

Where on Earth Do Users Say They Are?: Geo-Entity Linking for Noisy Multilingual User Input

 Tessa Masis (tmasis.github.io/) and Brendan O'Connor

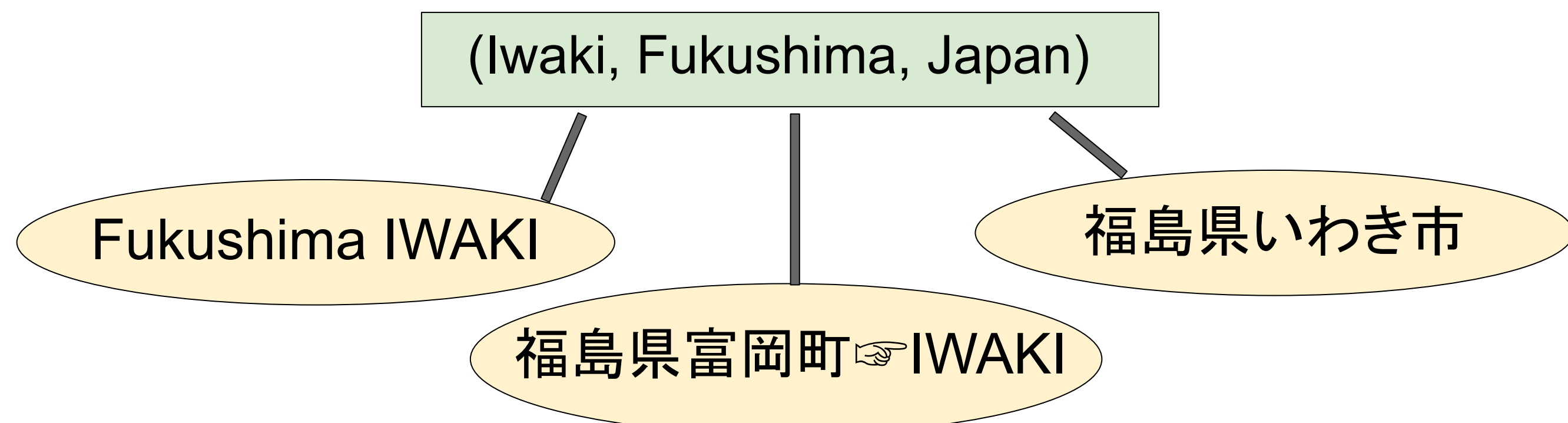
Geo-Entity Linking

brendan o'connor
@brendan642

Faculty @UMassCS | Natural language processing and computational social science | brenocon.com | @brenocon.bsky.social

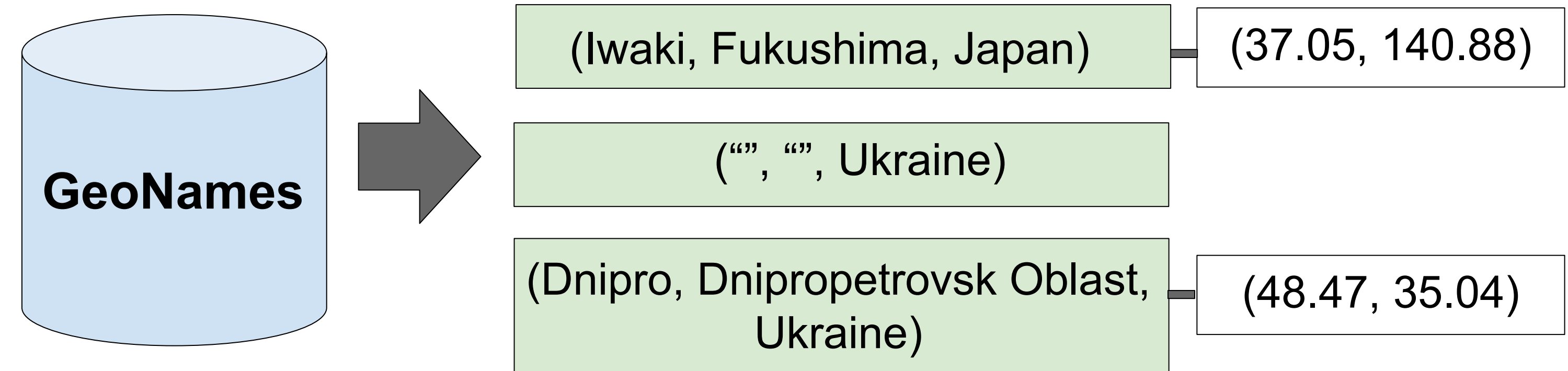
📍 western mass (northampton) 🌐 brenocon.com 📅 Joined June 2008

