



ÉCOLE NATIONALE
DES SCIENCES
GÉOGRAPHIQUES

France National School of
Geographic Sciences



Finnish Geospatial
Research Institute,
National Land Survey of
Finland



Funding bodies

Internship report

ENSG Engineering Cycle — 2nd year

Massive data contributions from the National Land Survey of Finland to the OpenStreetMap Database



Alexys Ren

May - August 2023

Non confidentiel Confidential IGN Confidential Industrie Jusqu'au ...

ÉCOLE NATIONALE DES SCIENCES GÉOGRAPHIQUES
6-8 Avenue Blaise Pascal - Cité Descartes - 77420 Champs-sur-Marne
Téléphone 01 64 15 31 00 Télécopie 01 64 15 31 07

Jury

Sponsor:

Prof Juha Oksanen, director, Geoinformatics and cartography, NLS/FGI

Internship supervision:

Juha Oksanen (and other department's senior scientists)

Referring teacher:

Victor Coindet, IGN/ENSG/TSI

Head of engineering studies:

Jean-François Hangouët, IGN/ENSG/DIAS

Internship management:

Isabelle Joyeux, IGN/ENSG/DSHEI

© ENSG

Second year Multidisciplinary Internship from 22/05/2023 to 11/08/2023

Webcasting: Internet Intranet Polytechnicum ENSG Intranet

Document status:

End-of-year internship report presented at the end of the 2nd year of the engineering cycle

Number of pages: 92 pages including 50 pages of appendices

Host system: L^AT_EX

Modifications:

EDIT	OVERHAUL	DATE	MODIFIED PAGES
1	0	06/2023	Created

Acknowledgements

First of all, I would like to thank Juha Oksanen, my internship supervisor. He, together with Valtteri Nikkinen, Project Coordinator, gave me a very warm welcome and made sure that I had everything I needed to feel comfortable. Thanks to them, I was able to quickly adapt to this new environment and try to make the best of it.

I'm very grateful to Ville Mäkinen, Research Group Manager, and Jaakko Kähkönen, Senior Research Scientist, who, along with Juha, were very helpful in untangling technical problems during my research and were themselves available whenever I needed support.

I would also like to thank Ms. Ana-Maria Olteanu-Raimond, Senior Researcher and Co-Director of the French Lastig Laboratory, who introduced me to Juha and provided me with the basic material I worked on, without which this internship wouldn't be possible.

I'm grateful to Victor Coindet, my referring teacher at ENSG, for his helpful feedback and suggestions that helped me make the best of my work.

I'd also like to thank the Fondation ENSG-Géomatique, Coexya SAS and the European Union for their financial support of this project, which allowed me to carry out my research under the best conditions.

And finally, I would like to reiterate my thanks to all of FGI. I've spent such a great moment in this atmosphere, every person I've met has been wonderful. Thank you — Panu, Pyry, Valtteri, Anssi, Ville, Camilo, Milla, Justus, Hilla, Alpo, Anna, Rebecca, Jaakko, Christian, Juha.

Résumé

(*French Abstract, English version next page*)

Le FGI, ou l’Institut finlandais de recherche géospatiale, est situé dans la ville d’Espoo, à l’ouest de la capitale Helsinki. Il mène des recherches novatrices et des travaux d’experts dans le domaine des données spatiales et, plus précisément, traite des services web interopérables, des normes techniques et de l’harmonisation des données spatiales, qui sont les principaux domaines du département de géoinformatique et de cartographie dans lequel je suis formé.

Dans le prolongement de l’harmonisation des données spatiales, qui vise à standardiser l’information géographique au sein de l’Union européenne, l’objectif de mon travail est d’identifier les contributions massives du NLS (organisme supérieur dont relève le FGI) à la base de données OpenStreetMap afin d’obtenir un état des lieux des données NLS sur Internet, ainsi qu’une évaluation préliminaire des risques, ceci dans l’espoir de contribuer à de futures études sur le potentiel des entités importées ayant subi des modifications géométriques pour améliorer la base de données topographiques du NLS.

Notre étude commence par l’obtention des données d’intérêt auprès de fournisseurs de données bien connus tels que Geofabrik.de et Planet OSM, avant de présenter quelques techniques de manipulation des données OSM à l’aide de Osmium Tool, le pilote OSM de la bibliothèque GDAL, Python et SQL (PostgreSQL), en préparation de la phase d’analyse des données.

Les premiers résultats permettent de supposer l’emplacement des bâtiments liés au NLS ayant des contributions massives dans la région finlandaise d’Uusimaa, ainsi que l’emplacement de ceux qui ont subi des modifications géométriques.

Mots clés : OpenStreetMap, importation massive d’OSM, analyse des données d’OSM, clé source d’OSM, import du fichier historique d’OSM, ouverture des données du NLS Finlande au grand public

Abstract

The Finnish Geospatial Research Institute (FGI) is located in the city of Espoo, west of the capital city of Helsinki. It conducts innovative research and expert work in the field of spatial data, and more specifically, deals with — interoperable web services, technical standards and harmonization of spatial data — which are the core areas of the Geoinformatics and Cartography department in which I train.

Following on from the harmonization of spatial data, which aims to standardize geographic information across the European Union, the goal of my work is to identify the massive contributions from the NLS (higher body under which the FGI is submitted) to the OpenStreetMap database to get a situation picture of the NLS data on the Internet, as well as part of a preliminary risk assessment, this in the hope of contributing to future studies on the potential of imported features that have undergone geometric modifications to improve the NLS topographic database.

Our study begins by obtaining the data of interest from well-known data providers such as Geofabrik.de and Planet OSM, before presenting some OSM data manipulation techniques using — Osmium Tool, GDAL OSM Driver, Python and SQL (PostgreSQL) — in preparation for the data analysis phase.

First results help to assume the location of NLS-related buildings with massive contributions in the Finnish Uusimaa region, as well as the location of those that have undergone geometric modifications.

Key words: OpenStreetMap, OSM massive import, OSM data analysis, OSM source tag, OSM history file import, NLS of Finland open data release

Contents

Glossary and useful acronyms	5
Introduction	7
1 Study site and data sets selection	9
1.1 Study site	9
1.2 OSM data structure	11
1.3 Identifying NLS data among the whole OSM data	11
1.4 Some related work	11
2 Data retrieval from the OSM database	13
2.1 QuickOSM QGIS plugin	13
2.2 Geofabrik OSM data extracts provider	14
2.3 Planet OSM	15
3 Data manipulations	17
3.1 Sub sets extraction from raw downloaded data files	17
3.2 Source tags values retrieval	19
3.3 Feature-level modifications collecting	23
4 Data analysis and Visualisations	25
4.1 Identifying the massive contributions of NLS data in OSM	25
4.2 Feature-level modifications underwent by NLS data since its import into OSM	32
Conclusion	35
A Source tags fetching with Python	45
B GDAL Configuration file for OSM import	49
C First UML representation for uusimaa database	53
D Python script for importing changesets into uusimaa database	55
E NLS-related data proportions with additional 'created_by' and 'other_tags' tags	59
F SQL queries for categorised feature table creation	61
G SQL query to select the most prominent occurrences of NLS-like source values	63
H Python script for OSM history file import into PostgreSQL	65
I Naive Python script for OSM history file import into PostgreSQL	73

