**BREAKING NEWS!!!**

Incident outside **Melbourne's Sully's Backstreet Bar**. Suspect taken to the **Brevard County Jail**.

# GOAL: Geolocation of text.

## Geoparsing Pipeline:

**NER GEOTAGGING**

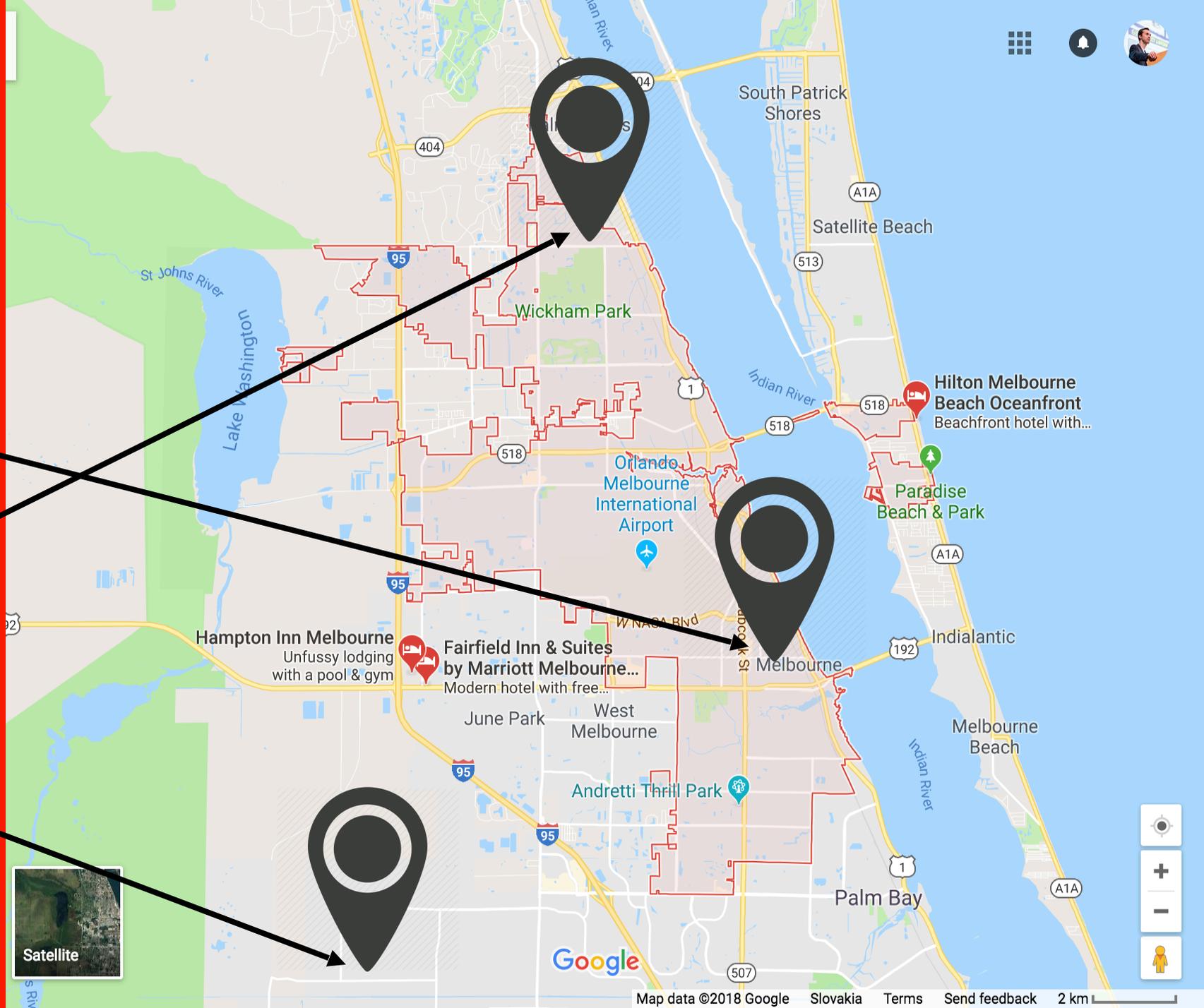**NEL + WSD GEOCODING**

Reference Geocoding

Document Geocoding

### Geocoding or Toponym Resolution

**BREAKING NEWS!!!**

Incident outside Melbourne's Sully's Backstreet Bar. Suspect taken to the Brevard County Jail.