Locations can be referred to on social media with varying writing styles, informal names, and in different languages or writing scripts



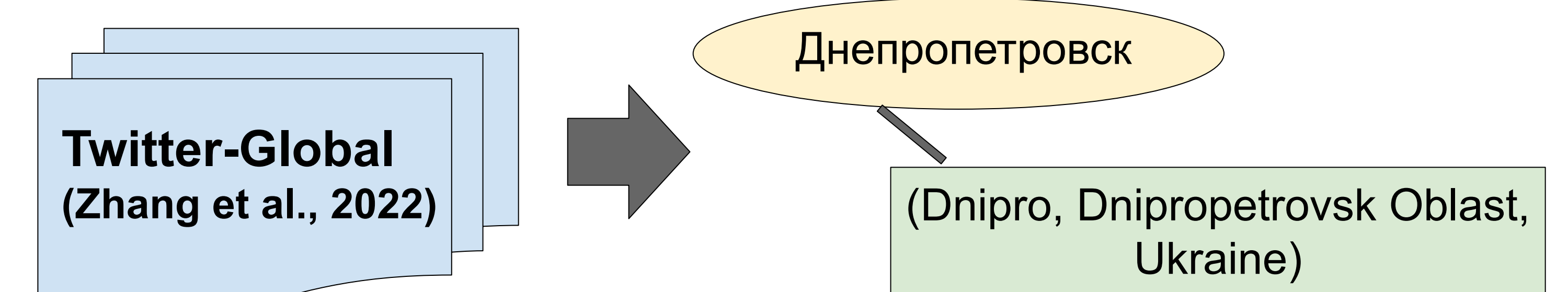
Goal: given a noisy free text Location field mention, link it to the correct real-world geographic entity

