# SRAI: Towards Standardization of Geospatial AI

Piotr Gramacki*
Kacper Leśniara*
Kamil Raczycki*
Szymon Woźniak*
Wrocław University of Science and Technology
Department of Artificial Intelligence / Kraina.AI
Wrocław, Poland
piotr.gramacki@pwr.edu.pl
{klesniara,kraczycki,swozniak}@kraina.ai

Marcin Przymus
Piotr Szymański
Wrocław University of Science and Technology
Department of Artificial Intelligence / Kraina.AI
Wrocław, Poland
mprzymus@kraina.ai
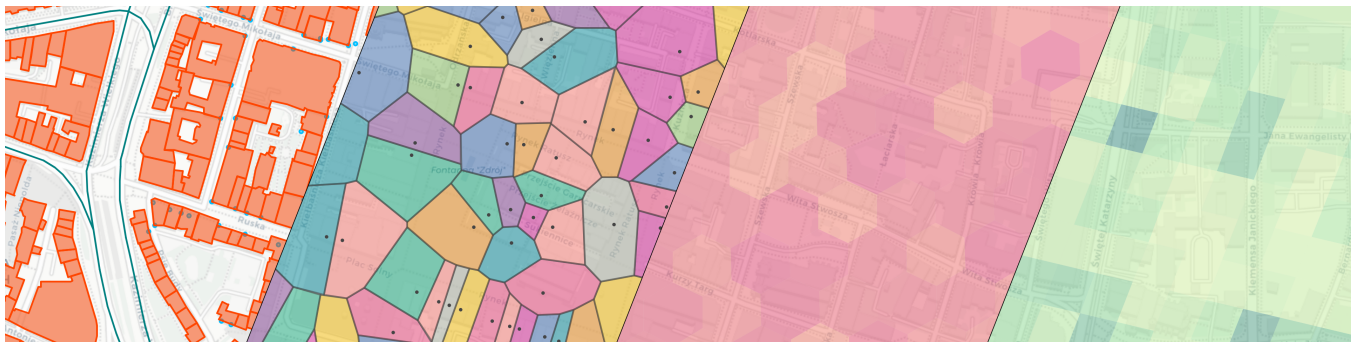piotr.szymanski@pwr.edu.pl

Figure 1: The *srai* library features overview: data acquisition, area regionalization, and shape agnostic embedding methods.

## ABSTRACT

Spatial Representations for Artificial Intelligence (*srai*) is a Python library for working with geospatial data. The library can download geospatial data, split a given area into micro-regions using multiple algorithms and train an embedding model using various architectures. It includes baseline models as well as more complex methods from published works. Those capabilities make it possible to use *srai* in a complete pipeline for geospatial task solving. The proposed library is the first step to standardize the geospatial AI domain toolset. It is fully open-source and published under Apache 2.0 licence.

## CCS CONCEPTS

• **Software and its engineering** → **Software libraries and repositories**; *Application specific development environments*; • **Computing methodologies** → **Spatial and physical reasoning**; **Machine learning**; **Learning latent representations**; • **Information systems** → **Geographic information systems**.

*Equal contribution, authors in alphabetic order.

## KEYWORDS

geospatial data processing, spatial embeddings, urban data embeddings, openstreetmap embeddings, spatial representation learning, standardization in geospatial domain, python library

## 1 INTRODUCTION

Spatial Representations for Artificial Intelligence (*srai*) is a fully open-source Python module integrating many geo-related algorithms in a single package with a unified API. It aims to provide simple and efficient solutions to geospatial problems - accessible to everybody and reusable in various contexts. Many tasks can benefit from enrichment with geospatial data. Our main goal is to unify geospatial data processing and models, making reproducing experiments and sharing new models easier. We have already included previously published models for geospatial representation learning. The library is available on GitHub[1] where one can find a link to pre-trained models download.

With *srai*, we want to lay the groundwork for unifying geospatial models for representation learning and make it possible to reuse them easily for different tasks. We hope that more researchers

[1]https://github.com/kraina-ai/srai

follow this trend and that the *srai* library will take the first steps to standardize the geospatial AI domain.

## 1.1 Geospatial data and tasks

The geospatial domain has recently gained popularity in academic research and business applications. The vast amount of data available becomes hard to analyze using traditional GIS methods. Companies try to leverage internal data and combine them with external sources to gain more insight for location intelligence tasks such as site selection, market analysis, geo-marketing, supply chain optimization, fraud detection, real-time traffic management, and navigation[2]. Governments and academic researchers try to improve safety and quality of life by modelling the risk of natural disasters[43], analyzing air pollution[16, 19], predicting forest fires[31], analysis of public transport networks[21] or traffic[44] and many more.

To be used in a geospatial context, data must contain information that either represents a geometry with coordinates or can be geocoded to a region or a location. That said, a lot of data available can be included in geospatial tasks. Those can take the form of text (tags, reviews, social media posts, documents), images (satellite multispectral and hyperspectral imagery, aerial photos, street view images, and videos), and numeric values (weather readings, number of residents in a household, average time spent in traffic, POI ratings).

The *srai* library focuses mainly on publicly available datasets with clearly defined vector geometry (points, lines, polygons) to increase reproducibility, simplify usage and democratize the geospatial domain for everyone without the need to be an expert in the GIS tools usage.

## 1.2 OpenStreetMap

Regarding geospatial data, one of the most critical open data sources is OpenStreetMap[]. It is a free, editable map service that allows people worldwide to post updates and improve data coverage. It contains a wide range of tags, which allows for selecting a subset suitable for a given task. We chose it as the primary data source available in *srai* and provided multiple ways of accessing it.

## 1.3 Representation learning

We designed the *srai* library around representation learning approaches to geospatial data. Representation learning is finding a function $f : D \rightarrow \mathbb{R}^d$, which transforms data $D$ into a $d$-dimensional vector. Similar to Natural Language Processing (NLP)[20, 26] and Computer Vision (CV)[6, 8], this approach has already proven helpful in solving geospatial tasks[3, 40, 42]. A significant benefit of such representations is that they are reusable. They can either be transferred from task to task or defined as general enough to be used in multiple further downstream tasks (similarly to pre-trained transformer models in NLP).

The remainder of this paper is organized as follows: Section 2 presents an overview of available solutions in the geospatial area and positions the *srai* library in comparison to them; Section 3 describes the library and explains the design choices we made; Section 4 presents the use cases for *srai*, to show how it can improve geospatial tasks solving; Section 5 discusses future improvements

in *srai* to show the vision for the library; Section 6 concludes this paper and discusses the limitations of *srai* in its current state.

