# KnowWhereGraph:
## Enriching and Linking Cross-Domain Knowledge Graphs using Spatially-Explicit AI Technologies to Address Pressing Challenges at the Human-Environment Nexus
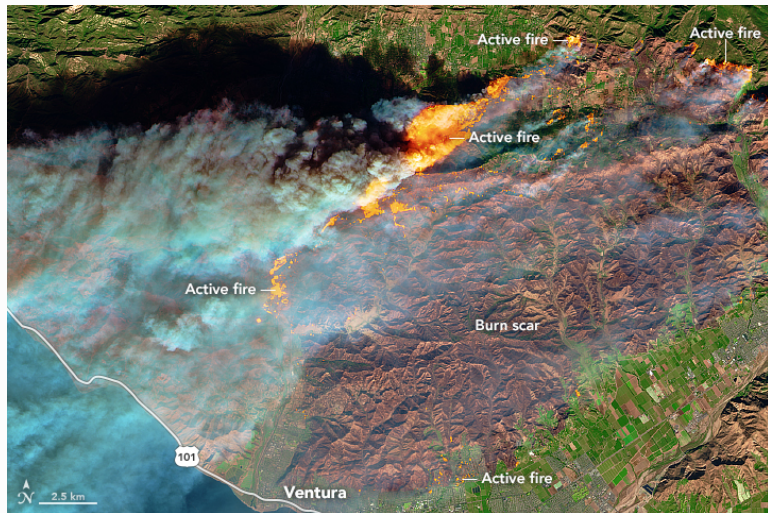
Krzysztof Janowicz (for A-6677)

STKO Lab, University of California, Santa Barbara, USA

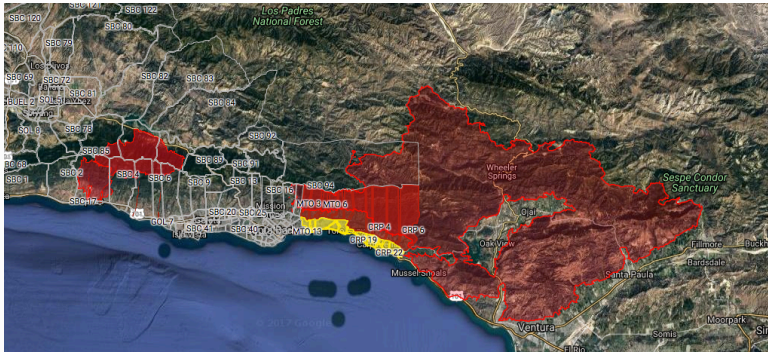April 2020
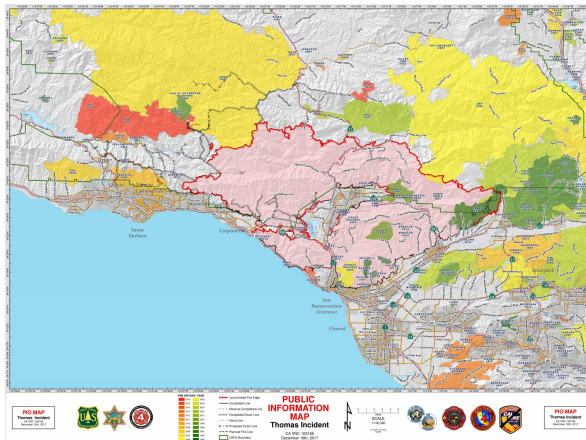
# Thomas Fire and Mudslide

# Thomas Fire and Mudslide



■ Screenshot of the Santa Barbara County post-fire hazards map shows the flooding risk areas in blue (shown on 1/5/18; days before the storm/mudslide).
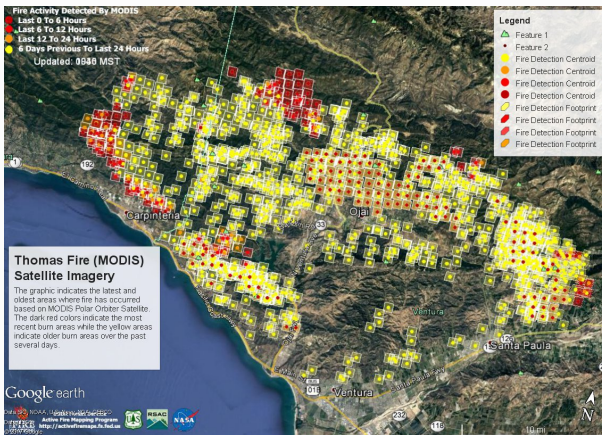
# Thomas Fire and Mudslide



■ 'But your own home, right next to a creek, was in a voluntary evacuation zone. Didn't that seem odd to you? No. To even be in a voluntary zone seemed significant to me. The problem with evacuating for a flood is that our evacuation zones aren't broken up by watersheds. They're broken up by roads. If you look at our zone map, it's broken up in boxes, and you evacuate people from that box. Part of that is so you can protect it from looting. Once you clear everybody out, you kind of own it, and now you have to keep it safe.'