4 CONTENTS

J OSM Feature-level modifications analyser Python script	79
K Python script for retrieving OSM ids from feature which underwent geometric modifications	91

Glossary and useful acronyms

Apothem Distance from the center of a regular polygon to one of its sides

Capital Region Consists of four municipalities — Helsinki, Vantaa, Espoo and Kauniainen

Contributor Refers to any user who contributes to the OpenStreetMap community by importing its data (personal data or from external sources)

ENSG *École Nationale des Sciences Géographiques* — FRENCH. Standing for France National School of Geographic Sciences

Feature Synonym for an OSM element

FGI *Finnish Geospatial Research Institute*. Conducts innovative research and expert work within the field of spatial data. Under direct supervision of the NLS and place of the internship (Geoinformatics and Cartography Department)

Helsinki Metropolitan Area Not firmly established term, varying in different contexts. Commonly refers to the Greater Helsinki, the metropolitan area surrounding the capital city of Helsinki, and including the smaller Capital Region

IGN *Institut national de l'information géographique et forestière* — FRENCH. Public administrative institution under the supervision of the French Ministries of Ecology and Forestry. Its mission is to produce and distribute reference data (open data) and representations (online and paper maps, geovisualisation) relating to knowledge of the national territory and French forests, as well as their evolution

NLS *National Land Survey of Finland*. Official body dealing with cartographic and cadastral matters in Finland under the direction of the Finnish Ministry of Agriculture and Forestry. Equivalent to the French IGN

OSM *OpenStreetMap*. Collaborative online mapping project that aims to build a free geographic database of the world, using GPS and other open data. Launched in July 2004 by Steve Coast at University College London

Introduction

In 2012, the National Land Survey (NLS) of Finland opened its topographic data to everyone, thus allowing a large amount of data downloading from all over the world. Yet, with the ongoing war in Ukraine, there has been a change in the way of thinking regarding the transparency of map data. Already last year, there were a lot of requests from Russia for municipal map services, in line with its strategy of using open source data in the offensive. This prompted some politicians to call for more restricted spatial data accesses, especially those about the locations of critical infrastructures [17].

Even though there is no sudden threat that could be directed at Finland, it is justified to take a closer look at what information is made available and how it is used. Our case study will focus on OpenStreetMap, one of the biggest open spatial database in the world, where interesting information is most likely to appear. As part of a preliminary risk assessment for NLS data and with this in mind, we will try to answer the following question: **What massive contributions from the NLS has been done to OSM in history?**

More specifically, the goal of my work is *to discover how the NLS's data has been used in contributing to OSM database since the NLS's open data release*. We will investigate what was the data sent to OpenStreetMap from the NLS's perspective and what modifications those data underwent.

In addition to being an important investigation of getting a situation picture of the NLS's data in OSM, this work is also an opportunity for Prof. Oksanen to strengthen the collaboration between the French IGN and the Finnish NLS FGI across exchanges with Ms. Ana-Maria Olteanu-Raimond (Senior Researcher and Co-Director of the French Lastig Laboratory), with who they both planned to collaborate on summer 2022. On the other hand, NLS will be glad to have a better understanding on this matter, as this will allow the organization *to get understanding on OSM's potential to improve NLS's data with its latest updates*.

To this extent, my whole work structure will be based on Le Guilcher's, Olteanu-Raimond's and Bailo Balde's paper published in 2022 [19]. *The idea is to try to follow as carefully as possible the method her team discovered by applying it to Finland*, in order to enable the possibility for making comparison with the results obtained in France. We'll look on introducing the study area and the data choices first, before retrieving data from OSM. Then, we'll proceed with cleaning up the data and the data handling techniques, and eventually with both analyzing the data and getting visualizations of NLS features.

STUDY SITE AND DATA SETS SELECTION

1.1 Study site

On Prof. Oksanen's suggestion, the study area will be limited to the *Uusimaa region and its surroundings*. The extent of the area is about 10,000 km² wide and includes the Helsinki Metropolitan area, which should be sufficient for preliminary tests.

