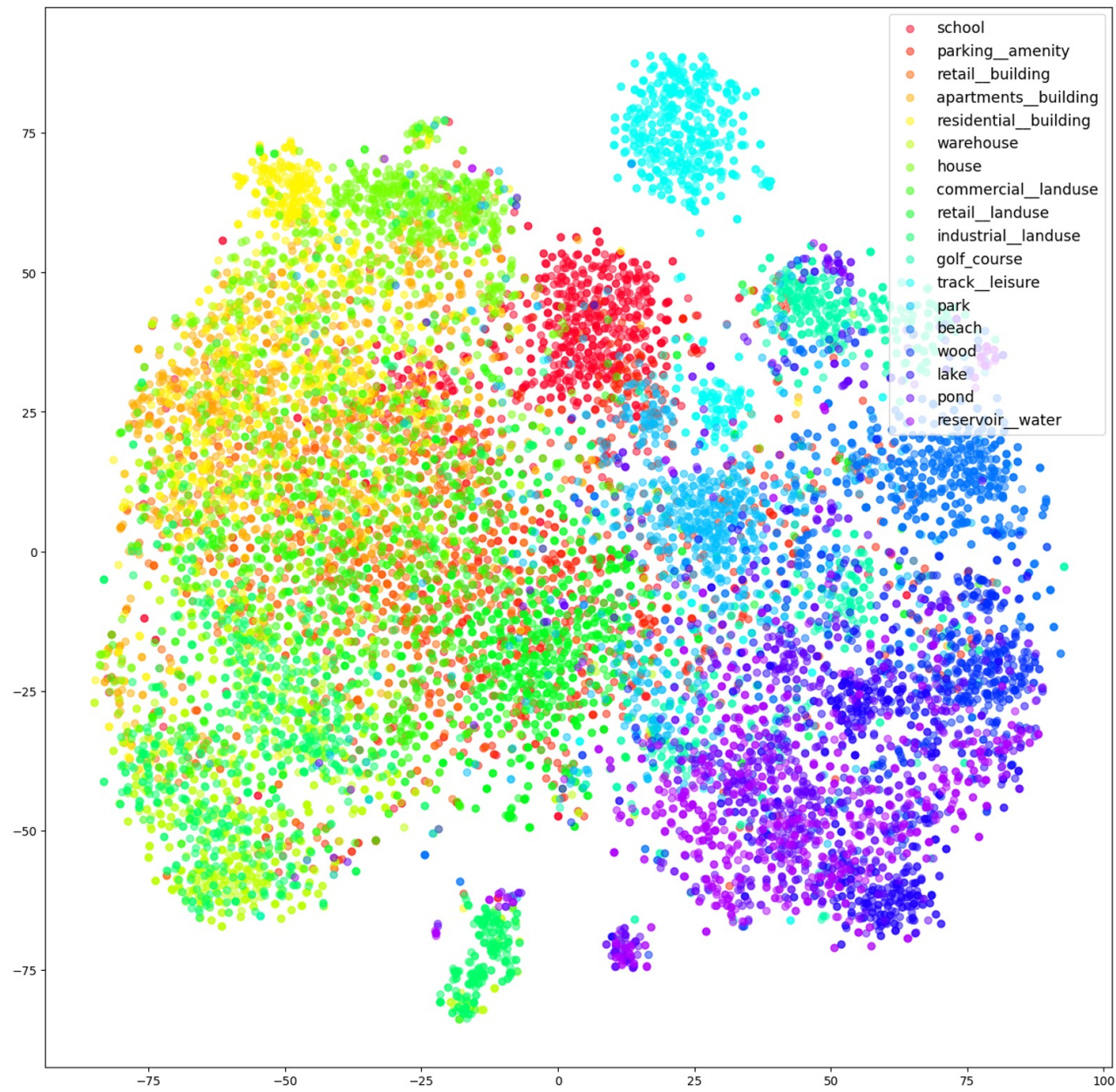
OpenStreetMap

Training: 200k CA OSM dump (2.3 tags avg)



Setting	WD-2k			OSM-16k		
	Precision	Recall	F_1	Precision	Recall	F_1
1 Ours _{shape}	0.497	0.506	0.501	0.473	0.512	0.492
2 Ours _{shape+spatial}	0.506	0.545	0.525	0.491	0.536	0.513
3 Ours _{full}	0.850	0.823	0.836	0.877	0.725	0.794
4 Ours _{full w/taxonomy}	0.849	0.852	0.850	0.858	0.854	0.856
GPT-3.5-Turbo	0.198	0.209	0.121	0.145	0.063	0.026
GeoVectors	0.819	0.834	0.826	0.833	0.815	0.824