# Thomas Fire and Mudslide as a Vector (Map)



- We know how to **partially represent** such 'object' data on a KG, e.g., using GeoSPARQL, SWEET, SSN.

# Thomas Fire as a Raster Product



- **Not so clear how to put such 'field' data on the graph.** E.g., annotation ontologies, SSN, object recognition/extraction, semantic image parsing, and so on.

# Why SpaceX is so Radically Different?



t⏎ Mary Retweeted

**SPadre**
@SpacePadreIsle

SpaceX fleet vessel GO Discovery journey to Port of Brownsville has paused at port of Fort Lauderdale. Very high wind and seas are expected the next several days in the Gulf of Mexico, placing cargo on the low profile deck at risk from overwash. #SpaceX #BocaChica #Starship

| | |
|---|---|
| Latitude: 26.09246 | |
| Longitude: -80.1115 | |
| Area: | US East Coast |
| Info received: | 9 Jan 04.51 UTC |
| Status: | Stopped |
| Speed/Course: | 0 kn / 229° |
| AIS Source: | 700 (TowBoatU.S. Ft. Lauderdale) |
| Wind: | 24 kn |
| Wind direction: | 64° ENE |
| Temperature: | 22°C |
| Distance from my location: | 821 NM 90° |

9:04 PM · Jan 8, 2020 from South Padre Island, TX · Twitter for iPhone
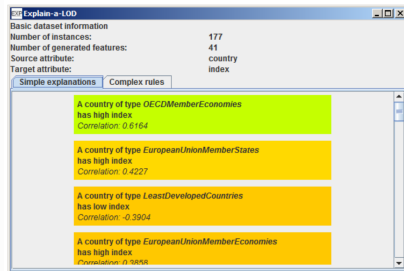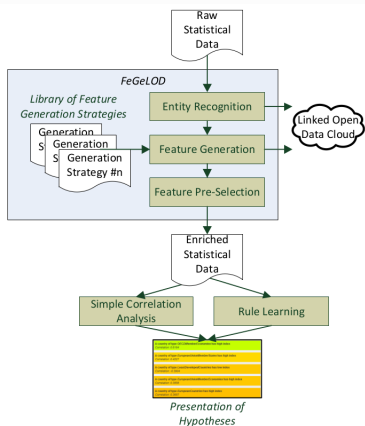
**3** Retweets **67** Likes

Imagine yourself as an analyst at

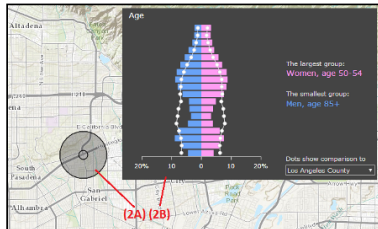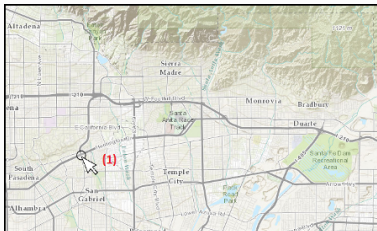- their **insurance** company
- their **competitors**
- …



Can this be **generalized**, e.g., supply chains more broadly?

# Explain-a-LOD!



'Linked Open Data as a means to interpret statistics [and collect] attributes to generate hypotheses for *explaining*[...]' (Paulheim, 2012)

# GeoEnrichment



'Enriches your data by adding demographic and landscape facts about the people and places that surround or are inside your data locations. The output is a duplicate of **your input with new attribute fields** added to the table.' (Esri)
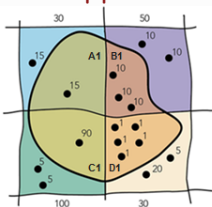
# Esri's GeoEnrichment

## Fetch data



## Data apportionment



**Weighted Block %**

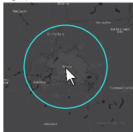| | Pop | |
|---|---|---|
| A1 | 15 | (15/30 = 50% × 30) |
| B1 | 30 | (30/50 = 60% × 50) |
| C1 | 90 | (90/100 = 90% × 100) |
| D1 | 5 | (5/30 = 16% × 30) |

(figures from Esri)

Imagine you could **supercharge GeoEnrichment**

- ■ 'Things not strings' (**semantic** graph data)

- ■ Open-ended, up-to-date **open content** across domains instead of pre-selected themes.