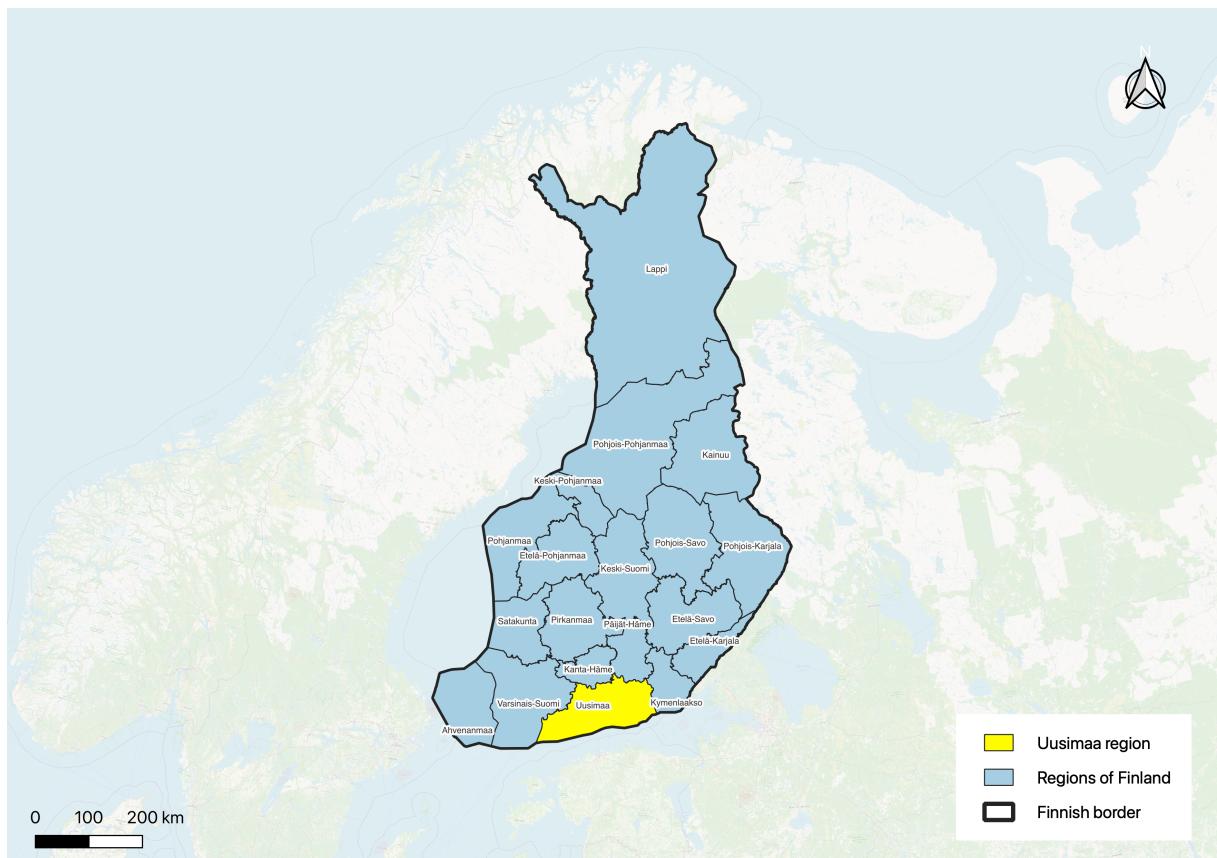
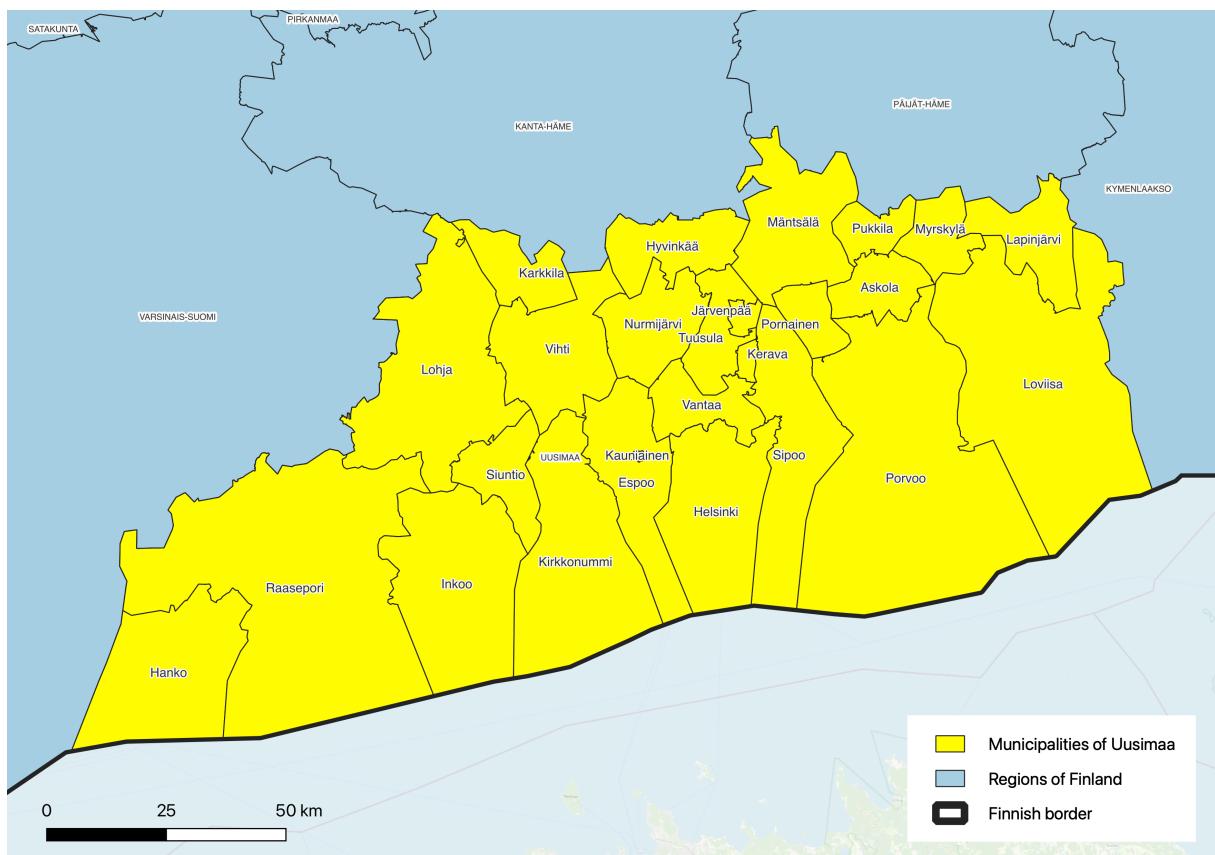


Figure 1.1: Finland's Uusimaa region (Map base credits: © OpenStreetMap contributors)



(b) Focus on the Uusimaa region

Figure 1.1: Finland's Uusimaa region (cont.)

For the sake of brevity, the study site will be referred to as *Uusimaa* for the rest of our reasoning.

1.2 OSM data structure

The OpenStreetMap data is made of core components which are the *elements* [5]. These features correspond to specific geometries which are among one of those three following types:

- node: points defined by their coordinates
- way: polylines defined as sequences of 'nodes'
- relation: features grouping one or more 'node', 'way' or 'relation'

Each element is also characterized by a list of tags (pairs of key/value) giving more information about it. Some elements may have mandatory tags (for instance the feature ID).

1.3 Identifying NLS data among the whole OSM data

Searching for NLS-related data is done via the *source* tag [9], as described in the first step of the methodology in [19]. It is used to indicate the source of some piece of information in OSM, and can help contributors who wish to understand the origins of the data they are editing.

In order to explore the most important types of features in one go, we are going to focus only on the *building* data [2] (represented by nodes, ways and relations) in the one hand, and on the *road network* data [6] (represented by ways and relations) in the other hand. The study of Points of Interest (POI) is deprecated, since it is not interesting for the NLS.

From now on, the building data and the road network data are referenced as *buildings* and *roads* respectively.

1.4 Some related work

As of now, few massive contributions on our selected themes are documented on the OSM wiki [3]. Each import has also its own tagging rules, without any standardized labelling. Then, our work is being meant to allow identification of the massive contributions which were not properly documented.

In the same spirit as said in the Introduction, some Italian researchers shed light on the pros and cons of integrating OSM data in the process of producing governmental data sets [26], touching the discussion where in the future the FGI need to think if OSM could be an important information source for the NLS topographic database.

DATA RETRIEVAL FROM THE OSM DATABASE

2.1 QuickOSM QGIS plugin

One very quick way to download the OSM data we are interested in here is to use the QuickOSM QGIS plugin. Based on the Overpass API [1], it features the "Quick query" tool, allowing the user to retrieve data on keys/values filtering (tags).