SotA

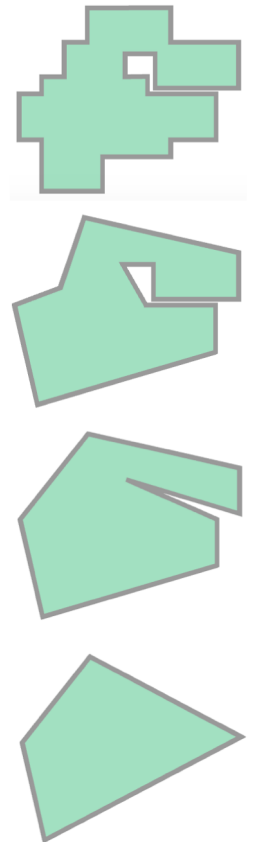


Embedding Geo-Entities for Semantic Typing


- Related Work
 - ML for Geospatial Classification (Castelluccio 2015, Klemmer 2023, Kaczmarek 2023, Xu 2022, Yan 2021)
 - Employ **CNNs**, **GNNs**, and **GCNs** for: building footprints & urban land-use classification
 - Do not address the incorporation of external (open) knowledge
 - Geospatial Embedding Techniques (Tempelmeier 2021, Jenkins 2019, Li 2022)
 - Develop **unsupervised embedding** such as *GeoVectors* & *SpaBert*
 - Do not address shape or explicit spatial data for enhanced geo-entity representation
 - OSM Embedding (Woźniak 2021)
 - Proposes embedding method for OSM regions using **hexagonal grids**
 - Does not address individual entities

Embedding Geo-Entities for Semantic Typing

- Takeaways
 - Method for **geo-referenced entity embedding** on the web
 - **self-supervised**
 - leverages **geometric, spatial, & semantic contexts**
 - **weighted** contrastive learning
 - enables seamless **semantic typing for integration** on the web
 - fuels further discovery & enrichment
 - Still, many challenges exist
 - **availability**
 - **granularity**
 - **quality**



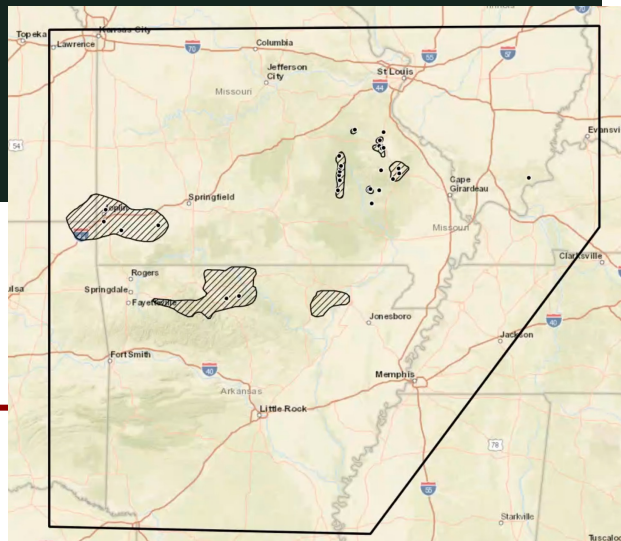
Agenda

- Basel's PhD Journey
- Intro
- Thesis Overview
- Approach:
 - Building Spatio-Temporal KGs from Digitized Maps
 - Embedding Geo-Entities for Semantic Typing
 - From Digitized Reports to Spatio-Temporal KGs 
- Conclusions & Future Directions

From Digitized Reports to Spatio-Temporal KGs

- Goals:
 - Integrating **geo-referenced textual & historical data** with quantitative information into a comprehensive, dynamic, & **spatio-temporal KG**
 - capture data & entity semantics, entity resolution, & accurate data modeling
 - Demonstrate via a KG of **historical mining data**

```
{  
  "MineralSite": [  
    {  
      "source_id": "https://w3id.org/usgs/z/4530692/2P29BJHV",  
      "record_id": 1,  
      "name": "NI 43-101 Technical Report for the Lantinen Koillismaa Project in Europe, Finland dated  
March 2017",  
      "location_info": {  
        "location": "POINT(28.17472 65.94611)",  
        "crs": "WGS84",  
        "country": "Finland",  
        "state_or_province": "Central Finland"  
      },  
    },  
  ],  
}
```



USGS
Mineral Resources / Online Spatial Data

Mineral Resources Data System (MRDS)

MRDS is a collection of reports describing metallic and nonmetallic mineral resources throughout the world. Included are deposit name, location, commodity, deposit description, geologic characteristics, production, reserves, resources, and references. It subsumes the original MRDS and MAS/MILS. MRDS is large, complex, and somewhat problematic. This service provides a subset of the database comprised of those data fields deemed most useful and which most frequently contain some information, but full reports of most records are available as well.

Current status: As of 2011, USGS has ceased systematic updates to MRDS, and is working to create a new database, focused primarily on the conterminous US. For locations outside the United States, MRDS remains the best collection of reports that USGS has available. For locations in Alaska, the Alaska Resource Data File remains the most coherent collection of such reports and is in continuing development.