## 2 BACKGROUND

This section provides an overview of existing methods for geospatial analysis. It consists of tools we integrated into *srai* and other libraries that simplify geospatial data processing. We also show how *srai* compares to those alternatives.

This review is composed of two parts: (i) a description of geospatial tools that power *srai* library and (ii) positioning of *srai* relative to geospatial tools and systems of different types.

## 2.1 Ecosystem behind srai

In this part, we describe the tools (both geospatial and not) we used to make *srai* versatile and capable. We divided them into three categories: (i) *Geospatial operations*, (ii) *Data acquisition and processing*, and (iii) *Machine Learning tooling*. We describe each of them in detail in the subsequent sections.

*2.1.1 Geospatial operations.* The main workhorse in *srai* is the GeoPandas[4] library. It is a library compatible with pandas, which extends it with geospatial operations. The ease of integration with other tools and versatile capabilities was why we selected this library. Geometries are defined using Shapely[9], which is the most popular geometry-related library in the Python domain and GeoPandas relies on it heavily. For regionalization (also referred to as tessellation in the literature), we use either our own implementations or pre-defined hierarchical spatial indexes - H3[35] and S2[10]. Both of them are widely utilized in the community, especially H3, which is why we included them.

*2.1.2 Data acquisition and processing.* The second group of tasks handled using external libraries in *srai* is downloading and processing spatial data. The primary data source in *srai* is OpenStreetMap, which we access using the OSMnx[2] library. For loading OSM data from PBF files, we use pyosmium[15] and DuckDB[27]. Reading a vast majority of geospatial file formats is enabled by GeoPandas.

Another data source available in *srai* is GTFS timetables, which we process using the GTFS Kit[29] library. It provides a wide range of operations, which lifts operating on relational GTFS schemes off our shoulders and provides calculations for basic public transport statistics. One of its more convenient capabilities is GTFS feed validation, which we also use in *srai*.

*2.1.3 Machine Learning tooling.* Apart from geospatial operations, we also use multiple libraries for data science and machine learning. We use widely adopted libraries to make *srai* compatible with the most popular tools and frameworks. Also, leveraging well-known libraries makes our library easier for first-time users. Our operations are backed by: pandas[34] and NumPy[14] for data processing and calculations; scikit-learn[25] for machine learning algorithms; NetworkX[13] for graph processing; PyTorch[24] with PyTorch Lightning[5] for deep learning models.

## 2.2 Positioning of srai

Let us juxtapose *srai* with existing tools for geospatial applications. We have identified a wide selection of libraries and systems in the geospatial domain and categorized them into four categories: (i) *GIS*

---

**Table 1: Geospatial toolboxes and their capabilities.**

| Library | Spatial files | OSM | Trajectories | GTFS | Raster | Visualization | Regionalization | Geocoding | ML | Datasets |
|---|---|---|---|---|---|---|---|---|---|---|
| geowrangler[1] | ✓ | ✓ | | | ✓ | | ✓ | | ✓ | |
| tesspy[2] | | ✓ | | | | ✓ | ✓ | ✓ | | |
| geomancer[3] | | ✓ | | | | | | | | |
| Mosaic[4] | ✓ | | | | ✓ | | ✓ | | ✓ | |
| PySal[30] | | | | | ✓ | ✓ | ✓ | | ✓ | |
| Verde[36] | ✓ | | | | | | ✓ | | ✓ | ✓ |
| WhiteboxTools[6] | ✓ | | | | ✓ | | | | ✓ | |
| Pandana[5] | ✓ | ✓ | | | | | | | | |
| MovingPandas[12] | | | ✓ | | | | | | | |
| Scikit mobility[23] | | | ✓ | | | | | | | |
| segment-geospatial[41] | ✓ | | | | ✓ | | | | ✓ | |
| TorchGeo[33] | ✓ | | | | ✓ | | | | ✓ | ✓ |
| srai | ✓ | ✓ | | ✓ | ✓* | ✓ | ✓ | ✓ | ✓ | |

[1] https://github.com/thinkingmachines/geowrangler, [2] https://github.com/siavash-saki/tesspy, [3] https://github.com/thinkingmachines/geomancer, [4] https://github.com/databrickslabs/mosaic, [5] https://github.com/UDST/pandana, [6] https://github.com/jblindsay/whitebox-tools, * only for data downloading and preparation

*Platforms*, (ii) *Location Intelligence Platforms*, (iii) *Spatial Databases and Computing Engines* and (iv) *Toolboxes*.

*GIS Platforms* are great, versatile tools for geospatial processing. Those are complete solutions that cover many use cases in the domain. The most popular and widely adopted, according to our review, are: ArcGIS[3], QGIS[4] and Google Earth Engine[5]. Those systems have a different scope than the *srai* library - ArcGIS and QGIS tend to be an all-purpose solution to geospatial data, while Google Earth Engine's focus is on handling raster data and providing access to a variety of geospatial datasets. These tools offer APIs in either Python or JavaScript or use them as scripting languages inside the tools themselves. Regarding pricing: QGIS is free and open-source, Google Earth Engine is free for noncommercial use, and ArcGIS is a paid solution. Compared to fully-fledged GIS applications, the *srai* library is free and open-source. It depicts itself more as a specific set of functionalities targeted at the geospatial ML domain and simplifying data downloading and wrangling. The advantage is that it is accessible for people without a GIS background, such as data scientists and Python developers.

*Location Intelligence Platforms* is a set of systems that collect data and enable analytics, mainly in the Location Intelligence domain. They prioritize cloud-based analytics, spatial data enrichment, and online visualizations. Examples are CARTO[6] and Foursquare[7], which focus more on location intelligence, as well as Mapbox[8] which is primarily a map and navigation provider. All solutions provide a data catalogue with geo-related datasets from domains such as traffic, movement, demographics, weather, retail, etc. Additionally, those platforms expose APIs for creating their map-based applications. Access to these platforms requires a paid subscription. Those platforms are very popular in the geospatial domain, especially in the corporate world, but need more machine learning functionalities. The *srai* library allows for creating more sophisticated machine learning models. Additionally, using paid, quality datasets from those platforms could create more advanced solutions in business applications and extend the capabilities of those platforms.