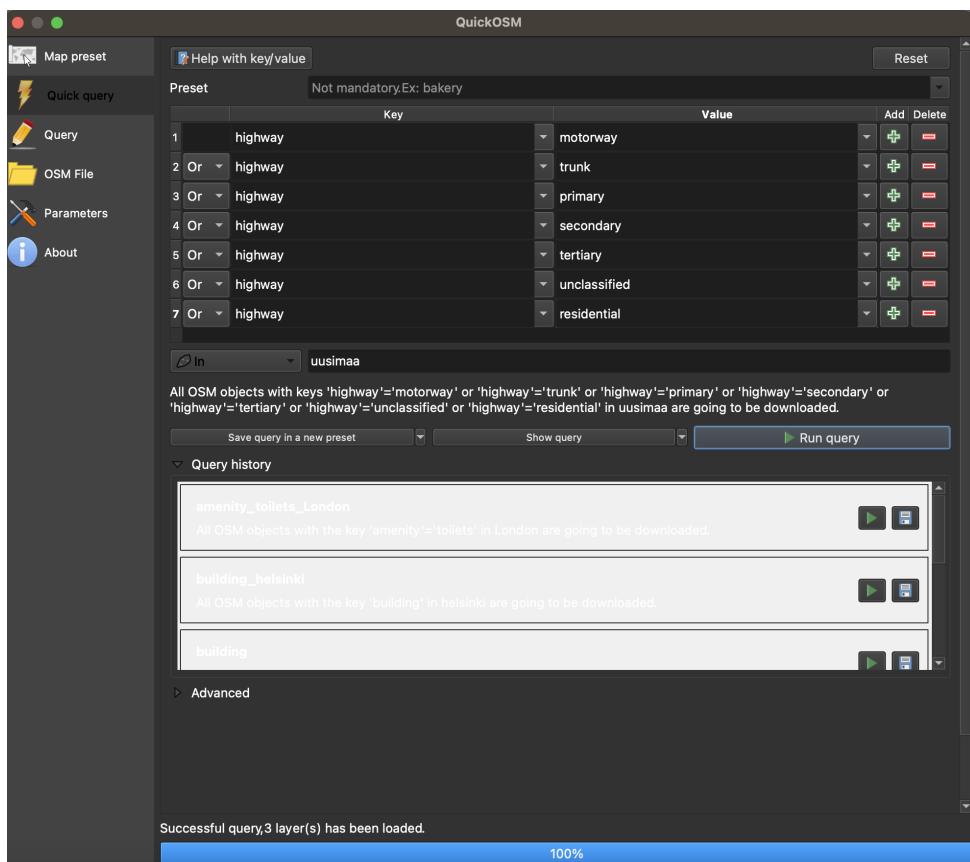


Figure 2.1: QuickOSM usage example in QGIS — *roads* from Uusimaa retrieval

As shown in Figure 2.1 above, several matching conditions are provided to get the desired output. In fact, simply putting the relevant key isn't enough. For example, *roads* features in OSM are identified with the *highway* key, but if no value is entered, the query doesn't return anything.

Then, this constraint force us to know the exact values that describe our *roads*, which is pretty irrelevant in our case: documentations on the OSM wiki [6] [2] don't specify clear boundaries for the *roads* tagging, same applying for the *buildings*'s. In the other hand, those queries don't give the upper hand on the data structure in the output, i.e. there is no control on the features tags and the

types of geometries. Some testing highlighted the missing of the source tag, which is indeed vital for our investigation.

Proceeding with this method does not guarantee to get all the wanted data, since there will always be some missing!

Eventually, the running time for loading the resulting attribute table in QGIS is quite long for the size of the actual data (only 69,207 features filtered by the query in Figure 2.1).

2.2 Geofabrik OSM data extracts provider

Following the previous discoveries, an alternative way to get our data is to download it directly from source. Geofabrik.de is a third party providing geographical extracts from the complete OSM database with up-to-date worldwide coverage.

Two downloading servers are available:

- *Public server*: contains regional data stripped from OSM contributors personal metadata (user names, user IDs, changeset IDs), as they are subject to data protection regulations in the European Union
- *Non-Public server*: reserved for OSM contributors and internal purposes, contains extracts with full metadata

With the permission of the Internship supervisor, **we will continue our study with data from the Non-Public server**, since *changeset IDs* are also required to identify NLS features (see section 3.2).

Download OpenStreetMap data for this region:

Finland

(one level up)

The OpenStreetMap data files provided on this server contain personal data of the OpenStreetMap contributors. Therefore, their usage is governed by data protection regulations in the European Union. These regulations apply even to data processing that happens outside the European Union because some people whose data is contained in these files live in the European Union.
The personal information contained in these files must only be used for OpenStreetMap internal purposes, e.g. quality assurance. You must ensure that derived databases and works based on these files are only accessible to OpenStreetMap contributors if they contain personal information.
Use our [public download server](#) to get the files without metadata.

Commonly Used Formats

- [finland-latest-internal.osm.bz2](#), suitable for Osmium, Osmosis, imposm, osm2pgsql, mkgmap, and others. This file was last modified 9 hours ago and contains all OSM data up to 2023-05-31T20:21:24Z. File size: 600 MB; MD5 sum: [5a08b8a3b9761f0b59473561122c6a4a](#).
- [finland-latest-free-shp.zip](#) Shape files are only available [without personal metadata](#).

Other Formats and Auxiliary Files

- [finland-latest.osm.bz2](#) is only available [without personal metadata](#).
- [finland-internal.osm.oif](#), a file that contains the full OSM history for this region for processing with e.g. osmium. This file was last modified 3 days ago. File size: 1.2 GB; MD5 sum: [ca807695c4592537e643bf9996d7699](#).
- [poly_file](#) that describes the extent of this region.
- [osm.gz](#) files that contain all changes in this region, suitable e.g. for Osmosis updates
- [raw_directory_index](#) allowing you to see and download older files

Sub Regions

No sub regions are defined for this region.

*) Shape files and bz2 compressed OSM XML files are only available [without personal metadata](#).

GEOFABRIK® downloads

Not what you were looking for? Geofabrik is a consulting and software development firm based in Karlsruhe, Germany specializing in OpenStreetMap services. We're happy to help you with data preparation, processing, server setup and the like. [Check out our website](#) and contact us if we can be of service.

Nicht das Richtige dabei? Die Geofabrik ist ein auf OpenStreetMap spezialisiertes Beratungs- und Softwareentwicklungsunternehmen in Karlsruhe. Gern helfen wir Ihnen bei der Datenaufbereitung, Datenkonvertierung, Serverinstallation und ähnlichen Aufgaben. [Besuchen Sie unsere Website](#) und sprechen Sie mit uns, wenn wir Ihnen helfen können.

Data/Maps Copyright 2018 Geofabrik GmbH and OpenStreetMap Contributors | Map tiles: Creative Commons BY-SA 2.0 Data: ODbL 1.0 | [Contact Us](#)

Figure 2.2: Finland's Non-Public data Geofabrik downloading webpage — URL: <https://osm-internal.download.geofabrik.de/europe/finland.html>

From this page, **finland-latest-internal.osm.pbf** and **finland-internal.osm.pbf** files were downloaded on the 1st of June 2023 at 2pm. The former is a *data file* containing OSM

data from a specific point in time (so as of the 1st of June here) *with at most one version per object (node, way, or relation)*. The second is a *history file* containing not only the current version of one object, but also all its history, i.e. there are zero or more versions for any object in this file.

Note that .pbf extension means the relevant file is in binary format. It can be read by most programs dealing with OSM data but not by a simple text editor. As for .osm or .osh extensions, these correspond to variants of the classical XML format. For more information on Geofabrik's data extracts, please see [14].

As you've maybe noticed in Figure 2.2, no subregions are available for Finland, i.e. data related to Uusimaa region is contained in our bigger file. Extracting sub sets from this downloaded data is the matter in the section 3.1.

2.3 Planet OSM

Planet OSM is a benchmark source allowing to download the whole OpenStreetMap database as a regularly-updated single file. Besides traditional .pbf and XML files, the website also offers a *changesets* file [4] containing complete metadata, for all the changes that ever happened in OpenStreetMap's history.