Resource descriptions here include an indication of the overall quantity and diversity of information they contain. Many records in this database are simple reports of commodity at some location, but some records provide substantial detail of the geological setting and industrial exploitation of the resource. To help users find these more thorough records, a map interface and search form are provided that rank results by overall quality, records graded A having more information about more aspects of the resource, records graded D having only summary information about the resource. Records graded B and C are intermediate between these, and records graded E generally lack bibliographic references.

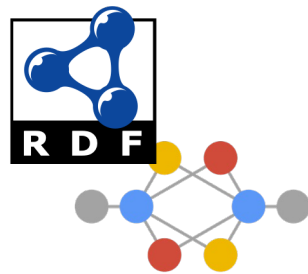
View

Show in a web browser window:

- MRDS records graded A
- Mineral Resources Data System: A collection of reports describing metallic and nonmetallic mineral resources throughout the world. Included are deposit name, location, commodity, deposit description, geologic characteristics, production, reserves, resources, and references. It subsumes the original MRDS and MAS/MILS. MRDS is large, complex, and somewhat problematic. This service provides a subset of the database comprised of those data fields deemed most useful and which most frequently contain some information, but full reports of most records are available as well.
- Mineral Resources (MRDS)
- A records (200)
- B records (24,200)
- C records (14,800)
- D records (21,195)
- E records (1,100)

Extracted
Textual &
Vector Data

**Data understanding &
KG construction**



From Digitized Reports to Spatio-Temporal KGs

- For a given commodity/deposit type/location/time-range:
 - Construct grade and tonnage models from the data on existing mines
 - Compile rich mineral site data

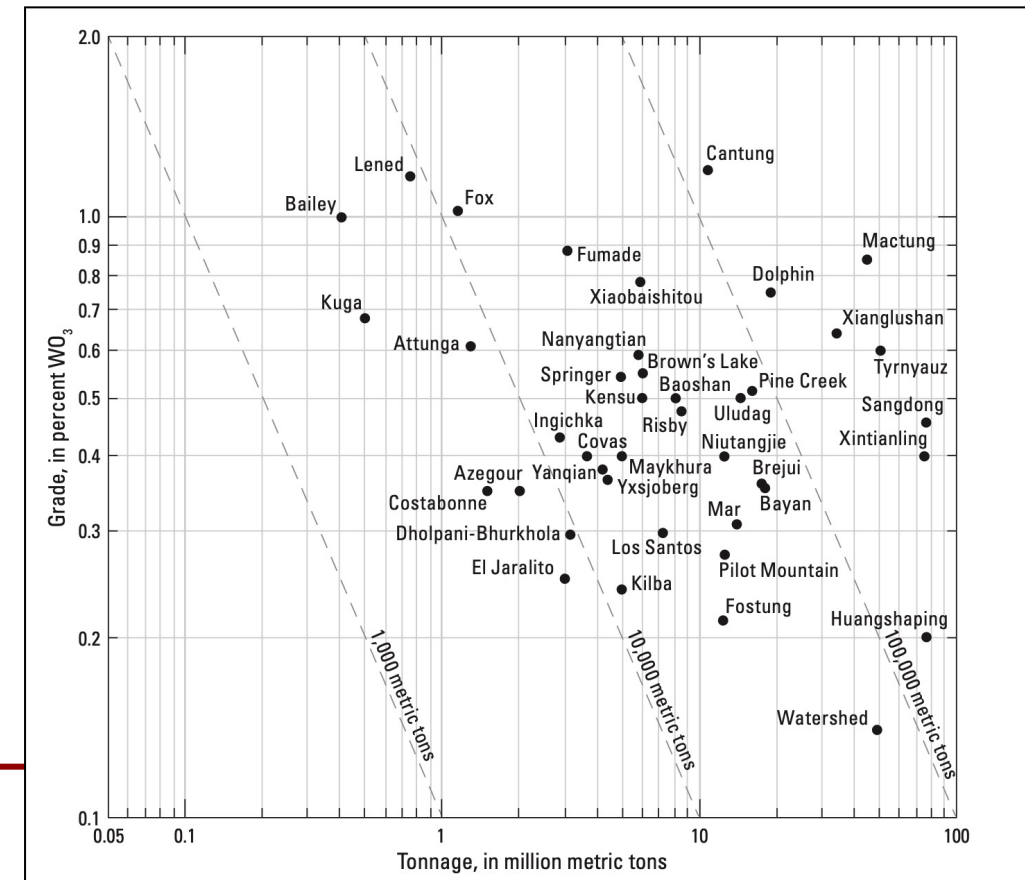
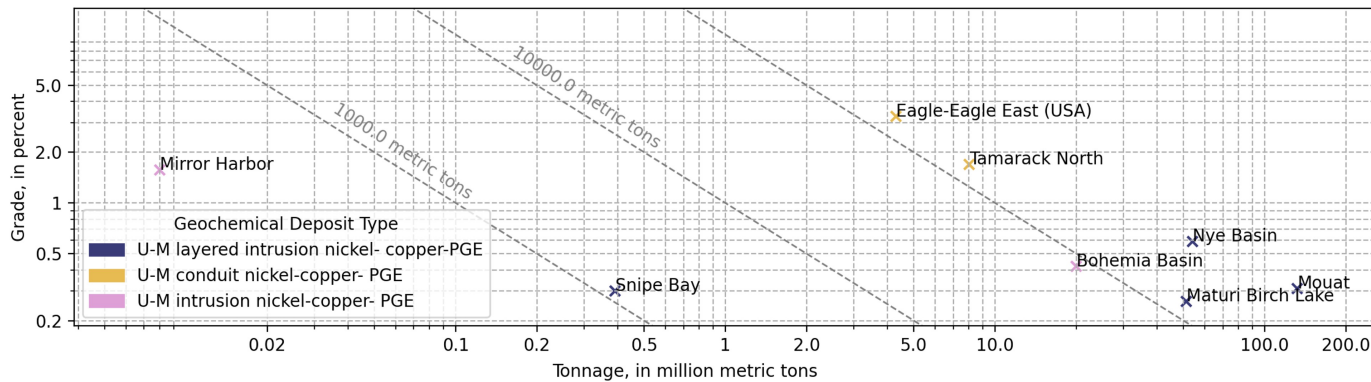