*Spatial Databases and Computing Engines* are dedicated solutions for large-scale geospatial computations and data storage. Examples include PostGIS[9] which is a spatial extension to PostgreSQL[10] and Apache Sedona[11] which is a highly efficient solution for spatial data based on Spark[12]. There are also cloud-based data warehouses capable of processing geospatial data, such as BigQuery[13], Athena[14] and Snowflake[15], which simplify set-up and maintenance. Alternatively, large-scale datasets can be handled using either parallelism-capable libraries (e.g., Dask-GeoPandas[16]) or more performant ones (e.g., GeoPolars[17]) All of the mentioned solutions could be treated as possible backends for the *srai* library regarding storage and/or computation, as it is with GeoPandas and DuckDB.

*Toolboxes* is a coherent collection of functionalities dedicated to a limited set of use cases. We have gathered a wide selection of such tools, compared in Table 1. Most of the depicted tools individually contain a subset of the *srai* functionalities. However, none offers complete coverage of the geospatial AI pipeline. The library most resembling *srai* is wrangler. It offers similar data loading and processing capabilities but is less ML-oriented and offers fewer utility functions such as plotting or geocoding. Mosaic also overlaps with *srai*. However, it is only available on the Databricks platform. We

have identified areas that are not covered in *srai* - trajectories (MovingPandas and scikit-mobility) and raster-focused data processing and machine learning (TorchGeo and segment-geospatial). We have also found libraries focused primarily on feature engineering and analytics - geomancer, PySal, Pandana, WhiteboxTools - or pure regionalization - tesspy. In comparison to the toolboxes mentioned above, *srai* introduces a needed enhancement to the geospatial tooling available to data scientists and developers.

## 3 THE SRAI LIBRARY

This section explains the design of the library and lists methods which are already available in *srai*.

### 3.1 General library design

The *srai* library design aims to bring different geospatial operations under a common interface. With that achieved, different algorithms will be easily interchangeable in the geospatial task pipeline. We define four main components:

(1) Loader - responsible for loading geospatial data (features) from a given source and pre-processing it to a format used in *srai*,
(2) Regionalizer - able to split a given area into micro-regions used in further analysis,
(3) Joiner - used to match loaded features with regions based on a given spatial predicate,
(4) Embedder - which embeds regions into a vector space based on features matched to them.
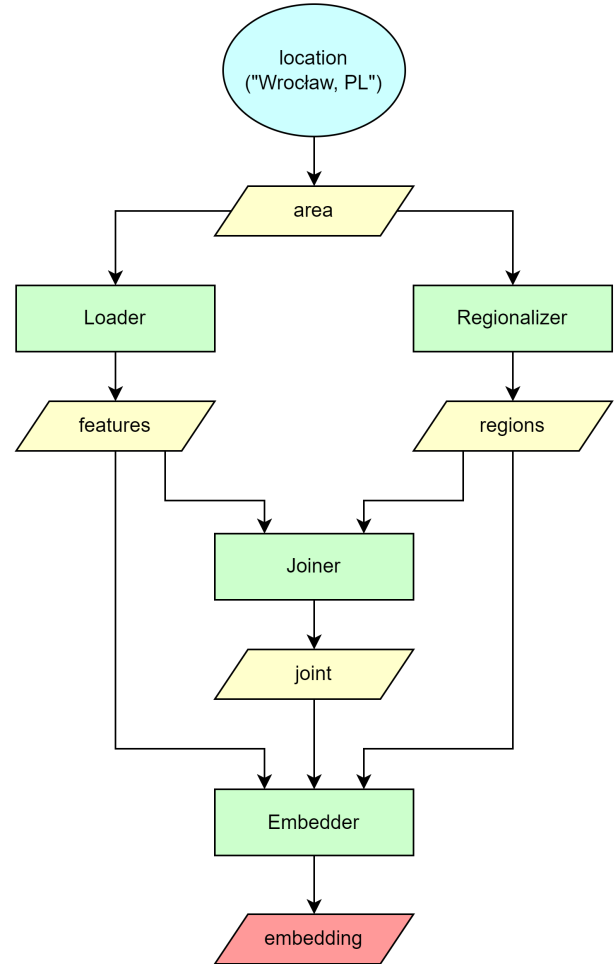
Those components create a complete pipeline for learning representations for geospatial data. Solving the final task can be obtained using any machine learning or deep learning library since *srai* returns embeddings as pandas[34] data frames, which are very versatile. The embedding step can be omitted, using a dedicated component that only summarises each region's features. A general overview of how those components are connected is presented in Figure 2. We start with location definition, which we can geocode using our utility functions. Next, we use a *Loader* to get geospatial features (for example, from OSM) from the obtained area and a *Regionalizer* to split this area into micro-regions (for instance, H3 cells). Regions are combined with features using a *Joiner*, and the results are passed over to an *Embedder*, which produces the final embeddings.

### 3.2 Functionalities in srai

The *srai* library implements all steps of the representation learning approach based on OSM data. In almost all steps, there are multiple options implemented, which are interchangeable thanks to unified interfaces.

*3.2.1 Loaders.* The first step is to gather data. In *srai*, our main source of data is OpenStreetMap, and we provide tools to download different data from there:

- OSM tags - downloading tags specified by a filter for a given area. We implement queries to OSM online (suitable for a small number of tags) and from *PBF* files (more efficient in cases with more tags). We provided a convenient functionality for users that automatically selects the smallest
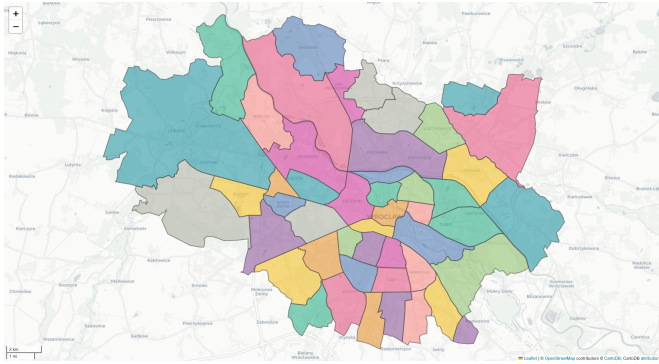


**Figure 2: Embedding calculations for a given location using the *srai* library.**

extract with OSM data to download, based on the given geometries, and downloads it either from Geofabrik[18] or OpenStreetMap[22] hosting services.
- OSM networks - downloading structured networks as a graph. It can download different types of networks (roads, bike paths, etc.),
- OSM tiles - downloading tiles with maps as images. It works with different zoom levels and tile providers,
- GTFS - loading data from GTFS feeds and pre-computations of public transport offer for stops (based on gtfs2vec[11] paper),
- GeoParquet - loading geospatial data from *geoparquet* files, with pre-processing required by other library components.