Since this file wasn't available on the Geofabrik's servers, it was downloaded from Planet OSM (**on 05/06/2023 at 5pm**) and we'll be of use later, especially for the collection of features modifications in the section 3.3.

Planet OSM

The files found here are regularly-updated, complete copies of the OpenStreetMap.org database, and those published before the 12 September 2012 are distributed under a Creative Commons Attribution-ShareAlike 2.0 license, those published after are Open Data Commons Open Database License 1.0 licensed. For more information, [see the project wiki](#).

WARNING Download speeds are currently restricted to 4096 KB/s due to limited available capacity on our Internet connection. [Please use torrents](#) or [a mirror](#) if possible.

Complete OSM Data	Using The Data	Extracts & Mirrors
Latest Weekly Planet XML File (torrent) (RSS) 127 GB , created 3 days ago. md5: 8f0689e4e4fb233d9bd65afe89ee0500.	You are granted permission to use OpenStreetMap data by the OpenStreetMap License , which also describes your obligations.	The complete planet is very large, so you may prefer to use one of several periodic extracts (individual countries or states) from third parties. GeoFabrik.de and BBBike.org are two providers of extracts with up-to-date worldwide coverage.
Latest Weekly Changesets (torrent) (RSS) 5.7 GB , created 3 days ago. md5: 2cd8e769a2518739e7ab9db773c58e79.	You can process the file or extracts with a variety of tools. Osmosis is a general-purpose command-line tool for converting the data among different formats and databases, and Osm2pgsql is a tool for importing the data into a Postgis database for rendering maps.	
Latest Weekly Planet PBF File (torrent) (RSS) 69 GB , created 3 days ago. md5: 9d37d1387b080970f6d996f893336cae.	Processed coastline data derived from OSM data is also needed for rendering usable maps.	
Each week, a new and complete copy of all data in OpenStreetMap is made available as both a compressed XML file and a custom PBF format file. Also available is the history file which contains not only up-to-date data		

Figure 2.3: Planet OSM website homepage — URL: <https://planet.openstreetmap.org/>

DATA MANIPULATIONS

3.1 Sub sets extraction from raw downloaded data files

As announced in chapter 1, our study will be limited exclusively to *buildings* and *roads* from Uusimaa region. In order to get those sub sets data from downloaded data, *Osmium Tool* we'll be used.

3.1.1 Osmium Tool

Osmium Tool [24] is a command-line tool based on the C++ library libosmium [22] and designed for processing OSM data. It features plenty of useful functions [23], as well as fast-reading properties, and often outperform similar tools.

Considering the unique shape of OSM data, the decision was made to carry on with this tool for the data extracting phase instead of QGIS, as it allows a greater flexibility, even so the learning curve is pretty steep in the early ages.

3.1.2 Data & History files

The data slicing was done in 2 stages: restricting the area covered by the data to the Uusimaa region, then extracting the themes of interest, i.e. *buildings* and *roads*.

Getting data from Uusimaa: we use the `osmium extract` command for creating geographical extracts from an OSM file. The geographical extent for the clip was given as a *bounding box* around Uusimaa, which includes a little more of neighboring regions to avoid edge effects.

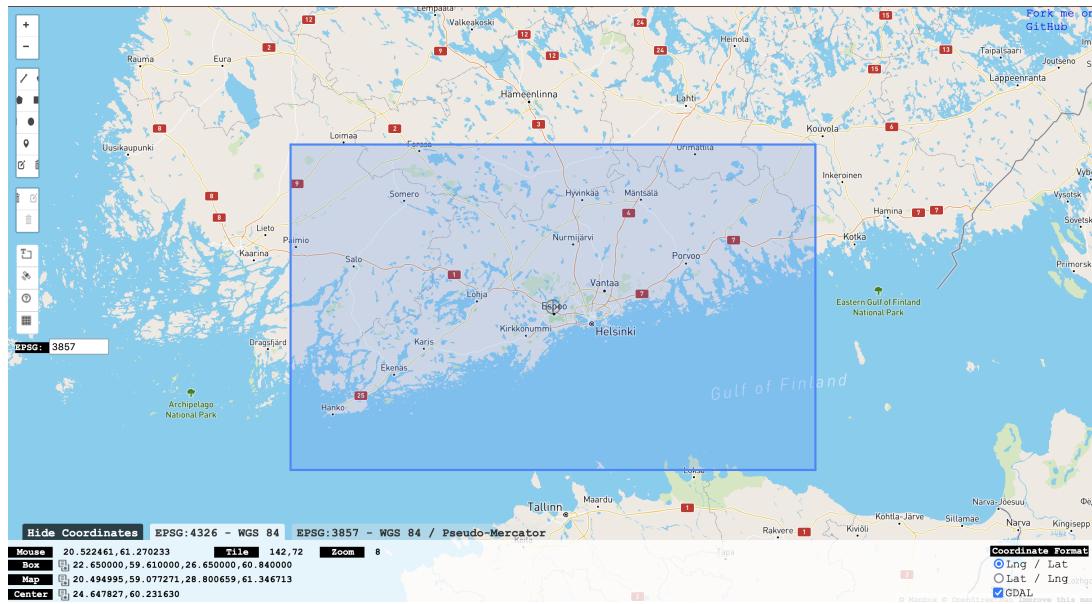


Figure 3.1: Bounding box used for Uusimaa's data clipping, *l*ng/*lat* coordinates (22.65,59.61,26.65,60.84) — Uusimaa boundaries in *light gray* inside the blue rectangle (credits: [bboxfinder.com](#))

Here's what the final command looks like:

```
1 osmium extract -b 22.65,59.61,26.65,60.84 finland-latest-internal  
    .osm.pbf -o uusimaa.osm.pbf
```

The result from this code launch is `uusimaa.osm.pbf`, a new binary file containing all the features in the Uusimaa region.

We proceed in much the same way for the history file, except that we add the `-H` option to make the command compatible with this particular type of file: in fact, most programs using OSM data expect features IDs to be unique, so they can't work with history data! As a result, `uusimaa-history.osh.pbf` is also created.

Getting buildings and roads from Uusimaa (**only for the data file**): we use the osmium tags-filter command for filtering objects matching specified keys/tags with the previous `uusimaa.osm.pbf` file as the input. Here, unlike in section 2.1, the filtering doesn't require specific values to be matched with the keys.

The process is then greatly simplified. Here's the command used for the *roads* extraction:

```
1     osmium tags-filter uusimaa.osm.pbf wr/highway -o roads-uusimaa.osm.pbf
```

Just above, we looked for all ways and relations (wr/) that have been tagged with the highway key, which identifies the *roads* features.

The same operation is realized to get the *buildings* data:

```
1     osmium tags-filter uusimaa.osm.pbf building -o buildings-uusimaa.osm.pbf
```

Here, geometry types are not specified, since we are interested in all of them (nodes, ways and relations).