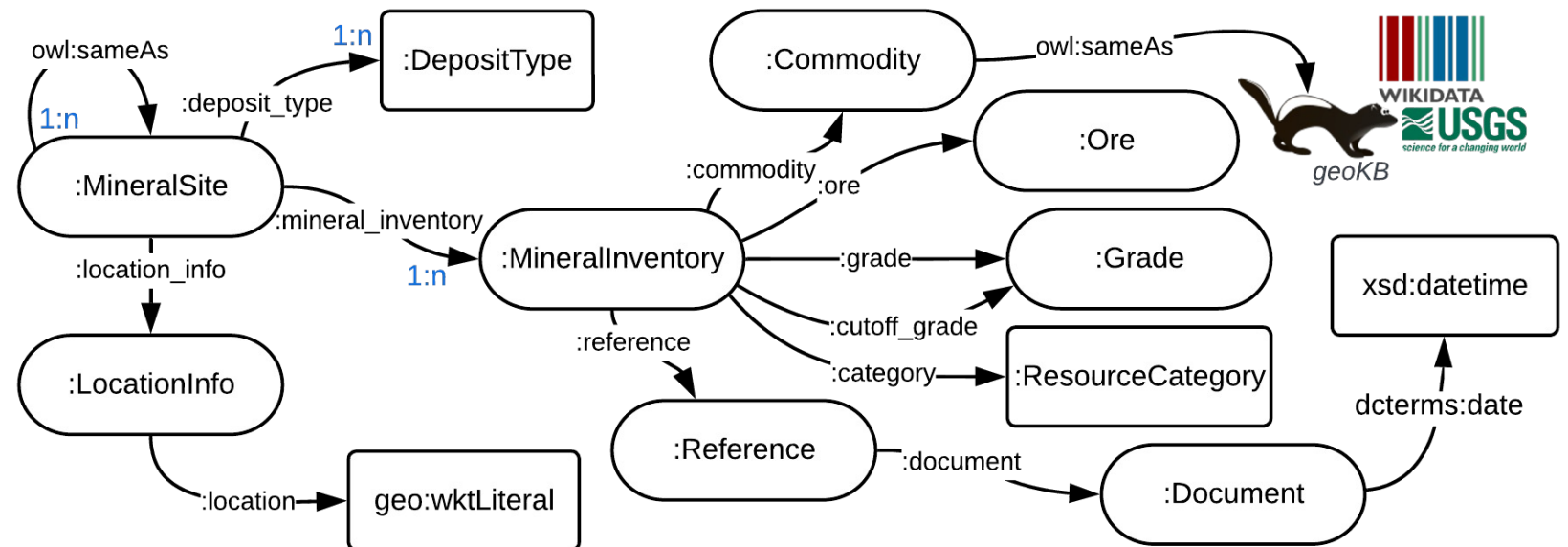
Figure from Green, C. J., Lederer, G. W., Parks, H. L., & Zientek, M. L. (2020). Grade and tonnage model for tungsten skarn deposits—2020 update (No. 2020-5085). US Geological Survey

From Digitized Reports to Spatio-Temporal KGs

- Approach

- Step 1. Semantic Modeling & URI assignment (Data Representation)

- Transform & materialize the data (construct KG)
 - **Generate entities** (URIs) based on unique identifiers
 - Provide a **useful semantic representation** supporting downstream tasks
 - Construct a **meaningful semantic model**
 - Follows W3C & OGC standards (**GeoSPARQL**)