*3.2.2 Regionalizers.* The next step is to split the area of interest into micro-regions, for which the data will be aggregated, and representations will be computed. For this task, we implement support for spatial indexes as well as data-based division approaches:

**Figure 3: Wrocław city districts downloaded by Administrative Boundary Regionalizer.**

- H3 and S2 - spatial indices with different shapes, hexagons, and squares, respectively. Both are deterministic and implemented hierarchically.
- Voronoi - allows the given areas to be divided into Thiessen polygons using seed geometries,
- Administrative boundaries - splits an area based on administrative boundaries on different levels available in OSM (see example in Figure 3),
- Slippy map - splits an area into regions, which match with OSM tiles from the corresponding loader.

*3.2.3 Joiners.* Features returned by the loader should be matched with regions if regionalization is used in the pipeline. This can be done, with different spatial predicates; in *srai*, we implement one which covers most of the use-cases - intersection.
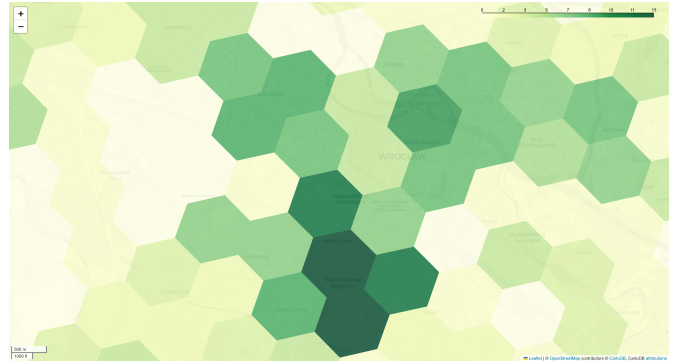
*3.2.4 Embedders.* The key functionality of *srai* is supporting geospatial data embedding. We provide multiple embedders:

- Count - baseline embedder, which simply counts feature occurrences in regions,
- Contextual Count - counts features in regions with added information from neighbour regions (used by [28]),
- Hex2Vec - embedder from hex2vec[39] paper, which works with OSM tags,
- Highway2Vec - implementation from highway2vec[17] paper of an embedder for road segments,
- GTFS2Vec - public transport offer embedding, based on the method described in the gtfs2vec[11] paper.

*3.2.5 Pre-trained models.* Some of the embedders mentioned above require model training. This process can be resource-consuming, making the experiment challenging to reproduce. To overcome this problem, we include an option to save a pre-trained embedder to a file and quickly load it. We also share a selection of pre-trained hex2vec models and intend to extend this list in time.

*3.2.6 Neighbourhoods.* Some of the geospatial tasks can benefit from including neighbourhood information. At the moment, two different neighbourhood implementations are available in *srai*:

- Adjacency - based on regions touching each other
- H3 - based on the H3 spatial index



**Figure 4: Sum of parks objects per H3 cell in Wrocław, Poland.**

*3.2.7 Utilities.* Apart from those main components, *srai* contains various utility functions, which speed up spatial data processing or unify data formats which are used. This includes geocoding, buffering, flattening, merging, plotting, and other operations.

### 3.3 Technical backend

We need a spatial operations backend to make all of the abovementioned operations possible in *srai*. We decided to use GeoPandas as our primary spatial backend. It is efficient for smaller use cases, which can be fitted into the available memory. One of the most important benefits is its compatibility with the pandas interface, which makes it easy to integrate with existing libraries. This high compatibility is crucial to making *srai* a standard geospatial AI library. However, this approach is limited by available memory since all data must be loaded into RAM. This limitation makes it impossible to process larger areas simultaneously, so we are actively working on alternative geospatial backends. We have already started working on utilization of DuckDB's spatial extension for most memory-consuming tasks. We provide more detail on that in Section 5.1.

### 4 USAGE EXAMPLES

This section presents examples of *srai* usage in various supported tasks. We show how the unification of interfaces makes solving geospatial tasks easier. We present examples for (i) geospatial data download and processing, (ii) unsupervised and (iii) supervised task solving, and (iv) *srai* usage for transfer learning. We present *srai* functionalities on a single task of rental price estimation. We assume that we have data from multiple cities with short-term rental offerings, containing location and price per person. We split the task into parts to present how to apply *srai* to each of them.

### 4.1 Processing pipeline

Firstly, we prepare both the input (features from OSM) and output (rental price) data. In this example, we try to predict the average price in an area defined as a single H3 cell based on leisure and amenities offered. Listing 1 presents how to do this using our library. The steps are as follows:

(1) download the boundaries of the analyzed city,
(2) define OSM tags filter and download the features from OSM,

(3) split an area into H3 cells,

(4) aggregate OSM features for cells, creating pre-training data.

In Figure 4 we can see an example result for these steps - the number of park objects in each of the hexagonal micro-regions we selected.

```python
import geopandas as gpd
from srai.loaders import OSMOnlineLoader
from srai.regionalizers import H3Regionalizer
from srai.joiners import IntersectionJoiner
from srai.utils import geocode_to_region_gdf

# download boundaries of the city
area = geocode_to_region_gdf("Wroclaw, PL")

# define which tags to download
query = {
    "leisure": ["park"],
    "amenity": ["restaurant"]
}

# use loader to download OSM data
loader = OSMOnlineLoader()
features = loader.load(area, query)

# split an area into H3 cells
regionalizer = H3Regionalizer(resolution=8)
regions = regionalizer.transform(area)

# match features to regions
joiner = IntersectionJoiner()
joint = joiner.transform(regions, features)
```

**Listing 1: Loading and processing data from OSM for a given area using srai library.**

## 4.2 Clustering

Secondly, we perform an exploratory analysis. This helps us understand how our data is distributed in the city and presents *srai* capabilities in solving unsupervised tasks. We use the hex2vec method and clustering for that. All of these steps are illustrated in Listing 2:

(1) create and configure hex2vec model,

(2) train it on OSM features downloaded earlier,

(3) use KMeans clustering to find similar areas.

