# Matching Flickr Photos to City Districts in Vienna

Germaine Pötgen (11738607)

Robert Simetzberger (11775655)

Simon Groß (12238274)

# Content

01 - Motivation

02 - Data

03 - Methods

04 - Results and Discussion

05 - Conclusion and Limitations

06 - Future Work

# Motivation

- Explore computation of document embedding

- Implicit geo-data: geo-data that do not have geographic coordinates in their initial form

- (How) Can we georeference Flickr photos (implicit geo-data) on district level using Doc2Vec?

- How can we identify references to location when trying to georeference Flickr photos?
  - Title
  - Description
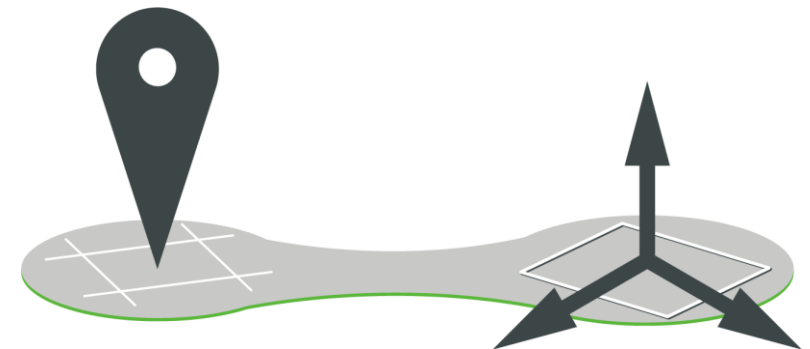  - Comments
  - Tags
  - Generated image description

Fig 1: Georeferencing. *Source*: https://developer.vuforia.com/library/develop-area-targets/geolocation-external-prior-area-targets/georeferencing-area-targets.

# Data

- *Flickr API*
  - 7250 geotagged photos inside the bounding box of Vienna
  - No limitations regarding number of photos per user
  - 'id', 'title', 'description', 'tags', 'latitude', 'longitude', 'context', 'url_c', 'owner', 'comments'


Fig 2: Flickr Logo. *Source*: https://www.flickr.com/.

- *Data.gv.at*
  - Tourist selection of the most important POIs in Vienna
  - 135 POIs in the categories:
    - Sightseeing, Museums, Gastronomy, Nightlife, Music, Shopping, Cafés and Restaurants
  - Districts of Vienna Polygon


Fig 3: data.gv.at Logo. *Source*: https://www.data.gv.at/.

- *Wikipedia API*
  - Wikipedia article for each district in Vienna
  - Save in dataframe


Fig 4: Wikipedia Logo. *Source*: https://de.wikipedia.org/wiki/Wikipedia:Hauptseite.

# Methods

## *Doc2Vec*

- Doc2Vec is a Model that represents each Document as a Vector

- Basic idea:
  - Act as if a document has another floating word-like vector, which contributes to all training predictions, and is updated like other word-vectors, but we will call it a doc-vector.
  - Gensim's Doc2Vec class implements this algorithm

- Our textual model will be used to compare each picture with all of the districts
  - → Get the most similar district

- Our model has to contain the text from all pictures and the wiki articles
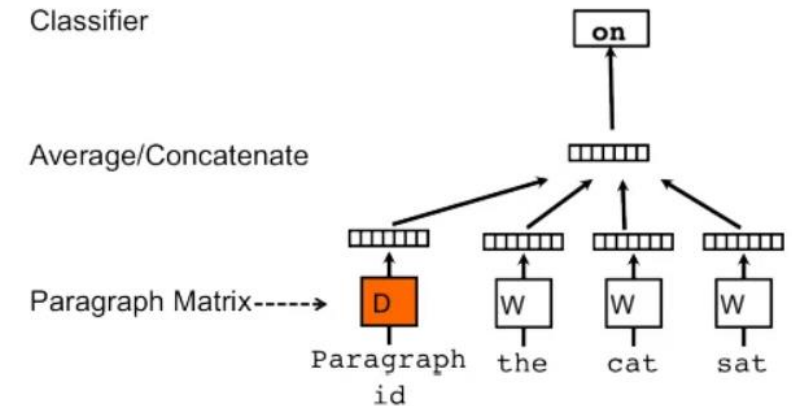


Fig 6: PV-DM model. *Source*: https://medium.com/wisio/a-gentle-introduction-to-doc2vec-db3e8c0cce5e.