From Digitized Reports to Spatio-Temporal KGs

- Approach

- Step 2. Entity Linking

- Link the generated entities to a domain data-rich vocab (i.e., GeoKB)

- Determine similarity by textual similarity (i.e., Jaccard)
- Directly within SPARQL

```

1 SELECT ?entity ?entityLabel WHERE {
2   ?entity rdfs:label ?entityLabel.
3   ?entity gkbt:P1 gkbi:Q406. # instance of mineral commodity
4   FILTER(CONTAINS(LCASE(?entityLabel), "nickel")) }

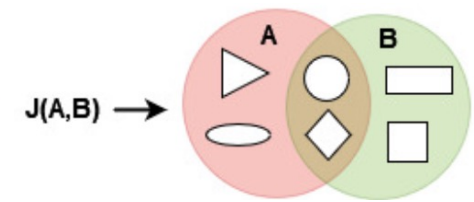
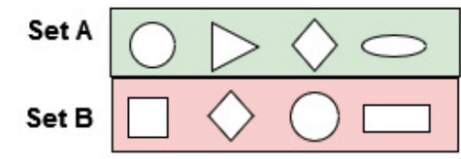
```

"Clay, Fire (Refractory)"

	uri	entity.value \
0	https://geokb.wikibase.cloud/entity/Q413	
1	https://geokb.wikibase.cloud/entity/Q424	
2	https://geokb.wikibase.cloud/entity/Q423	
3	https://geokb.wikibase.cloud/entity/Q421	
4	https://geokb.wikibase.cloud/entity/Q162319	
5	https://geokb.wikibase.cloud/entity/Q425	
6	https://geokb.wikibase.cloud/entity/Q426	
7	https://geokb.wikibase.cloud/entity/Q428	
8	https://geokb.wikibase.cloud/entity/Q427	
9	https://geokb.wikibase.cloud/entity/Q429	

	entityLabel.xml:lang	entityLabel.type	entityLabel.value
0	en	literal	high alumina clay aluminum
1	en	literal	bloating material clay
2	en	literal	brick clay
3	en	literal	clay
4	en	literal	Clay
5	en	literal	bentonite clay
6	en	literal	chlorite clay
7	en	literal	fire (refractory) clay
8	en	literal	fullers earth clay
9	en	literal	glaucanite clay

Jaccard



From Digitized Reports to Spatio-Temporal KGs

- Approach
 - Step 3. Triplify!

```

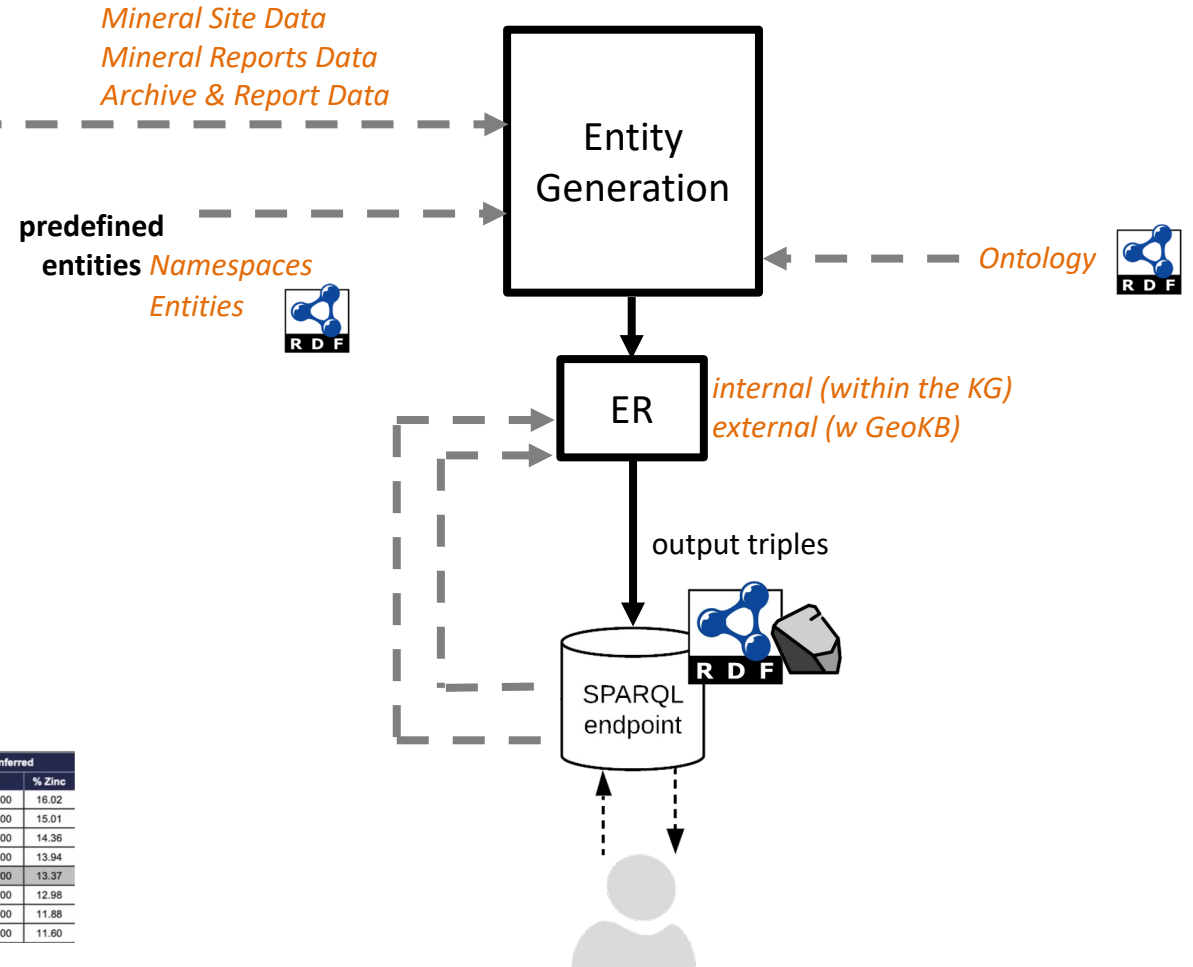
id: "Site103"
name: "Balmat - Edwards District"
location_info:
  location: "POINT (-75.33292 44.28331)"
  country: "United States"
  location_source_record_id: "22"
  location_source: "MRDS_Zinc"
  crs: "NAD83"
same_as:
  MRDS_Zinc:
    id: 22
    Attributes:
      country: "United States"
      ogc_fid: 71087
      dep_id: 10073156
      site_name: "Balmat - Edwards District"
      dev_stat: "Producer"
      url: "https://mrdata.usgs.gov/mrds/show-mrds.php?dep_id=10073156"
      code_list: "ZN PB AG HG"
      wkb_geometry: "0104000020E6100000010000000101000000CC0BB08"
      commodity: "Zinc, Lead, Silver, Mercury"
      geometry: "POINT (-75.33292 44.28331)"
            
```

```

id: 0
name: ""
location_info: "44o14'51" N, 75o23'50" W"
geology_info: ""
same_as: ""
MineralInventory:
  id: 1
  commodity: "zinc"
  category: "Indicated"
  ore:
    grade:
      cutoff_grade:
        contained_metal: 174604.65
        reference:
          date: "09-19-2017"
            
```