- ■ **Traverse** the graph instead of flat tables

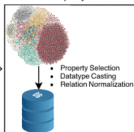▶ **Reenvision GeoEnrichment by fusing it with Linked Data!**

# Regions Studied by Explorers Influenced by von Humboldt
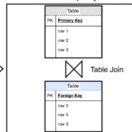
# Deeply Integrating Linked Data with GIS



(Mai, G., Janowicz, K., Yan, B., & Scheider, S. (2019).
Deeply integrating linked data with geographic
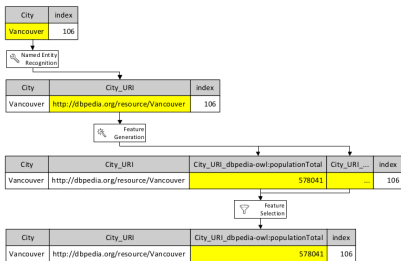information systems. Transactions in GIS, 23(3))

**Many well-known challenges:**

- Ontologies (ODP) + microtheories vs. learning representations
- Ontology alignment
- Co-reference resolution/ deduplication
- Data fusion / conflation
- KG summarization
- Link prediction
- **Spatial and temporal scopes**
    - What is an 'essential service' varies by state (e.g., COVID-19)
    - How long do we have to socially distance

Do not worry, we are not trying to solve the **semantic interoperability** problem ;-)
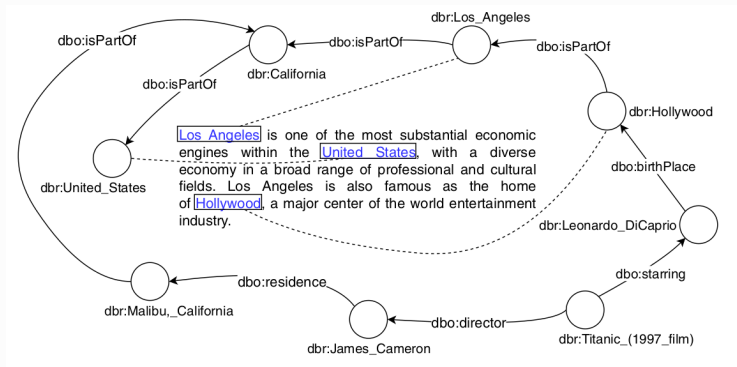
# Named Entity Recognition, Disambiguation, and Linking

| City | index |
|------|-------|
| Vancouver | 106 |

↓ Named Entity Recognition

| City | City_URI | index |
|------|----------|-------|
| Vancouver | http://dbpedia.org/resource/Vancouver | 106 |

↓ Feature Generation

| City | City_URI | City_URI_dbpedia-owl:populationTotal | City_URI_… | index |
|------|----------|--------------------------------------|------------|-------|
| Vancouver | http://dbpedia.org/resource/Vancouver | 578041 | … | 106 |

↓ Feature Selection

| City | City_URI | City_URI_dbpedia-owl:populationTotal | index |
|------|----------|--------------------------------------|-------|
| Vancouver | http://dbpedia.org/resource/Vancouver | 578041 | 106 |

(Paulheim; ESWC 2012)

| Geographic Area | Confirmed Cases as of 4/28/2020 | Number of Deaths |
|-----------------|--------------------------------|------------------|
| SOUTH COUNTY UNINCORPORATED AREA includes communities of Montecito, Summerland and the City of Carpinteria | 21 | 0 |
| CITY OF SANTA BARBARA and the unincorporated area of Mission Canyon | 56 | 0 |
| CITY OF GOLETA | 7 | 2 |
| COMMUNITY OF ISLA VISTA | 1 | 0 |
| UNINCORPORATED AREA OF THE GOLETA VALLEY AND GAVIOTA | 13 | 1 |
| SANTA YNEZ VALLEY including the Cities of Solvang & Buellton, and the communities of Santa Ynez, Los Alamos, Los Olivos and Ballard | 5 | 0 |
| CITY OF LOMPOC and the communities of Mission Hills and Vandenberg Village | 79 | 1 |
| People incarcerated at the Federal Prison in Lompoc | 104 | 1 |
| CITY OF SANTA MARIA | 128 | 1 |
| COMMUNITY OF ORCUTT | 36 | 0 |
| UNINCORPORATED AREAS of Sisquoc, Casmalia, Garey, Cuyama, New Cuyama, and the City of Guadalupe | 24 | 1 |
| Pending | 3 | |
| Total | 477 | 7 |