Data



Target location database

28,767 distinct locations; cities labeled with coordinates

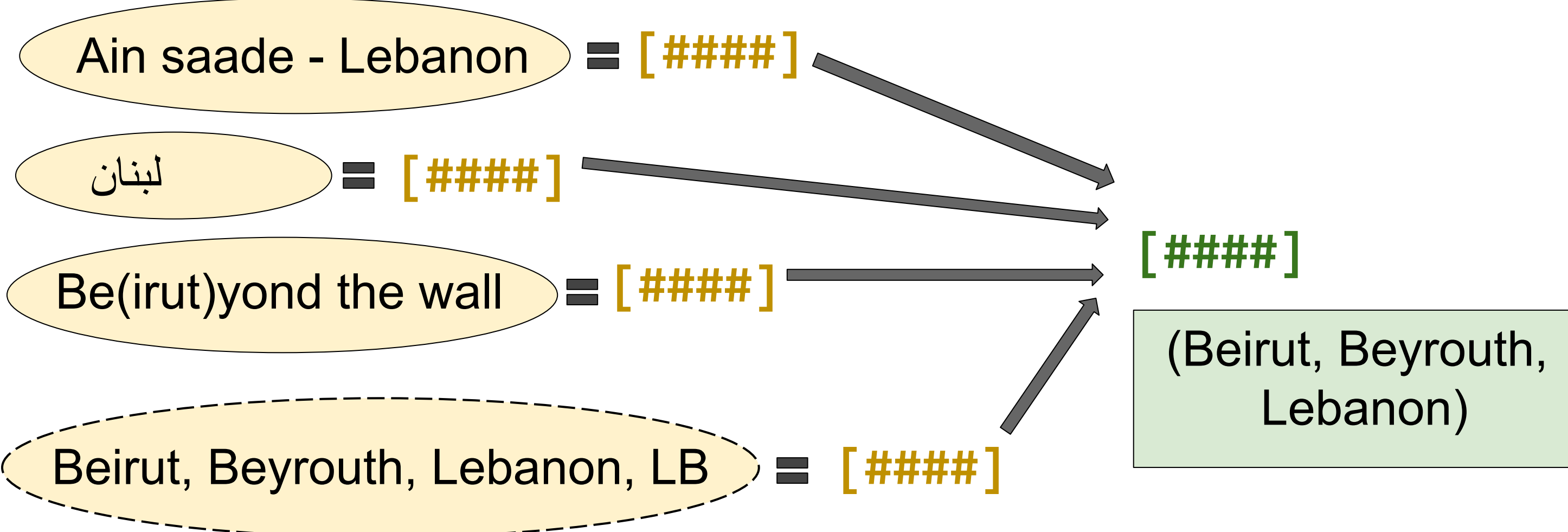


Labeled geo-entity linking dataset

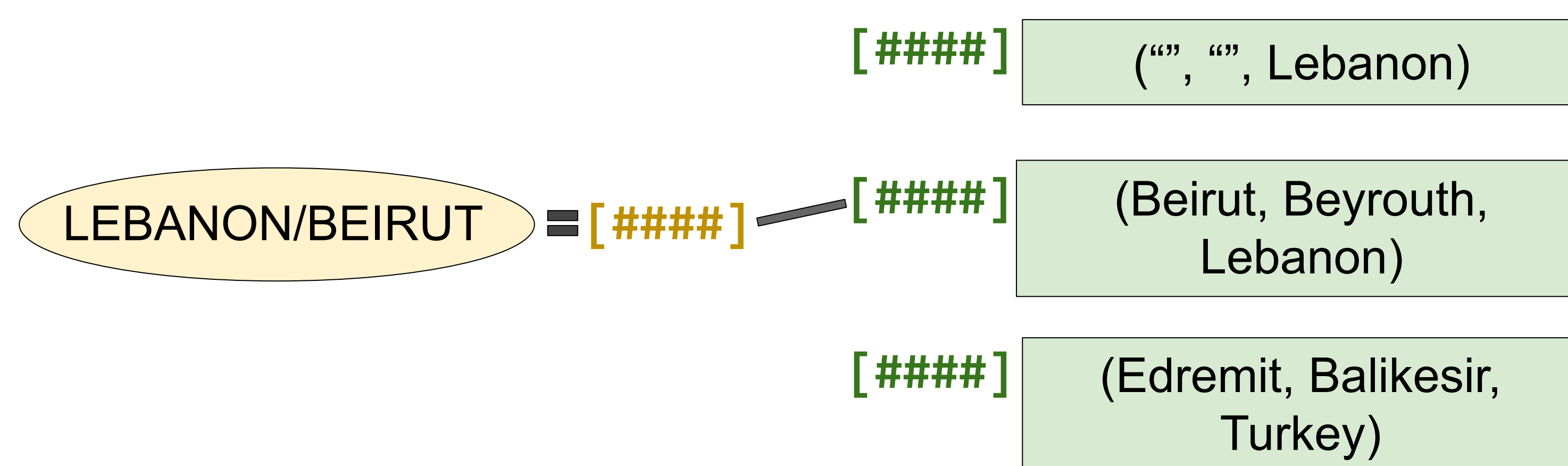
4.1M geocoordinate-tagged tweets; we link each poster's Location field to a ground truth location, defined as the closest city in GeoNames database

Proposed Method: UserGeo

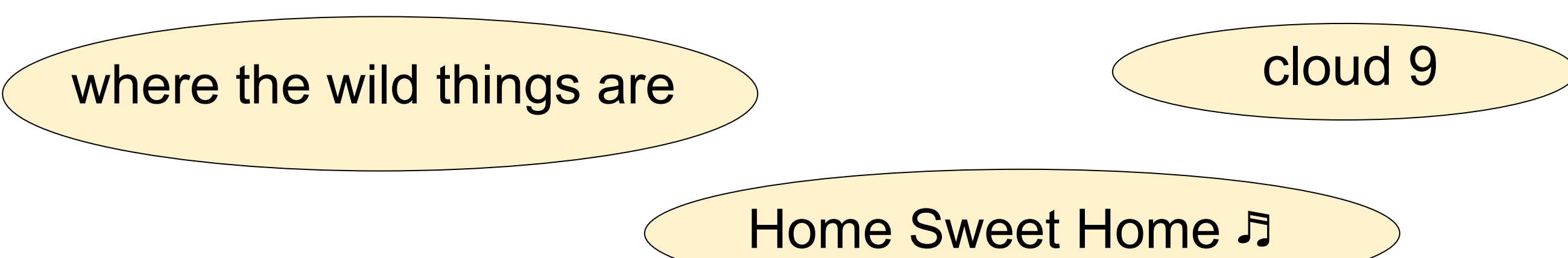
1) **Training:** For each location in GeoNames database, create a soft-alias location name representation by averaging SBERT embeddings of all linked Location fields in Twitter-Global