Cut-Off (% Zinc)	Measured		Indicated		M&I		Inferred	
	tons	% Zinc	tons	% Zinc	tons	% Zinc	tons	% Zinc
>10%	543,000	16.15	840,600	16.27	1,383,600	16.22	1,499,200	16.02
>9%	617,500	15.34	962,500	15.42	1,580,000	15.39	1,772,600	15.01
>8%	696,100	14.57	1,080,000	14.67	1,776,100	14.63	1,970,400	14.36
>7%	770,200	13.89	1,200,500	13.96	1,970,700	13.93	2,100,600	13.94
>6%	850,100	13.19	1,307,900	13.35	2,158,000	13.29	2,276,000	13.37
>5%	932,800	12.51	1,416,700	12.76	2,349,500	12.66	2,393,400	12.98
>4%	1,004,900	11.94	1,524,400	12.18	2,529,300	12.08	2,887,100	11.88
>3%	1,074,300	11.39	1,612,400	11.70	2,686,700	11.58	2,824,300	11.60

Cut-Off (% Zinc)	Measured	Tons	% Zinc
>10%	Measured	543,000	16.15
	Indicated	840,600	16.27
	M&I	1,383,600	16.22
	Inferred	1,499,200	16.02
>9%	Measured	617,500	15.34
	Indicated	962,500	15.42
	M&I	1,580,000	15.39
	Inferred	1,772,600	15.01



From Digitized Reports to Spatio-Temporal KGs

- Evaluation

- Data completeness (SHACL)
- EL
- RDF Query Performance

Data: 2.4m triples // 135 commodities // focus on 2 critical mineral: nickel, zinc

Characteristic	Count
Total Triples	2,397,708
Distinct Classes	16
Instances (Non-literals)	226,267
Geospatial Instances	2,884
Blank Nodes	1,518,981

Method	MRR	Hits@1	Hits@3	Hits@5
String search, then Jaro	0.557	0.459	0.659	0.659
String search, then Jaccard	0.648	0.637	0.659	0.659
Instance search, then Jaro	0.801	0.689	0.926	0.956
Instance search, then Jaccard (proposed)	0.940	0.904	0.978	0.978

```

1 ?ms :location_info/:location ?loc_wkt .
2 FILTER(geof:distance(?loc_wkt, "POINT(-118.57 47.56)"^^geo:wktLiteral, unit:mile) < 500)
    
```

Query Constraint Type	Avg	Min	Max
Textual	450	369	649
Temporal/Numeric	438	388	607
Spatial	708	501	811