**Background**

**Geoparsing Systems**

**Geocoding similar to WSD but…**

- Ambiguity of toponyms greater (e.g. 10+ Melbournes in the world)
- Contextual clues not adequate or missing for small (local) places
- Often difficult for humans to judge
- 50% - 75% resolved by population

1 TEXT DOCUMENT, 2 SETS OF FEATURES

Lexical Footprint

Geographic Footprint

The Map Vector

(reshape to) 1D Map Vector

**Bag of words.**

| BAR | SUSPECT | OUTSIDE | PLACE | INCIDENT | SULLY'S | TAKEN | MELBOURNE |
|---|---|---|---|---|---|---|---|
| 0.5 | 0.6 | 1.0 | 0.1 | 0.4 | 0 | 0.2 | 0.9 |

LONGITUDE

LATITUDE

180 DEGREES

| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|---|
| 0 | 0.8 | 0 | 0.2 | 1.0 | 0.6 | 0 | 0 |
| 0 | 0.2 | 0.3 | 0 | 0 | 0 | 0.9 | 0 |
| 0 | 0.1 | 0 | 0 | 0.1 | 0 | 0 | 0 |
| 0 | 0 | 0.4 | 0 | 0.9 | 0.5 | 0.6 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

360 DEGREES

**Bag of locations.**

## Evaluation Datasets

- **LOCAL GLOBAL LEXICON (LGL)** by (Lieberman et al. 2010) – packaged with our code.

- 588 local news articles from global sources

- 4460 annotated places, Medium Difficulty Test

- **WIKIPEDIA TOPONYM RETRIEVAL** (WikToR) by (Gritta et al. 2017) – also packaged with our code.

- Wikipedia-based geoparsing of 5,000 articles

- High Difficulty Test, 25,000+ locations in total

- Other corpora available (De Lozier et al. 2010), (Wallgrun et al. 2017), (Buscaldi and Rosso 2008), (De Oliveira et al. 2017), (Mani et al. 2010), (Eisenstein et al. 2010) but issues with cost, scope, annotation, size, type of task, completeness, etc.

- **OR RESOURCES NOT PUBLISHED WITH PAPER**

**GeoVirus.xml**

**New Dataset**

**DOWNLOAD**

**GPL V3** Free Software

**Free as in Freedom**

**WIKINEWS**

- 229 articles (August, September 2017)
- NER/Geotagging and Geocoding
- KEYWORDS: Ebola, Bird Flu, Swine Flu, AIDS, Mad Cow Disease, many more. (Medisys JRC)
- Locations: 2,167, Word Count: 63,205
- **https://github.com/milangritta**

# OVERALL PERFORMANCE COMPARISON

| Geocoder | System configuration | | Dataset | | | Average |
|---|---|---|---|---|---|---|
| | **Language Features** | **+** **MapVec Features** | **LGL** | **WIK** | **GEO** | |
| **CamCoder** | CNN | MLP | **0.22** | **0.33** | **0.31** | **0.29** |
| Lexical Only | CNN | — | 0.23 | 0.39 | 0.33 | 0.32 |
| MapVec Only | — | MLP | 0.25 | 0.41 | 0.32 | 0.33 |
| Context2vec[†] | LSTM | MLP | 0.24 | 0.38 | 0.33 | 0.32 |
| Context2vec | LSTM | — | 0.27 | 0.47 | 0.39 | 0.38 |
| Random Forest | MapVec features only, no lexical input | | 0.26 | 0.36 | 0.33 | 0.32 |
| Naive Bayes | MapVec features only, no lexical input | | 0.28 | 0.56 | 0.36 | 0.40 |
| Population | — | — | 0.27 | 0.68 | 0.32 | 0.42 |

| Geocoder | | Area Under Curve[†] | Average Error[‡] | Accuracy@161km |
|---|---|---|---|---|

| | | | | |
|---|---|---|---|---|
| Context2vec | LSTM | — | 0.27 0.47 0.39 | 0.38 |
| Random Forest | MapVec features only, no lexical input | | 0.26 0.36 0.33 | 0.32 |
| Naive Bayes | MapVec features only, no lexical input | | 0.28 0.56 0.36 | 0.40 |
| Population | — | — | 0.27 0.68 0.32 | 0.42 |

GeoVirus.xml

Melbourne + Perth + Newcastle

The Map Vector
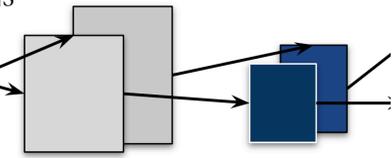
CamCoder

Lexical CNN Geocoder

CONTEXT WORDS CONVOLUTIONS PC

pyramid complex, archeological site, outskirts
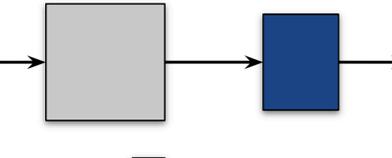
LOCATION MENTIONS