(publichealthsbc.org)

Key problem for Explain-a-LOD like tools, e.g., what if there is no 1:1 correspondence for NER? How would data apportionment work for graphed data?
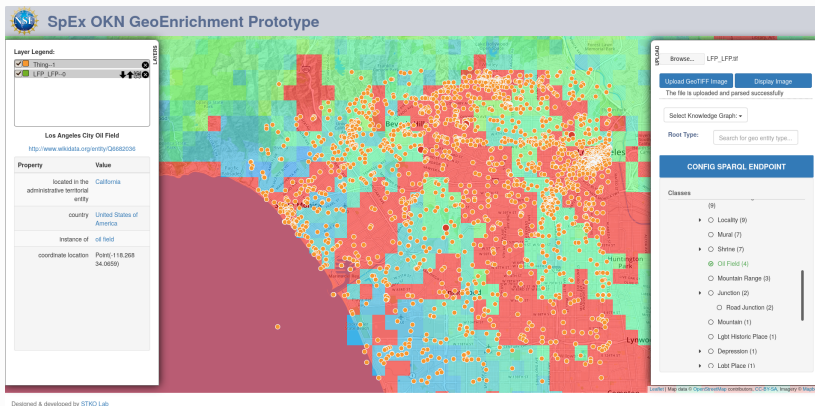
# KG Summarization



**Many steps required:** Data apportionment, relevant features –> **KG summarization**; n-degree path queries –> requires (multi-hop) **link prediction** to handle sparsity

(Yan, B., Janowicz, K., Mai, G., & Zhu, R. (2019). A spatially explicit reinforcement learning model for geographic knowledge graph summarization. Transactions in GIS, 23(3), 620-640.)

# Phase I Experiments



Designed & developed by STKO Lab

Can we **enrich** remotely sensed information or model outputs with auxiliary data (not visible from the sky) such as events that took place in a region?

# Phase I Experiments



Can we improve the **geographic** components of successful commercial systems such as **OpenCalais**, e.g., by introducing the concept of a **region**?

# Phase I Experiments


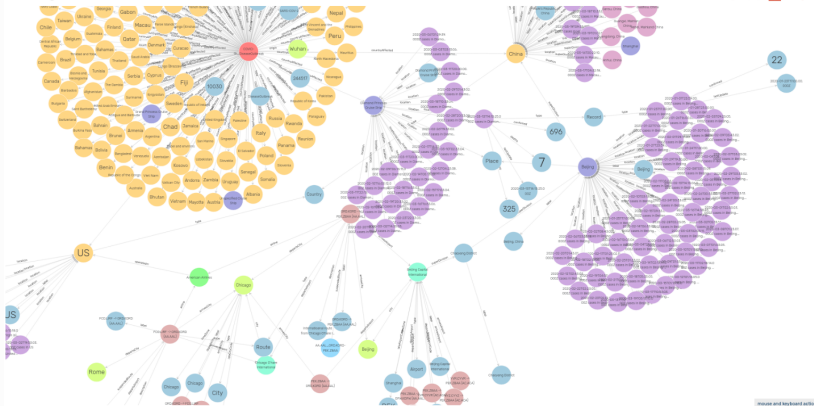
Knowledge graphs are about direct access to AI-ready, data-level statements (not dataset search). How can we guide industry in understanding whether their data is **graph-ready**?

# Our COVID-19 Knowledge Graph



**Geographically integrating data across domains** with a focus on social distancing measures, transportation, supply chain disruptions, canceled local events, geographic regions...

# KnowWhereGraph Vision

- (Geographic) space and time matter not only for the obvious reason that everything happens somewhere and at some time, but because knowing **where** and **when** things happen is critical to understanding **why** and how they happened or will happen.

- With 'KnowWhere' Graph we hope to take **GeoEnrichment** to the next level, by providing **open** graph-based linking and semantic enrichment technologies far beyond pre-defined data themes and silos.

- The ultimate goal of our project is to understand how to **engineer** meaningful **features** (independent variables) via knowledge graph-based GeoEnrichment for **downstream models** such as supply chain forecasting or soil health mapping

- This remains a **key roadblock** for our partners in industry, governments, and non-profits.