Based on that result, we can make some preliminary observations on how rental price is affected by different leisure and amenities. Such an approach can also be useful in exploratory analysis and embedding method evaluation. In Figure 5 we can see a visualization of our model's embedding space and Figure 6 presents clustering results. Both of these visualizations can be useful in exploratory analysis before predictive model training.

```python
from srai.embedders import Hex2VecEmbedder
from srai.neighbourhoods import H3Neighbourhood
# use any library for clustering you want
from sklearn.cluster import KMeans

# configure hex2vec embedder
embedder = Hex2VecEmbedder(encoder_sizes=[42, 13])
neighbourhood = H3Neighbourhood(regions)

# pre-train hex2vec model on OSM data
embeddings = embedder.fit_transform(
    regions,
```

```python
    features,
    joint,
    neighbourhood,
    batch_size=64,
)

# run clustering to explore differences between regions
clustering = KMeans(5)
clustering.fit(embeddings)
regions["cluster"] = clustering.labels_
```

**Listing 2: Training a hex2vec embedder and clustering in obtained embedding space**

## 4.3 Prediction

Next, we train a regression model to estimate rental prices. We use our pre-trained hex2vec model from the previous step and combine it with the regression head. Then we train these models on rental prices, which we preprocessed using our library. This is presented in Listing 3 and consists of the following steps:

(1) load prices dataset and preprocess it using *srai*,

(2) generate train and test data,

(3) create Regressor based on hex2vec embeddings,

(4) train the model and evaluate it.

```python
import geopandas as gpd
from srai.constants import FEATURES_INDEX, REGIONS_INDEX
# use the library and regressor of your choice
from sklearn.svm import SVR
import sklearn.metrics.mean_squared_error as mse

# load rental prices and calculate avg per region
prices = gpd.read_file("prices.geojson")
prices.index.name = FEATURES_INDEX
pj = joiner.transform(regions, prices)
pj = pj.join(on=FEATURES_INDEX).reset_index(inplace=True)
prices_avg = pj.groupby(REGIONS_INDEX).agg(
    {"price": "avg"}
)

# create training data
y_train, y_test = train_test_split(prices_avg)
x_train = embeddings[y_train.index]
x_test = embeddings[y_test.index]

# create regressor from pre-trained hex2vec
regressor = SVR()

# fit regressor to predict rental price from OSM data
regressor.fit(x_train, y_train)

# test results
result = mse(y_test, regressor.predict(x_test))
```
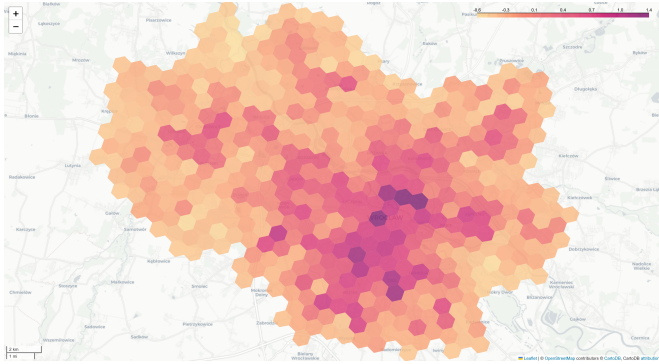
**Listing 3: Using pre-trained hex2vec embeddings to train a regressor for price estimation**

## 4.4 Transfer learning

Finally, thanks to the *srai* standardized approach we can easily transfer this model to any other city. We want to find an area with the highest potential for expensive rental in a new city, without rental data. We can do this easily with our model trained in previous steps. We only need to download data for the new region and use a pre-trained model. In this example (presented in Listing 4) we show

**Figure 5: Hex2Vec model embedding layer visualization.**



**Figure 6: Clustering based on embeddings from *srai* library.**

how to use the pre-trained hex2vec model from the previous step to run predictions on a different city, with the following steps:

1. download boundaries for a new city,
2. load OSM tags based on the filter matching pre-trained model,
3. preprocess the data following the recipe from Section 4.1,
4. make predictions using the model from the previous step.

Those are simple toy examples, but we hope they show that standardization in the GeoAI domain leads to easier access to advanced modeling tools. Ease of use and unified interfaces mean re-using models becomes almost effortless and available to everyone. We can use models pre-trained at the scale of a whole country/continent without having to create one ourselves. This also may lead to a lower carbon footprint in our experiments. Furthermore, a unified format can be an enabler for benchmarking geospatial models.

```
# get boundaries of new city
area_new = geocode_to_gdf("Hamburg, DE")

# extract OSM filter from pre-trained model
query = embedder.get_query()

# download and process OSM data, same as in stage 1
features_new = ...
regions_new = ...
joint_new = ...
neighbourhood_new = ...

# run predictions
embeddings_new = embedder.transform(
    regions_new,
    features_new,
    joint_new,
    neighbourhood_new,
)

# predict prices
prices_new = regressor.transform(embeddings_new)
```
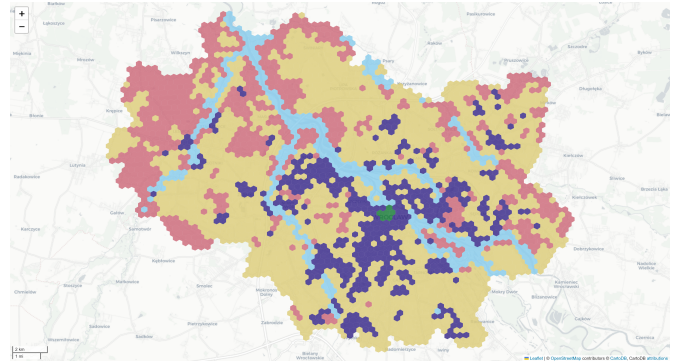**Listing 4: Transfering price regression model to another city**

## 5  FUTURE WORKS

This section provides insight into *srai* future releases. We want to show the vision which drives *srai* forward. We intend to extend to other data modalities and focus on making the reproducibility

of experiments even more accessible. Below, we present a non-exhaustive list of future functions, which will find their way to the *srai* library.

### 5.1  DuckDB with spatial extension as geospatial operations backend

To overcome the memory limitations and speed up the processing of larger datasets, we are currently implementing the most memory-intensive operations using DuckDB with spatial extension. We use it as an in-memory database (which speeds up processing) with secondary file storage (which lifts RAM constraints). This approach has a downside in decreased performance on operations requiring more memory than available. However, it is still a reasonable middle-ground between RAM-limited operations and high-performance clusters because running it locally without extensive configuration is possible.

To keep compatibility with pandas, we implement only the most memory-consuming operations with DuckDB and convert results back to GeoPandas; therefore, switching a backend is transparent from the user's perspective.

### 5.2  Pre-trained models and pre-calculated embeddings

We aim to provide pre-trained models (Hugging Face style[38]) and pre-calculated embeddings (word2vec style[20]), which will make geospatial problem-solving accessible for individuals or organizations with limited computation budgets. Another benefit is the opportunity to use transfer learning or benchmark methods using exact embeddings.