# Methods

## Confusion Matrix

- Visualization of the performance of the model

- Compute Precision, Recall, and F1-score

Precision: $Precision = \dfrac{|Correctly\ identified\ place\ names|}{|all\ place\ names\ returned\ by\ the\ model|} = \dfrac{true\ positives}{true\ positives + false\ positives}$

Recall: $Recall = \dfrac{|Correctly\ identified\ place\ names|}{|all\ correct\ place\ names\ identified\ by\ humans\ |} = \dfrac{true\ positives}{true\ positives\ +\ false\ negatives}$

F1 score: $F_1 = 2 \cdot \dfrac{1}{\frac{1}{recall} + \frac{1}{precision}} = 2 \cdot \dfrac{precision \cdot recall}{precision + recall}$

# Methods

## *ML Image Processing*

- Attached metadata of a Flickr photo could be insufficient

- Image captioning model
  - Pretrained ML model from huggingface.co (Salesforce/blip-image-captioning-large)
  - Get a simple image description
  - The function pipeline from the transformers module creates a pipeline that only requires a link to the image as an input

Fig 5: An example for the pretrained image captioning model. *Source*: https://huggingface.co/nlpconnect/vit-gpt2-image-captioning.

# Results & Discussion



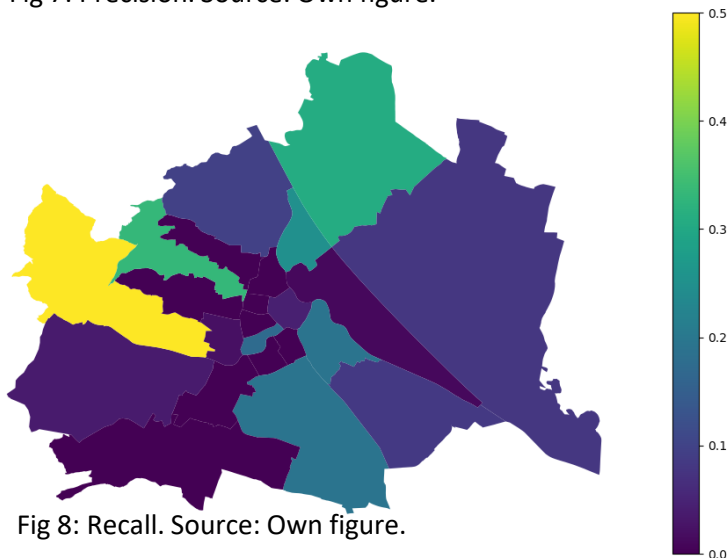Fig 7: Precision. Source: Own figure.
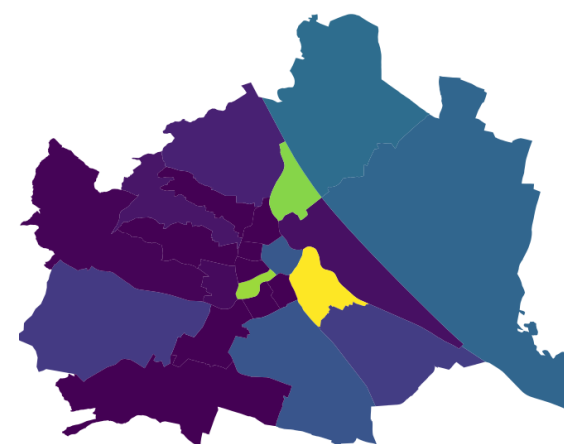


Fig 8: Recall. Source: Own figure.