At the end of this step, `roads-uusimaa.osm.pbf` & `buildings-uusimaa.osm.pbf` are created, corresponding to Uusimaa road network features and Uusimaa buildings features respectively.

3.1.3 Changesets file

Likewise, a bounding box is again used here to cut Uusimaa's changesets from whole planet's, but with the specific changeset command `osmium changeset-filter`.

```
1      osmium changeset-filter -B 22.65,59.61,26.65,60.84 -o uusimaa
          -changesets.osm.bz2 changesets-230529.osm.bz2
```

The output is `uusimaa-changesets.osm.bz2`, a bzip2-compressed XML file.

3.2 Source tags values retrieval

3.2.1 Strategy

As introduced in section 1.3, the key for browsing NLS features is based on the source tag. Historically, it was the usual practice for adding this source information as an attribute on each relevant feature. However with modern editors, the source tag is more and more put on the changeset where several features edited in the same session are regrouped [9].

The aim of the following paragraphs is then to get the features' sources by looking at the *feature's attributes level*, but also at the *changesets level* in which the features are included.

3.2.2 Python-based method

The first idea was to take a look at our extracted data in QGIS (`roads-uusimaa.osm.pbf` & `buildings-uusimaa.osm.pbf`). Despite QGIS's ability to handle .pbf files by adding it as a Vector Layer, the files' sizes were way too big for the program to load correctly. Then, one solution was to convert our files into *SpatiaLite* databases using *GDAL ogr2ogr command-line tools*, and connect those to QGIS [1]:

```
1      ogr2ogr -f SQLite -dsco SPATIALITE=YES buildings-uusimaa.db
          buildings-uusimaa.osm.pbf
```

Here is an example for *buildings*, resulting in a .db file, readable in QGIS through a SpatiaLite connection.

Unfortunately, an attribute analysis in QGIS of the features in our .db files again revealed the absence of the source tag in our data, thus forcing our investigation to focus on changesets.

Then, the imagined processing chain is based on the *changeset id*, a unique identifier for each changeset entity. Also considering that every feature in OSM is associated to a changeset, this link allows source information on changeset to be assigned to each feature. To do so, we use a homemade Python script which generates a .csv file on an initial data file input, containing the following information (fields):

- type: type of feature (node, way or relation)
- id: OSM unique identifier for the feature
- changesetId: OSM unique identifier for the associated changeset
- source: source information about the changeset

The script is based on OSMCha [10] (short for OpenStreetMap Changeset Analyzer) Python library used for advanced changesets analysis, and PyOsmium [18] which is an adaptation of the previous libosmium C++ library for Python support, used here for OSM file reading. For more details, please check the relevant code and the preliminary results overview both available in appendix A.

Note that the .csv output format is justified as this will allow an easier import into a PostgreSQL database for future string-targetted queries on the source spelling (i.e. references to NLS).

However, the proper functioning of this code depends on the input file's size. In fact, the source information retrieval corresponds to single requests to the OSM server, meaning the number of requests is equal to the number of the input features. This way of proceeding is not acceptable for our actual files considering its sizes, as it will also lead to Internet network saturation with an endless code launching.

This method is eventually deprecated.

3.2.3 GDAL Configuration file for OSM import

Thanks to a hint from the internship supervisor, **the feature-level source tag has finally been explicated by modifying the osmconf.ini file**. It is a configuration text file delivered with the GDAL installation packages, allowing to add the exact attributes wanted to be appearing in the output file generated from a ogr2ogr command (configuration file content available in appendix B).

In our case, 2 additional tags are considered:

- source: the feature-level source information
- created_by [8]: information on the software used for the feature edition

In addition to potential NLS-tagged features, the last tag allows broader investigations regarding features which might have a link to the NLS in case they are wrongly tagged for instance. In the same spirit, the converted_by tag could have been interesting, but has no known use according to the OSM wiki [7].

Also, in anticipation of the very next step, the field osm_changeset has been modified to yes value in the configuration file, enabling the changeset ids display.

Eventually, we use ogr2ogr commands to convert our data file from one format to another, in our case from .pbf to .gpkg (Geopackage file) for QGIS reading. Here is an example for the *buildings* data from Uusimaa:

```
1  ogr2ogr -f GPKG buildings-uusimaa.osm.gpkg buildings-uusimaa.osm.pbf
```

Attribute table verifications in QGIS have shown the correct addition of the desired attributes, which will be the subject of more advanced analysis in the next chapter.

3.2.4 PostgreSQL method

Following the previous discoveries, another way for getting changeset-level source tag is to look at our *Uusimaa's changesets file* directly. The idea is to import this file into a PostgreSQL database, and then to make SQL JOIN operations with imported feature tables on the changeset id. The link is pretty much the same as in subsection 3.2.2, but now the data fetching will be done *locally* instead.

ogc_fid	osm_id	osm_way_id	osm_version	osm_timestamp	osm_uid	osm_user	osm_changeset	name
1	1	4126	[null]	6	2022-06-19 14:02:54+03	8105430	Larmax	122380934 Trololi
2	2	4127	[null]	5	2021-03-25 21:26:43+02	145231	woodpeck_repair	101746091 Tilastoke
3	3	4138	[null]	5	2022-02-02 09:21:18+02	125897	overfloran	116903427 Ukraina
4	4	4141	[null]	3	2020-04-17 11:30:01+03	3343548	lkm-22	83704078 kasivuhi
5	5	4148	[null]	11	2023-05-15 10:53:43+03	8848422	Emil Bürger	136120806 Aalto-ylic
6	6	4150	[null]	3	2022-04-17 17:08:17+03	5654160	Riku V	119824851 [null]
7	7	4151	[null]	3	2022-06-19 09:34:36+03	12009838	Oliaq	122572372 [null]
8	8	4198	[null]	5	2020-05-20 06:15:09+03	2507399	aperezdc	85477809 [null]
9	9	4200	[null]	3	2017-03-15 22:15:46+02	691099	jleh	46882650 [null]
10	10	4202	[null]	5	2017-04-27 10:34:47+03	665748	sebastian	48190077 [null]
11	11	4701	[null]	5	2014-12-31 12:49:44+02	139957	jl-	27822009 [null]
12	12	5603	[null]	6	2018-05-03 12:04:34+03	345975	ViliYur	58645851 [null]
13	13	5605	[null]	6	2019-03-16 09:06:08+02	2507399	aperezdc	68251524 [null]
14	14	5606	[null]	5	2021-11-21 18:54:40+02	968323	susannaanas	114069349 [null]
15	15	5608	[null]	10	2022-07-03 11:53:33+03	8980540	pyyhtu	123148358 Marski_bj
16	16	6062	[null]	5	2023-01-04 15:35:23+02	11941474	JiriG	130872112 [null]
17	17	6065	[null]	5	2021-01-08 19:02:59+02	968323	susannaanas	97186493 [null]
18	18	6066	[null]	8	2022-11-29 10:23:21+02	4517549	HSL_HRT	129514173 Kansallinen
19	19	6077	[null]	3	2021-01-13 13:45:08+02	968323	susannaanas	97378072 [null]
20	20	6518	[null]	6	2021-09-09 11:49:05+03	365087	Tuukkash	110971957 Keskuspi
21	21	6519	[null]	3	2017-04-27 10:34:47+03	665748	sebastian	48190077 [null]

Figure 3.2: *pgAdmin 4* interface showing contents from the *uusimaa* database

The *uusimaa* database was created under PostgreSQL 15, with *postgis* and *hstore* extensions. The former allows storing features' geometry information, while the latter is a special data type for storing different information of undefined length as sets of key/value pairs within a single column [15]. It will host the feature data tables as well as the changesets table.