It is already possible to load pre-trained models using the *srai* library and use them for any task. However, it is still required to download a model manually. With sufficient resources for hosting and storage, we want to provide hosting for such models and make it possible to download and load them all from the code. We will explore integration with existing solutions (like Hugging Face) or design a specific solution for our library.

Pre-calculated embeddings would allow end-users to solve geo tasks without or with minimal contact with model training. Such embeddings store much information about a region they describe (as shown in hex2vec[39] or loc2vec[32] with vector addition and

subtraction) and may be sufficient to solve some geospatial tasks (e.g., similarity search or unsupervised clustering, as shown in gtfs2vec[11]). We intend to create a platform with pre-calculated embeddings for a wide selection of countries and cities, just like pre-calculated word2vec models are available for different languages.

### 5.3 Fine-tuning

Solving a downstream task is possible by using embeddings to enrich existing data or fine-tuning the whole model for this specific task. In *srai*, we plan to provide an option to add a classification/regression head to any of the available embedders. Those who are trainable could then be fine-tuned for the end task. This way, not only pre-training but also adaptation with fine-tuning will be possible using our library.

### 5.4 Other modalities of geospatial data

Another direction in which *srai* will grow is the support for different data modalities. Right now, we focus on tabular data, mainly from OSM. We already include an option to download map tiles from OSM, so an extension to computer vision (CV) embedding methods is on the roadmap. We will explore integrations with existing libraries like TorchGeo[33] or implement CV models into *srai*.

The next modality which we want to include is graphs. Geospatial data naturally form graph structures, and some models already rely on graph information for embeddings[40, 45]. We will include support for graph data download (some are already available with road networks) and dedicated models for graph embeddings designed for geospatial data.

Trajectories are another data modality that can be integrated into srai and benefit from a unified interface. This data type can use embedding methods already existing in the library and link regions that form a trajectory. Trajectories also can be used to build autoregressive models, which is an exciting area of research and would serve as a great addition to the *srai* library.

### 5.5 Spatial representations models

We will continue our work on embedding models for geospatial data and release all of our results compatible with the *srai* library. We hope that more researchers will follow this trend, and just like Hugging Face standardized NLP, the *srai* library will make the first steps to standardize the geospatial AI domain.

### 5.6 Benchmark datasets

One of the biggest challenges in any domain is the availability of diverse benchmark datasets. They are crucial for experiment reproducibility and comparison between methods. We want to introduce some benchmarks and make them accessible using *srai* library. A good option would be to utilize the Huggingface datasets[18] with dedicated support for geometries. We hope that this will lay the foundations for standardized geospatial benchmarks.

### 5.7 New sources of geospatial data

We will continue to monitor new emerging platforms for geospatial data. One auspicious example is Overture Maps Foundation[19], which provides OSM-like data in PBF format but with finer granularity and more accuracy. Limited data is available as of the time of writing this paper, but if it grows, we will extend our library. Accessing a wide selection of PBF sources under a standard interface will be highly beneficial.

## 6 CONCLUSIONS

In this paper, we introduced a new library for geospatial data processing and representation learning - Spatial Representations for Artificial Intelligence (srai). It is a Python module that integrates many geospatial operations and models under a unified API. Its main goal is to unify the process of geospatial data processing and set standards for geospatial model distribution. Hopefully, our library will take the first steps to standardize the geospatial AI domain.

We extensively reviewed available tools for the geospatial domain and categorized them based on their primary usage. We presented how the *srai* library is positioned against them, which gaps it fills, and how wide the range of tasks it covers. We see that there is a need for this library thanks to its versatility and unified API.

We presented in detail how we designed our library and justified all architectural choices. We described each functionality group and listed available algorithms and models already implemented in the *srai* library. They cover the whole pipeline for most geospatial tasks based on OSM data. We followed that with usage examples demonstrating the capabilities of the *srai* library.

Finally, we shared our view on the future of the library. We intend to expand it in multiple areas, bringing us closer to the primary goal of geospatial AI domain standardization. We believe our work will benefit the whole community and be the foundation for geospatial algorithms and models' unification, reproducibility, and reusability.

### 6.1 Limitations

Right now, due to the costs of computing and hosting, we do not provide pre-calculated embeddings. We hope to solve those issues and provide an option to share and use pre-computed vectors using the *srai* library.

As mentioned in this paper, *srai* currently supports tabular data, which may not always be sufficient for solving all geospatial tasks. However, we made sure that our methods were expandable by using widely adopted data types and libraries.

The *srai* library relies on GeoPandas[4] for most of the operations, which creates some limitations regarding RAM usage. Currently library must fit the whole spatial dataset into the memory for operations. An area under analysis can be limited by the available resources. There are available multiple solutions that can allow for out-of-memory calculations: dask-geopandas[37] which joins geopandas functionalities with dask capabilities created by the geopandas developers team; Apache Sedona[1] library with geospatial operations and the option to connect it to the Apache Spark cluster; DuckDB[27] in-memory database with spatial extension[7];

---

[19]https://overturemaps.org/

and GeoArrow[20] backend for zero-copy operations on geospatial objects. The final option, be it using only one option or implementing multiple backends under a unified API is still to be decided.