Giza, Egypt, Giza Plateau

TARGET ENTITY

CAIRO

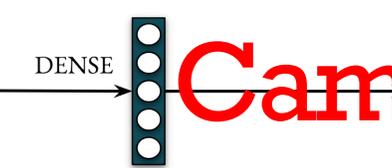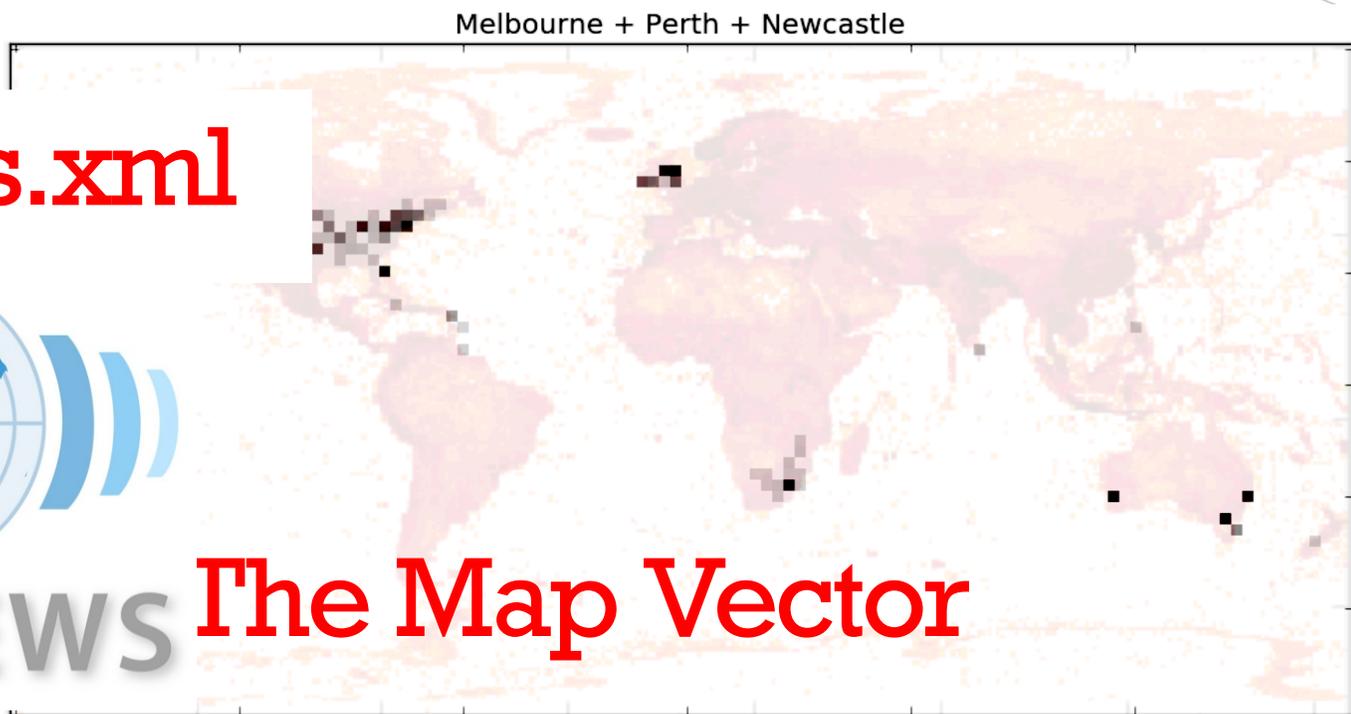MAPVEC

DENSE

DENSE SOFTMAX (LOCATION PREDICTION)

| Geocoder | Area Under Curve† | | | Average Error‡ | | | Accuracy@161km | | |
|---|---|---|---|---|---|---|---|---|---|
| | LGL | WIK | GEO | LGL | WIK | GEO | LGL | WIK | GEO |
| **CamCoder** | **22 (18)** | **33 (37)** | **31 (32)** | 7 (5) | **11 (9)** | **3 (3)** | **76 (83)** | **65 (57)** | **82 (80)** |
| Edinburgh | 25 (22) | 53 (58) | 33 (34) | 8 (8) | 31 (30) | 5 (4) | **76** (80) | 42 (36) | 78 (78) |
| Yahoo! | 34 (35) | 44 (53) | 40 (44) | **6 (5)** | 23 (25) | **3 (3)** | 72 (75) | 52 (39) | 70 (65) |
| Population | 27 (22) | 68 (71) | 32 (**32**) | 12 (10) | 45 (42) | 5 (**3**) | 70 79 | 21 14 | 0 80 |
| CLAVIN | 26 (20) | 70 (69) | 32 (33) | 13 (9) | 4 | 80 | 16 | 80 |
| GeoTxt | 29 (21) | 70 (71) | 33 (34) | 14 (9) | 47 (45) | 6 (5) | 68 (80) | 18 (14) | 79 (79) |
| Topocluster | 38 (36) | 63 (66) | NA | 12 (8) | 38 35 | NA | 63 (71) | 2 (20) | NA |
| Santos et al. | NA | NA | NA | 8 | NA | NA | 71 | NA | NA |

WIKINEWS