For importing the feature data, there are several possibilities, and the first one is to use *Osm2pgsql* tool. It is command-line software especially designed for OSM imports into PostgreSQL/PostGIS databases. However, as documentations state, getting to grips with Osm2pgsql is rather difficult for beginners [21]. No significant results were obtained after several attempts.

An alternative to Osm2pgsql can be found in the GDAL library. In fact, GDAL features an *OSM Vector driver* [13] which can also provide supports for OSM file import, but in a much simpler way. Through a single-line bash command, this driver imports any .pbf data file by splitting it into 5 distinct layers/tables (*points*, *lines*, *multilinestrings*, *multipolygons* and *other_relations*), while relying on custom parameters supplied by the user in the *osmconf.ini* configuration file as in section 3.2.3. Here is an example for importing *buildings* from Uusimaa in our database:

```
1 ogr2ogr -f PostgreSQL "PG:dbname='test_osm' host='localhost'
  port='5432' user='postgres' password='postgres'" buildings-
  uusimaa.osm.pbf -lco COLUMN_TYPES=other_tags=hstore
```

In fact, the actual command template used to import our data contains an additional option, designed to allow imports into a single database, by assigning a unique name to each table imported from a data file. No better solutions were found but to repeat the following code 10 times in total (5 times on each data file for its 5 different layers):

```
1 ogr2ogr -f PostgreSQL "PG:dbname='uusimaa' host='localhost'
  port='5432' user='postgres' password='postgres'" buildings-
  uusimaa.osm.pbf multipolygons -lnl buildings_multipolygons
  -lco COLUMN_TYPES=other_tags=hstore
```

This was an example for importing the *multipolygons* layer of the *buildings* from Uusimaa in our database.

At the end of the import, some tables were empty and were therefore deleted, resulting in a final number of **7 feature data tables** (see appendix C for UML representation).

As for the changesets, they were imported using a Python script (see appendix D) based on the ElementTree XML API [25] for parsing XML data, as well as on the Psycopg library [16], which lets you access a PostgreSQL database from Python. For this to work properly, the creation of a table with a very specific structure is required beforehand, in order to receive incoming data as shown in Figure 3.3 below:

```
1 CREATE TABLE public.changesets
2 (
3     changeset_id integer NOT NULL,
4     created_at timestamp without time zone,
5     closed_at timestamp without time zone,
6     open boolean,
7     "user" text,
8     uid integer,
9     min_lat double precision,
10    min_lon double precision,
11    max_lat double precision,
12    max_lon double precision,
13    num_changes integer,
14    comments_count integer,
15    other_tags hstore,
16    PRIMARY KEY (changeset_id)
17 );
18
19 ALTER TABLE IF EXISTS public.changesets
20   OWNER to postgres;
```

Figure 3.3: changesets table creation SQL code

Note that a particular attention was given to the changesets' tags, where the source information is stored. Since there is not a fixed number of tags on any changeset, all of them were put in a single column using the `hstore` data type.

Eventually there are still JOIN operations to make, but they'll be of interest in the very next chapter, alongside SQL queries for NLS-like strings pattern matching.

3.3 Feature-level modifications collecting

Although it was originally planned to carry out a full analysis of the changes as done in [19], it was decided to concentrate on the *number of geometric modifications* that the features have undergone since their integration into OSM, which is more interesting and useful from the point of view of the NLS.

3.3.1 Intensity rather than quality

This section's work does not measure the magnitude of how a given OSM feature has changed. In fact, a big geometry change will be treated just as a small one, but instead the willingness of the contributors to make changes to the OSM database is assessed.

3.3.2 About the changesets

The natural way for studying the modifications as described in [19] is to look into the *changesets* which were introduced in section 2.3. To be more precise, any changeset is supposed to regroup all the changes one contributor can make during an editing session, i.e. from the moment he opens his editor to the one he commits his changes to the OSM server [4]. In reality, changesets don't provide targeted information on modified features, but only general information about the scope of the changes (opening/closing timestamps of the changeset, geographical extent of the changes, overall number of modifications, etc.), while being referenced at the level of each feature with the *changeset id* to ensure a proper logging.

3.3.3 Digging into the OSM history file

The OSM history file is more promising, as it contains concrete data about all the versions that have ever existed for any OSM feature, thus allowing comparison between the versions to monitor the changes throughout time.

As stated in section 3.1.2, little programs are able to deal with OSM history files. The content of our history file from Uusimaa (`uusimaa-history.osh`) was then imported into the *uusimaa* database as `history` table using a custom Python script (cf. Appendix H), based on the `Ixml` XML parsing Python library [12].

Superior to the native ElementTree API (used in 3.2.4 for the changesets import) in terms of processing time and functionalities, `Ixml` allows an iterative approach, meaning the XML-based history file is read on the fly, thus reducing the memory usage considerably. To free up even more memory, intelligent variable management was put in place with outdated variables' deletion at each iteration. This constant attention to memory management stems from the fact that initial tests were based on storing the entire XML file into memory before reading it. This strategy do work for the history import of small areas like the Otaniemi neighborhood (famous R&D district, where the FGI is located) within a matter of a few seconds, but completely overwhelm the Python kernel capacities with a way bigger region as Uusimaa's (which account for 223.2 MB in the raw binary format, but for 7.73 GB when converted into XML!). For more information, the former code implementing this latter naive approach is available in Appendix I.

Once the history is imported, it needs to be analyzed to locate the geometric modifications. This task was done using again a Python script (cf. Appendix J), which core of its operation is based on the `geometry_modification` recursive function (from line 131 in the code). Due to the OSM data structure, geometric information is exclusively attached to the `node` elements (see *OSM data structure* in section 1.2). As a result, references need to be resolved in order to access

geometric information on *ways* and *relations*, which is handled by the recursive behavior of the `geometry_modification` function, while the `geometry_references_resolver` function offers a *version resolving system* suggestion, as referenced members of ways and relations might also have several versions to consider.

DATA ANALYSIS AND VISUALISATIONS

CHAPTER
4

4.1 Identifying the massive contributions of NLS data in OSM

4.1.1 NLS references at the feature-level source tag

Following the highlighting of the feature-level source tag field in the section 3.2.3, it should now be possible to identify the NLS data. However, as mentioned at the beginning of our study in 1.4, few massive contributions from the NLS have official documentation, adding that there is not a particular set of source tag values provided by the NLS for us to search the exact information we want.