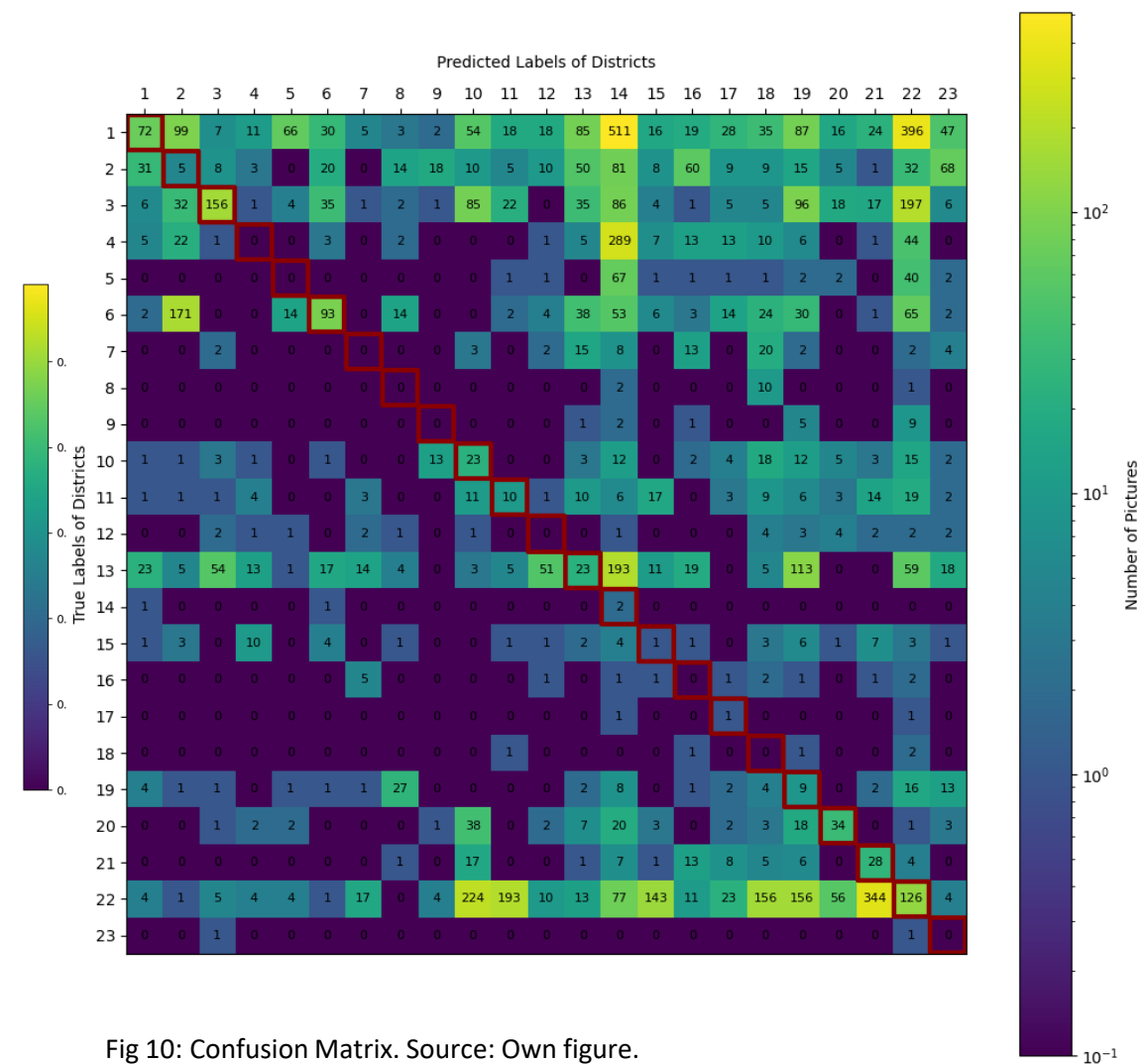Fig 9: F1 Score. Source: Own figure.



Fig 10: Confusion Matrix. Source: Own figure.

# Results & Discussion

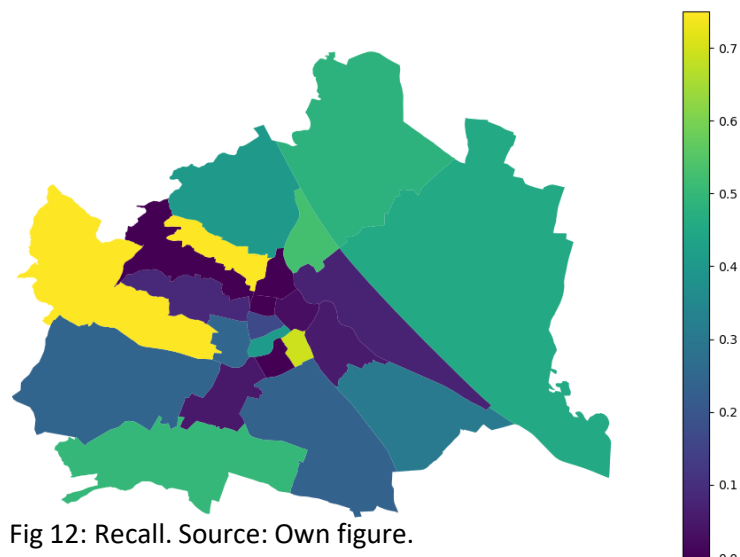Fig 11: Precision. Source: Own figure.