## ACKNOWLEDGMENTS

## REFERENCES

[1] Apache Software Foundation. [n. d.]. *Apache Sedona*. https://sedona.apache.org/
[2] Geoff Boeing. 2017. OSMnx: New methods for acquiring, constructing, analyzing, and visualizing complex street networks. *Computers, Environment and Urban Systems* 65 (2017), 126–139. https://doi.org/10.1016/j.compenvurbsys.2017.05.004
[3] Buru Chang, Yonggyu Park, Donghyeon Park, Seongsoon Kim, and Jaewoo Kang. 2018. Content-Aware Hierarchical Point-of-Interest Embedding Model for Successive POI Recommendation. In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence*, Vol. 2018-July. International Joint Conferences on Artificial Intelligence Organization, California, 3301–3307. https://doi.org/10.24963/ijcai.2018/458
[4] Joris Van den Bossche, Kelsey Jordahl, Martin Fleischmann, James McBride, Jacob Wasserman, Matt Richards, Adrian Garcia Badaracco, Alan D. Snow, Jeff Tratner, Jeffrey Gerard, Brendan Ward, Matthew Perry, Carson Farmer, Geir Arne Hjelle, Mike Taves, Ewout ter Hoeven, Micah Cochran, rraymondgh, Sean Gillies, Giacomo Caria, Lucas Culbertson, Matt Bartos, Nick Eubank, Ray Bell, sangarshanan, John Flavin, Sergio Rey, maxalbert, Aleksey Bilogur, and Christopher Ren. 2023. *geopandas/geopandas: v0.13.2*. GeoPandas. https://doi.org/10.5281/zenodo.8009629
[5] William Falcon and The PyTorch Lightning team. 2019. *PyTorch Lightning*. https://doi.org/10.5281/zenodo.3828935
[6] Christoph Feichtenhofer, Haoqi Fan, Bo Xiong, Ross Girshick, and Kaiming He. 2021. A Large-Scale Study on Unsupervised Spatiotemporal Representation Learning. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 3298–3308. https://doi.org/10.1109/CVPR46437.2021.00331
[7] Max Gabrielsson. 2023. *DuckDB Spatial Extension*. DuckDB Labs. https://github.com/duckdblabs/duckdb_spatial
[8] Rahul Ghosh, Xiaowei Jia, Leikun Yin, Chenxi Lin, Zhenong Jin, and Vipin Kumar. 2022. Clustering Augmented Self-Supervised Learning: An Application to Land Cover Mapping. In *Proceedings of the 30th International Conference on Advances in Geographic Information Systems* (Seattle, Washington) *(SIGSPATIAL '22)*. Association for Computing Machinery, New York, NY, USA, Article 3, 10 pages. https://doi.org/10.1145/3557915.3560937
[9] Sean Gillies, Casper van der Wel, Joris Van den Bossche, Mike W. Taves, Joshua Arnott, Brendan C. Ward, and others. 2023. *Shapely*. Shapely. https://doi.org/10.5281/zenodo.5597138
[10] Google Inc. 2023. S2 Geometry. https://s2geometry.io/. Accessed: 2023-06-23.
[11] Piotr Gramacki, Szymon Woźniak, and Piotr Szymański. 2021. Gtfs2vec: Learning GTFS Embeddings for Comparing Public Transport Offer in Microregions. In *Proceedings of the 1st ACM SIGSPATIAL International Workshop on Searching and Mining Large Collections of Geospatial Data* (Beijing, China) *(GeoSearch'21)*. Association for Computing Machinery, New York, NY, USA, 5–12. https://doi.org/10.1145/3486640.3491392
[12] Anita Graser. 2019. MovingPandas: Efficient Structures for Movement Data in Python. *GI_Forum* 1 (2019), 54–68. https://doi.org/10.1553/giscience2019_01_s54
[13] Aric Hagberg, Pieter Swart, and Daniel S Chult. 2008. *Exploring network structure, dynamics, and function using NetworkX*. Technical Report. Los Alamos National Lab.(LANL), Los Alamos, NM (United States).
[14] Charles R. Harris, K. Jarrod Millman, Stéfan J van der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J. Smith, Robert Kern, Matti Picus, Stephan Hoyer, Marten H. van Kerkwijk, Matthew Brett, Allan Haldane, Jaime Fernández del Río, Mark Wiebe, Pearu Peterson, Pierre Gérard-Marchant, Kevin Sheppard, Tyler Reddy, Warren Weckesser, Hameer Abbasi, Christoph Gohlke, and Travis E. Oliphant. 2020. Array programming with NumPy. *Nature* 585 (2020), 357–362. https://doi.org/10.1038/s41586-020-2649-2
[15] Sarah Hoffmann. 2023. pyosmium. https://github.com/osmcode/pyosmium.
[16] Hsun-Ping Hsieh, Su Wu, Ching-Chung Ko, Chris Shei, Zheng-Ting Yao, and Yu-Wen Chen. 2022. Forecasting Fine-Grained Air Quality for Locations without Monitoring Stations Based on a Hybrid Predictor with Spatial-Temporal Attention Based Network. *Applied Sciences* 12, 9 (2022). https://doi.org/10.3390/app12094268

[17] Kacper Leśniara and Piotr Szymański. 2022. Highway2vec: Representing OpenStreetMap Microregions with Respect to Their Road Network Characteristics. In *Proceedings of the 5th ACM SIGSPATIAL International Workshop on AI for Geographic Knowledge Discovery* (Seattle, Washington) *(GeoAI '22)*. Association for Computing Machinery, New York, NY, USA, 18–29. https://doi.org/10.1145/3557918.3565865
[18] Quentin Lhoest, Albert Villanova del Moral, Patrick von Platen, Thomas Wolf, Mario Šaško, Yacine Jernite, Abhishek Thakur, Lewis Tunstall, Suraj Patil, Mariama Drame, Julien Chaumond, Julien Plu, Joe Davison, Simon Brandeis, Victor Sanh, Teven Le Scao, Kevin Canwen Xu, Nicolas Patry, Steven Liu, Angelina McMillan-Major, Philipp Schmid, Sylvain Gugger, Nathan Raw, Sylvain Lesage, Anton Lozhkov, Matthew Carrigan, Théo Matussière, Leandro von Werra, Lysandre Debut, Stas Bekman, and Clément Delangue. 2021. Datasets: A Community Library for Natural Language Processing. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*. Association for Computational Linguistics, 175–184. https://aclanthology.org/2021.emnlp-demo.21
[19] Yijun Lin, Nikhit Mago, Yu Gao, Yaguang Li, Yao-Yi Chiang, Cyrus Shahabi, and José Luis Ambite. 2018. Exploiting Spatiotemporal Patterns for Accurate Air Quality Forecasting Using Deep Learning. In *Proceedings of the 26th ACM SIGSPATIAL International Conference on Advances in Geographic Information Systems* (Seattle, Washington) *(SIGSPATIAL '18)*. Association for Computing Machinery, New York, NY, USA, 359–368. https://doi.org/10.1145/3274895.3274907
[20] Tomas Mikolov, Kai Chen, Greg S. Corrado, and Jeffrey Dean. 2013. Efficient Estimation of Word Representations in Vector Space. http://arxiv.org/abs/1301.3781
[21] Masanao Ochi, Yuko Nakashio, Yuta Yamashita, Ichiro Sakata, Kimitake Asatani, Matthew Ruttley, and Junichiro Mori. 2016. Representation Learning for Geospatial Areas Using Large-Scale Mobility Data from Smart Card. In *Proceedings of the 2016 ACM International Joint Conference on Pervasive and Ubiquitous Computing: Adjunct* (Heidelberg, Germany) *(UbiComp '16)*. Association for Computing Machinery, New York, NY, USA, 1381–1389. https://doi.org/10.1145/2968219.2968416
[22] OpenStreetMap contributors. 2023. About OpenStreetMap . https://wiki.openstreetmap.org/wiki/About_OpenStreetMap. Accessed 2023-05-24.
[23] Luca Pappalardo, Filippo Simini, Gianni Barlacchi, and Roberto Pellungrini. 2019. *scikit-mobility*. https://doi.org/10.5281/zenodo.3273053
[24] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In *Advances in Neural Information Processing Systems 32*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché Buc, E. Fox, and R. Garnett (Eds.). Curran Associates, Inc., 8024–8035. http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf
[25] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. 2011. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research* 12 (2011), 2825–2830.
[26] Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep Contextualized Word Representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*. Association for Computational Linguistics, New Orleans, Louisiana, 2227–2237. https://doi.org/10.18653/v1/N18-1202
[27] Mark Raasveldt and Hannes Mühleisen. 2023. *DuckDB*. DuckDB Labs. https://github.com/duckdb/duckdb
[28] Kamil Raczycki and Piotr Szymański. 2021. Transfer Learning Approach to Bicycle-Sharing Systems' Station Location Planning Using OpenStreetMap Data. In *Proceedings of the 4th ACM SIGSPATIAL International Workshop on Advances in Resilient and Intelligent Cities* (Beijing, China) *(ARIC '21)*. Association for Computing Machinery, New York, NY, USA, 1–12. https://doi.org/10.1145/3486626.3493434
[29] Alex Raichev. 2019. GTFS Kit. https://github.com/mrcagney/gtfs_kit.
[30] Sergio Rey, Phil, Taylor Oshan, Charles Schmidt, jlaura, Levi John Wolf, Dani Arribas-Bel, David C. Folch, mhwang4, Nicholas Malizia, Wei Kang, Pedro Amaral, James Gaboardi, Luc Anselin, eli knaap, Qunshan, Stefanie Lumnitz, Andrew Winslow, Bas Couwenberg, Marynia, Omar Khursheed, Martin Fleischmann, KSai, yogabonito, Andrew Annex, Filipe, Stuart Lynn, Peter Quackenbush, Karl Dunkle Werner, and Caleb Robinson. 2023. *pysal/pysal: Release v23.01*. https://doi.org/10.5281/zenodo.7587827
[31] Fatih Sivrikaya and Ömer Küçük. 2022. Modeling forest fire risk based on GIS-based analytical hierarchy process and statistical analysis in Mediterranean region. *Ecological Informatics* 68 (2022), 101537. https://doi.org/10.1016/j.ecoinf.2021.101537