From Digitized Reports to Spatio-Temporal KGs

- Related Work
 - General geo KBs (Zhu 2017, Brodaric 2020)
 - Mostly encompasses **conceptual knowledge & data**
 - Does not address: quantitative data integration
 - GeoKGs related to mineral data (Qun 2023)
 - Tailored for **geochemical data**
 - Does not address: quantitative data integration
 - Information extraction for geo KGs (Wang 2018)
 - Focus is on the **data extraction**
 - Does not address: data integration & entity linking

From Digitized Reports to Spatio-Temporal KGs

- Takeaways

- Paradigm for **geospatial data integration** from digitized textual & geo-referenced archive data
- Method to **identify & retrieve** instances of a given type from a **publicly available KG**
 - **Automated**
 - **Incremental** semantic model for simple & efficient **querying**
 - Follows LD & SW principles
 - **Linked to Web**
 - Fuels further discovery & enrichment
- Still, many challenges exist
 - **Quality** of extracted data
 - Encoding **domain knowledge**



Putting it together

Grade	Assays				
	Cu (%)	Ni (%)	S (%)	Au (g/t)	Pt (g/t)
Concentrate	7.16-10.1	1.66-2.20	18.4-21.5	0.65-1.28	1.17-1.59

```

qudt:hasDimension: "",
qudt:abbreviation": "g tonne-1",
..
qudt:hasPart: [
{
  qudt:hasDimension: "M",
  qudt:quantityKind: "http://data.nasa.gov/qudt/owl/unit#Gram",
  qudt:conversionMultiplier: 0.001,
  qudt:conversionOffset: 0.0,
  qudt:symbol: "g"
},
{
  ccut:exponent: "-1",
  qudt:hasDimension: "M",
  qudt:quantityKind: "http://data.nasa.gov/qudt/owl/unit#MetricTon",
  qudt:conversionMultiplier: 1000.0,
  qudt:conversionOffset: 0.0,
  qudt:symbol: "t"
}
]

```

2nd aggregation:
total grade & tonnage
computation

```

3      ?total_contained_indicated ?total_contained_inferred
4      (?total_tonnage_measured + ?total_tonnage_indicated
5      (?total_contained_measured + ?total_contained_indicated
6      (IF(?total_tonnage > 0, ?total_contained_metal / ?total_tonnage
7      WHERE {
8      {
9      SELECT ?ms ?ms_name ?deposit_name ?country ?loc_wkt
10     (SUM(?tonnage_measured) AS ?total_tonnage_measured
11     (SUM(?tonnage_indicated) AS ?total_tonnage_indicated
12     (SUM(?tonnage_inferred) AS ?total_tonnage_inferred
13     (SUM(?contained_measured) AS ?total_contained_measured
14     (SUM(?contained_indicated) AS ?total_contained_indicated
15     (SUM(?contained_inferred) AS ?total_contained_inferred
16     WHERE {
17     ?ms :deposit_type [ rdfs:label ?deposit_name ] .
18     ?ms :mineral_inventory ?mi .
19     OPTIONAL { ?ms rdfs:label:name ?ms_name . }
20     ?ms :location_info/:location ?loc_wkt .
21     ?mi :category ?mi_cat .

```

```

FILTER(geof:sfWithin(?loc_wkt, "POLYGON(...)" ) )

```

```

FILTER(geof:distance(?loc_wkt, POINT(-118.57 47.56)^^geo:wktLiteral, unit:mile) < 500)

```

```

FILTER(?date >= "2000"^^xsd:gYear && ?date <= "2010"^^xsd:gYear) .

```

```

?mi :ore [ :ore_value ?ore_val_raw; :ore_unit ?ore_unit ] .
?mi :grade [ :grade_value ?grade_val; :grade_unit ?grade_unit ] .
BIND(IF(bound(?ore_val_raw), ?ore_val_raw, 0) AS ?ore_val_pre)

```

```

BIND(IF(?ore_unit = <http://data.nasa.gov/qudt/owl/unit#MetricTon>, ?ore_val_pre / 1e6, ?ore_val_pre)) AS ?ore_val

```

```

29 BIND(IF(CONTAINS(LCASE(STR(?mi_cat)), "measured"), ?ore_val, 0) AS ?tonnage_measured)
30 BIND(IF(CONTAINS(LCASE(STR(?mi_cat)), "indicated"), ?ore_val, 0) AS ?tonnage_indicated)
31 BIND(IF(CONTAINS(LCASE(STR(?mi_cat)), "inferred"), ?ore_val, 0) AS ?tonnage_inferred)
32 BIND(IF(CONTAINS(LCASE(STR(?mi_cat)), "measured") && ?grade_val > 0, ?ore_val * ?grade_val, 0) AS ?contained_measured)
33 BIND(IF(CONTAINS(LCASE(STR(?mi_cat)), "indicated") && ?grade_val > 0, ?ore_val * ?grade_val, 0) AS ?contained_indicated)
34 BIND(IF(CONTAINS(LCASE(STR(?mi_cat)), "inferred") && ?grade_val > 0, ?ore_val * ?grade_val, 0) AS ?contained_inferred)
35 }

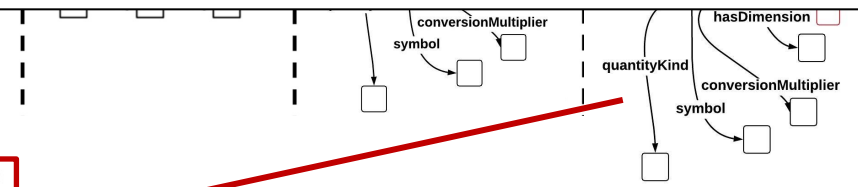
```