Fig 12: Recall. Source: Own figure.


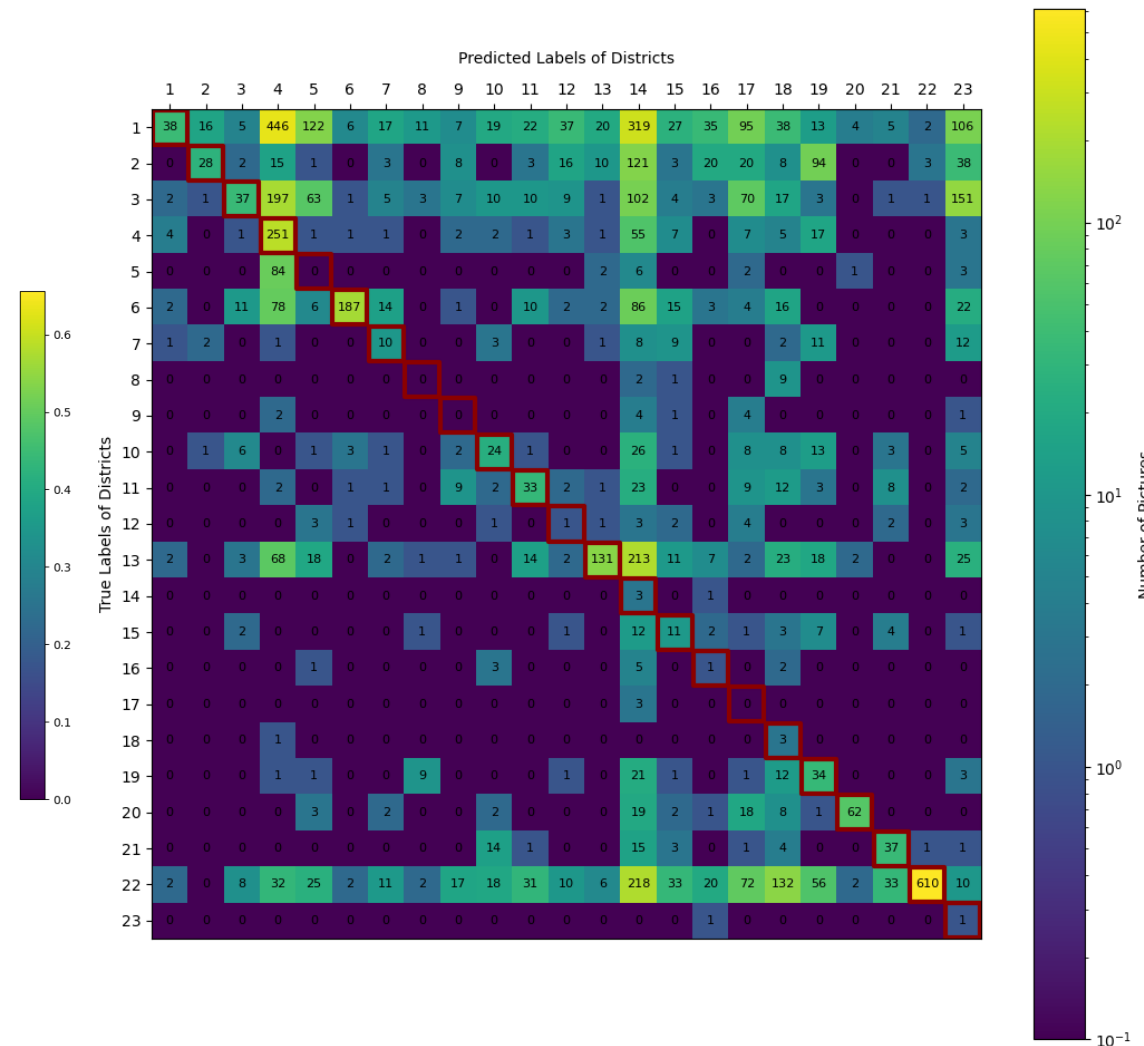Fig 13: F1 Score. Source: Own figure.


Fig 14: Confusion Matrix. Source: Own figure.