---

[20]https://github.com/geoarrow/geoarrow

[32] Vincent Spruyt. 2018. Loc2Vec: Learning location embeddings with triplet-loss networks. https://sentiance.com/loc2vec-learning-location-embeddings-w-triplet-loss-networks. Accessed: 2023-06-22.

[33] Adam J. Stewart, Caleb Robinson, Isaac A. Corley, Anthony Ortiz, Juan M. Lavista Ferres, and Arindam Banerjee. 2022. TorchGeo: Deep Learning With Geospatial Data. In *Proceedings of the 30th International Conference on Advances in Geographic Information Systems (SIGSPATIAL '22)*. Association for Computing Machinery, Seattle, Washington, 1–12. https://doi.org/10.1145/3557915.3560953

[34] The pandas development team. 2023. *pandas-dev/pandas: Pandas*. https://github.com/pandas-dev/pandas

[35] Uber Technologies Inc. 2023. H3: Uber's Hexagonal Hierarchical Spatial Index. https://eng.uber.com/h3/. Accessed: 2023-06-23.

[36] Leonardo Uieda. 2018. Verde: Processing and gridding spatial data using Green's functions. *Journal of Open Source Software* 3, 30 (2018), 957. https://doi.org/10.21105/joss.00957

[37] Joris Van Den Bossche, Martin Fleischmann, Thomas Statham, Daniel Jahn (Dahn), Tom Augspurger, Julia Signell, Pete Gadomski, Ray Bell, Stefanie Lumnitz, Ali Abbas Zaidi, Fred Bunt, Ian Rose, Irina Truong, James Bourbeau, Jason Baker, Matt Morris, Raphael Hagen, RichardScottOZ, and Bernardpazio. 2023. *geopandas/dask-geopandas: v0.3.1*. https://doi.org/10.5281/ZENODO.7875807

[38] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Perric Cistac, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. Transformers: State-of-the-Art Natural Language Processing. Association for Computational Linguistics, 38–45. https://www.aclweb.org/anthology/2020.emnlp-demos.6

[39] Szymon Woźniak and Piotr Szymański. 2021. Hex2vec: Context-Aware Embedding H3 Hexagons with OpenStreetMap Tags. In *Proceedings of the 4th ACM SIGSPATIAL International Workshop on AI for Geographic Knowledge Discovery* (Beijing, China) *(GEOAI '21)*. Association for Computing Machinery, New York, NY, USA, 61–71. https://doi.org/10.1145/3486635.3491076

[40] Ning Wu, Xin Wayne Zhao, Jingyuan Wang, and Dayan Pan. 2020. Learning Effective Road Network Representation with Hierarchical Graph Neural Networks. In *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. Association for Computing Machinery, 6–14. https://doi.org/10.1145/3394486.3403043

[41] Qiusheng Wu and Lucas Prado Osco. 2023. *samgeo: A Python package for segmenting geospatial data with the Segment Anything Model (SAM)*. https://doi.org/10.5281/zenodo.7966658

[42] Zijun Yao, Yanjie Fu, Bin Liu, Wangsu Hu, and Hui Xiong. 2018. Representing Urban Functions through Zone Embedding with Human Mobility Patterns. In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence*, Vol. 2018-July. International Joint Conferences on Artificial Intelligence Organization, California, 3919–3925. https://doi.org/10.24963/ijcai.2018/545

[43] Manzhu Yu, Chaowei Yang, and Yun Li. 2018. Big Data in Natural Disaster Management: A Review. *Geosciences* 8, 5 (2018). https://doi.org/10.3390/geosciences8050165

[44] Sen Zhang, Shaobo Li, Xiang Li, and Yong Yao. 2020. Representation of Traffic Congestion Data for Urban Road Traffic Networks Based on Pooling Operations. *Algorithms* 13, 4 (2020). https://doi.org/10.3390/a13040084

[45] Zheng Zhang and Liang Zhao. 2021. Representation Learning on Spatial Networks. In *Advances in Neural Information Processing Systems*, M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (Eds.), Vol. 34. Curran Associates, Inc., 2303–2318. https://proceedings.neurips.cc/paper_files/paper/2021/file/12e35d9186dd72fe62fd039385890b9c-Paper.pdf