2) **Predicting:** For a new free text location mention, predict the location with the highest cosine similarity



3) If cosine similarity is below a certain threshold, make no guess i.e. NULL



Results

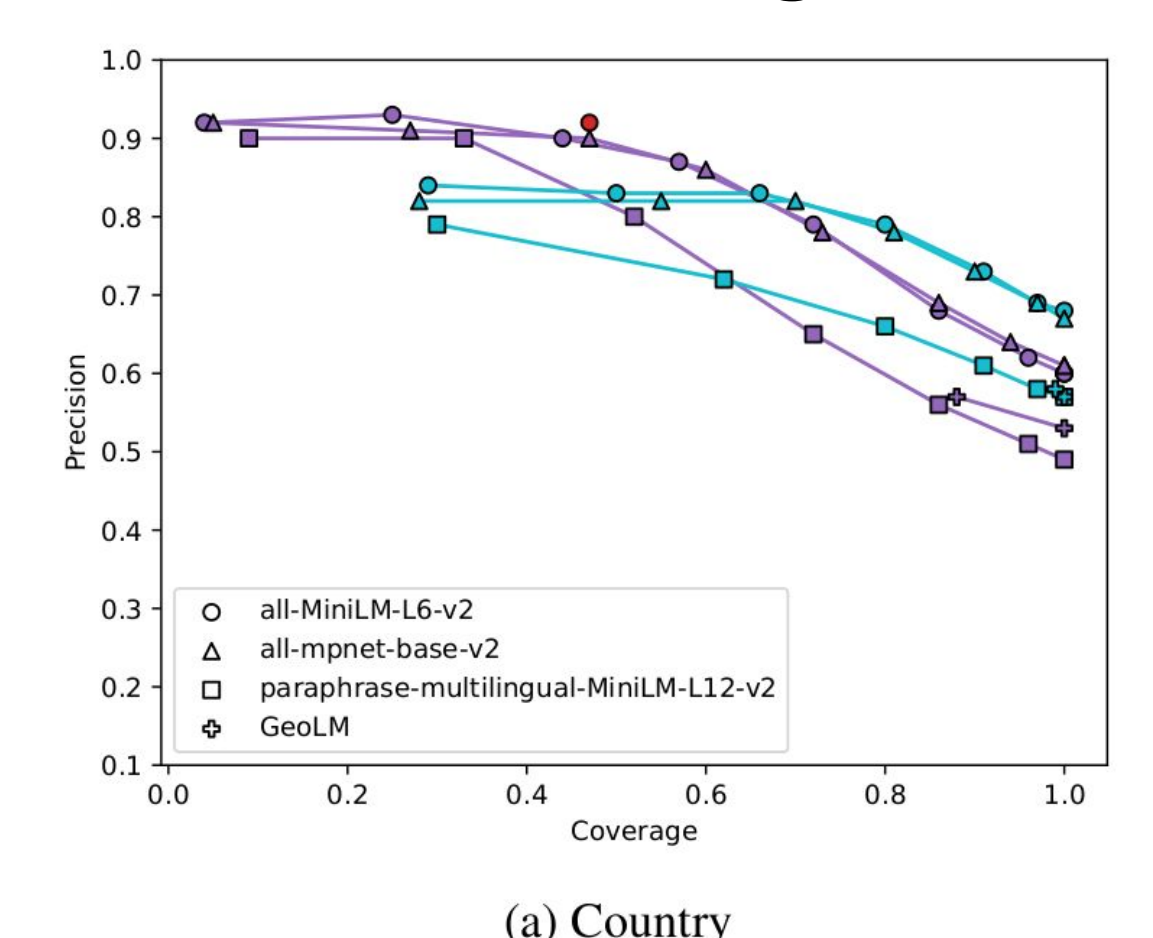
Qualitative Analysis

Location field	UserGeo	NameGeo	Carmen
福島県いわき市	(Iwaki, Fukushima, Japan)	(Zhongshu, Yunnan, China)	NULL
Catskills	(Hyde Park, New York, US)	(Catalca, Istanbul, Turkey)	NULL