# Conclusion & Limitations

- It's possible to assign photos on district level using their textual information

- ML Image Processing wasn't used
  - ML Image Processing takes about 20-30 seconds per photo
  - Not enough computational power

- The dataset "Most important POIs in Vienna" doesn't contain POIs for every district

- A big proportion of our photos is from few users
  - 7250 photos from 225 users
  - 1604 photos from one user
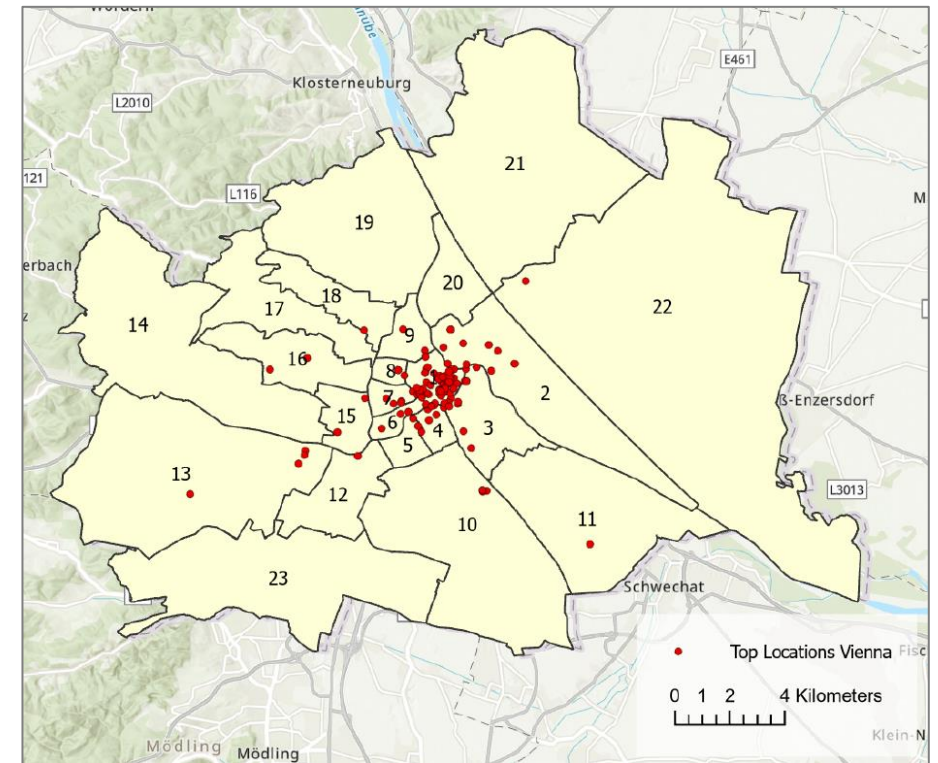  - Train-test split should be conducted without splitting users



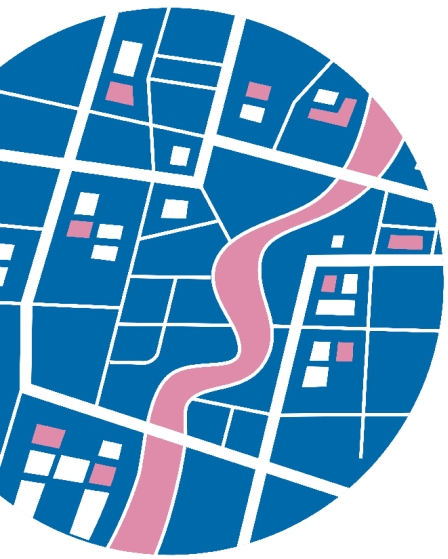Fig 15: Most important POIs in Vienna. *Source*: Own figure.

# Future Work

- Geographic Entitiy Recognition
  - For Flickr photo description, tags, title etc.
- Explore further how ML Image Processing works in this context
- Use further datasets (apart from Wikipedia articles and POIs) for the similarity comparison
  - E.g. OSM data
- Use geotagged photos to train a model to georeference photos without a geotag

# References

GENSIM (Hrsg.) (n.d.). Doc2Vec Model. Available online at https://radimrehurek.com/gensim/auto_examples/tutorials/run_doc2vec_lee.html (accessed 1/21/2024).

Hugging Face (Hrsg.) (n.d.). nlpconnect/vit-gpt2-image-captioning. Available online at https://huggingface.co/nlpconnect/vit-gpt2-image-captioning (accessed 1/21/2024).

Questions?