1st aggregation:
tonnage computation


```

36 GROUP BY ?ms ?ms_name ?deposit_name ?loc_wkt }

```



Agenda

- Basel's PhD Journey
- Intro
- Thesis Overview
- Approach:
 - Building Spatio-Temporal KGs from Digitized Maps
 - Embedding Geo-Entities for Semantic Typing
 - From Digitized Reports to Spatio-Temporal KGs
- Conclusions & Future Directions 

Conclusions

- Presented my thesis, based on 3 main contributions
 - Paradigm for **automatic transformation** of **historical maps** into dynamic **spatio-temporal KGs**
 - evaluated via **change analysis** of 2 different topographic features over time
 - Approach for accurate **geo-entity embedding, classification & integration**
 - evaluated via **semantic typing** of geo-referenced digitized instances to **two different Open KBs**
 - Method for the construction of a **spatio-temporal KG** from geo-referenced **spatial entities in archive reports**
 - demonstrated & evaluated on different tasks in the domain of **historical mining data**
- Presented 2 additional (auxiliary) contributions
 - Method & tool to **generate geo-feature taxonomies** from public OSM data
 - Approach & tool to **identify, ontologize & transform units of measurement** for quantitative data

Future Directions

- Advanced **data modeling**
 - More modalities
 - More data (e.g., rapidly changing geographies)
 - Hyperparameter optimization
- **Enhanced embedding** techniques
 - Utilize subword information and deep learning attention mechanisms
 - Expand integration of textual data
- **KG** expansion
 - Extend KG linkage to cover a broader range of knowledge bases & LOD
 - Apply & integrate with additional domains like archaeology & environmental sciences
- Dynamic **semantic modeling**
 - Create more sophisticated & evolving semantic models for accurate representation across multiple domains

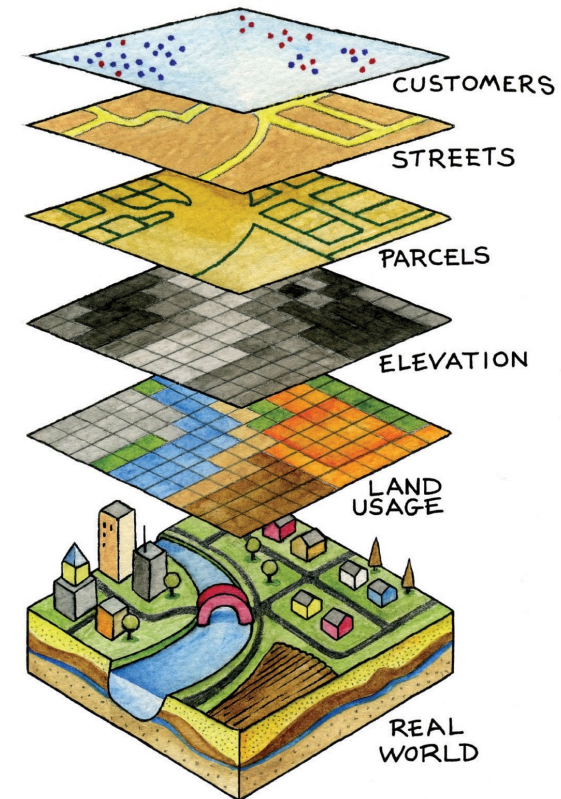


figure from *Essentials of Geographic Information Systems, Ch 7, Saylor Academy, 2012*

Final Remarks

- Thanks to my collaborators!
 - Craig A. Knoblock (**advisor**)
 - Cyrus Shahabi, John P. Wilson, Jay Pujara, Yao-Yi Chiang (**committee**)
 - Pedro Szekely, Filip Ilievski, Muhammad Rostami, Jon May (**USC/ISI collaborators**)
 - Johannes H. Uhl, Stefan Leyk (**University of Colorado Boulder**)
 - Anna Lisa Gentile, Pengyuan Li, Guang-Jie Ren (**IBM Research**)
 - Abha Moitra (**GE Research**)
 - Shui Hu (**Amazon Research**)
 - **USC/ISI** colleagues & friends

Thank you for listening!
Questions?