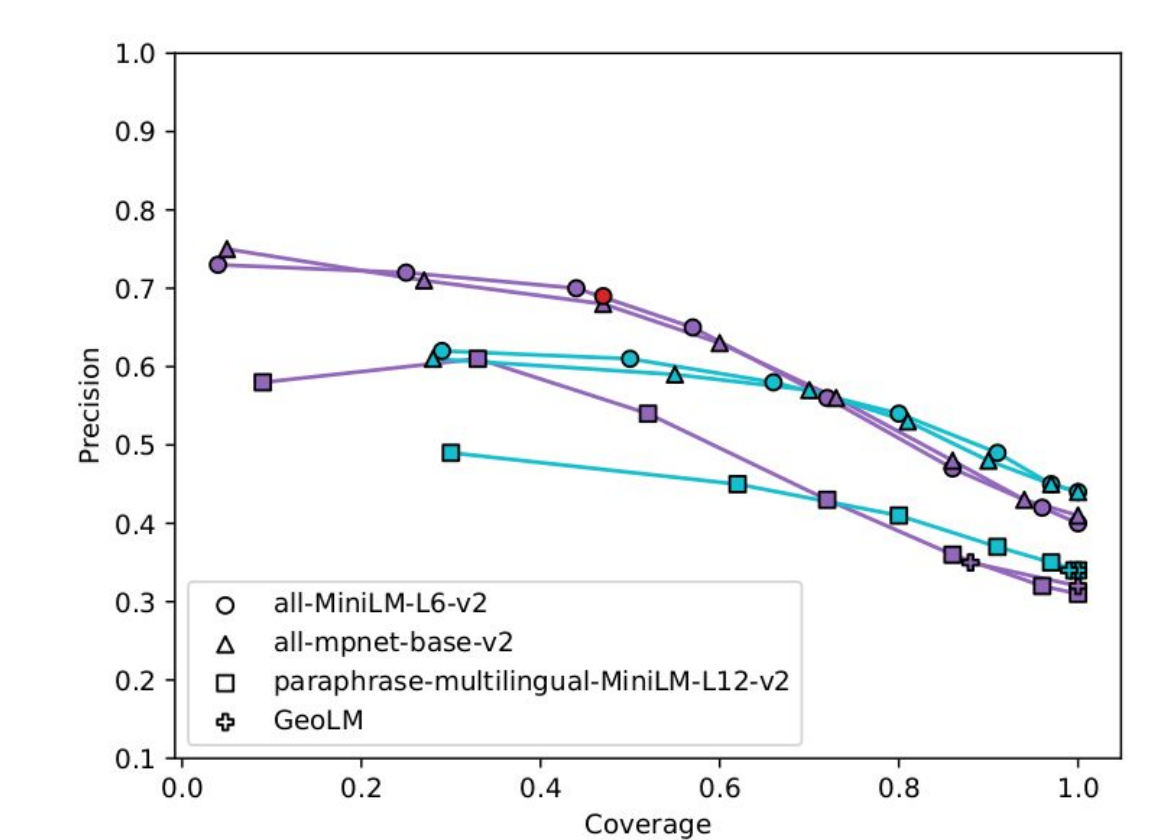
Accuracy

Method	Country	Admin.	City
CARMEN 2.0	43.5	27.3	9.8
<i>all-MiniLM-L6-v2</i>			
NameGeo	59.8	37.7	14.3
UserGeo	67.8	44.2	14.8
<i>all-mpnet-base-v2</i>			
NameGeo	60.9	38.3	14.9
UserGeo	67.4	43.7	13.9
<i>paraphrase-multilingual-MiniLM-L12-v2</i>			
NameGeo	48.7	28.9	8.1
UserGeo	57.0	34.3	9.4
GeOLM			
NameGeo	52.5	30.5	12.1
UserGeo	57.4	33.9	10.7

Precision-Coverage Curves



(a) Country



(b) Admin.

How Many Have a Real Location?

Manual examination of proportion of Location fields that reference an actual location; we use this as an **upper bound of accuracy**

	Country	Admin.	City
Upper bound	72.5	58.3	49.2
NameGeo+variants	62.0	40.9	17.0
UserGeo	67.8	44.2	14.8

Summary & Future Work

- Proposed methods for geo-entity linking noisy multilingual social media data with selective prediction
- Of two best performing methods, UserGeo achieves SOTA performance at country and administrative levels while NameGeo+variants doesn't require training data
- Identified problems with geo-entity linking at the city level for social media data
- Hope to extend to broader task of geoparsing unstructured text