It was therefore necessary to search blindly, considering for example particular values that the source tag would be likely to take to identify the data coming from the NLS. A first NLS source typology was then designed as shown in Figure 4.1, including the two official languages of Finland which are Finnish and Swedish, but also English, this with a view to taking into account as much data as possible.

```
"source" ILIKE '%mml%' OR
"source" ILIKE '%nls%' OR
"source" ILIKE '%maanmittauslaitos%' OR
"source" ILIKE '%national land survey of finland%' OR
"source" ILIKE '%finnish geospatial research institute%' OR
"source" ILIKE '%fgi%' OR
"source" ILIKE '%Paikkatietokeskus%' OR
"source" ILIKE '%lmv%' OR
"source" ILIKE '%Lantmäteriverket%' OR
"source" ILIKE '%Geodatacentralen%'
```

Figure 4.1: First NLS source tag value typology

These values, framed by '%' signs, correspond in fact to patterns, which we search for within the values of the source tags. The Figure 4.2 below shows the results we get for both *buildings* and *roads* data sets, by counting the features which fall under our typology.

Overall, as we can see, the vast majority of features from the 2 datasets have no source information (i.e. source tag left blank). The percentages for NLS data are thus very low, while further tests have shown that also considering *created_by* and *other_tags* tags (introduced in section 3.2.3, with the latter being a tag for everything as described in [13]) doesn't bring much better numbers, with a slight variation of a few dozen features at most (see Appendix E for related numbers).

Sources / Type of features	Points (49,353 features)	Lines (85 features)	Multipolygons (502,786 features)	All features (total)
NLS	0.008% (4 features)	0% (0 features)	6.6% (33,149 features)	6 %
Other	0.7% (363 features)	11.8% (10 features)	6.1% (30,576 features)	5.6%
No source	99.3% (48,987 features)	88.2% (75 features)	87.3% (439,061 features)	88.4 %

Table 1 – Source key values proportions for the “buildings” dataset (552,224 features)

Sources / Type of features	Points (105,452 features)	Lines (468,982 features)	Multipolygons (4,214 features)	All features (total)
NLS	0.2% (199 features)	0.7% (3,410 features)	0.2% (9 features)	0.6%
Other	3.5% (3,718 features)	5.6% (26,092 features)	2% (86 features)	5.2%
No source	96.3% (101,535 features)	93.7% (439,480 features)	97.8% (4,119 features)	94.2%

Table 2 – Source key values proportions for the “roads” dataset (578,648 features)

Figure 4.2: Feature-level source tag values proportions

4.1.2 Adding changeset-level source tag information

Deeper investigation is then needed on the changesets, where the most source information should be contained. From the 7 feature tables and the changesets table created in the section 3.2.4, 3 classes of new tables are defined:

- feature tables beginning with 'nls': contain features falling under the typology introduced in the section 4.1.1 and including matching source information from the changesets
- feature tables beginning with 'no_source': contain features with both feature-level and changeset-level source tag empty
- feature tables beginning with 'other': contain features that do not fall into any of the categories listed above

As a result, 21 additional tables are added to the *uusimaa database*, mostly created by SQL LEFT JOIN operations and string pattern analysis, so to allow quick access to relevant data without complex querying if needed in the future. For more insights on the SQL queries used to build those tables, please take a look at Appendix F.

Here are then the updated results from Figure 4.2 by including additional source information from the changesets:

The inclusion of sources related to changesets seems to be bearing fruit, with overall about 3.6 times more NLS buildings being taken into account (6% previously), and about 16.3 times more

Sources / Type of features	Points (49,353 features)	Lines (85 features)	Multipolygons (502,786 features)	All features (total)
NLS	2.6% (1,305 features)	2.4% (2 features)	23.6% (118,810 features)	21.8%
Other	35.1% (17,300 features)	68.2% (58 features)	27.6% (138,453 features)	28.2%
No source	62.3% (30,748 features)	29.4% (25 features)	48.8% (245,523 features)	50 %

Table 5 – Source key values proportions for the “buildings” dataset (552,224 features)

Sources / Type of features	Points (105,452 features)	Lines (468,982 features)	Multipolygons (4,214 features)	All features (total)
NLS	4.5% (4,692 features)	11.1% (51,861 features)	4.8% (201 features)	9.8%
Other	48.5% (51,168 features)	47.1% (221,146 features)	55.4% (2,335 features)	47.5%
No source	47% (49,592 features)	41.8% (195,975 features)	39.8% (1,678 features)	42.7%

Table 6 – Source key values proportions for the “roads” dataset (578,648 features)

Figure 4.3: Complete source tag values proportions

NLS roads being considered (0.6% before). As more than the half of the features from both data sets have the source tag filled, we consider these numbers as sufficient to continue our study.

4.1.3 Focus on NLS-like source values spelling

A more detailed analysis on the source values spelling for features being categorized as *NLS* can also be interesting. The 20 most recurrent values from the multipolygons layer of the *buildings* data set are shown below in Figure 4.4.

	source character varying	occurrences numeric
1	MML	37229
2	MML/NLS	17648
3	aerial imagery;MML background map	12734
4	MML2017	10322
5	MML Orthophoto	8739
6	MML Topographic Map	8601
7	NLS.fi	4649
8	MML2016	4335
9	MML_2012	2524
10	local knowledge, aerial, MML	2010
11	aerial imagery;MML Ortophoto	1937
12	mml2018/helsingin karttapalvelu.	1925
13	Mapbox Satellite / National Land Survey of Finland	1788
14	MML2018	1261
15	NLSF	1151
16	MML2019/Helsingin kiinteistökartta	1121
17	MML Background Map	959
18	Addresses: MML/DVV.fi; aerial imagery; map imagery;	945
19	MML Background Map; MML Orthophoto; MML Topographic Map	799
20	MML Orthophoto; MML INSPIRE Buildings Maastotietokanta (WMS)	690

Figure 4.4: Source tag value most prominent occurrences from nls_buildings_multipolygons

As we can see, the NLS can be mentioned in numerous ways, with the 'MML' value being the most used without big surprise. However, a significant number of values are being composed of sub-values in accordance with [9], but actually showing the complexity in identifying proper NLS data. Indeed, the Figure 4.4 provide us with just an overview on how NLS is mentioned in OSM, as more than 400 additional different values have been filtered by our SQL query (cf. Appendix G).

From that, **it is impossible to say if one feature is really coming from the NLS**, as no official documentation related to NLS features tagging in OSM exist at this moment to confirm our assumptions. Then, it also supports the fact that the OpenStreetMap project is anarchic with no real constraints on data contributions, in contrary to authoritative data providers such as National agencies, where a lot of standards govern the formatting of data.

As a result, the following paragraphs are intended to provide some insights into features that have a high probability of belonging to imports from the NLS, but with no absolute certainty that this is the case. For the sake of brevity and misuse of language, the features concerned will be referred to as *NLS's*.

4.1.4 NLS massive contributions visualizations

From the identified NLS features in 4.1.2, it is now possible to highlight areas of massive contributions through a spatial pattern analysis, which relies on NLS feature proportions in different geographical locations. This can be achieved with the use of hexagon grids.

A hexagon grid (created using MMQGIS QGIS Plugin [20]) is defined by two parameters:

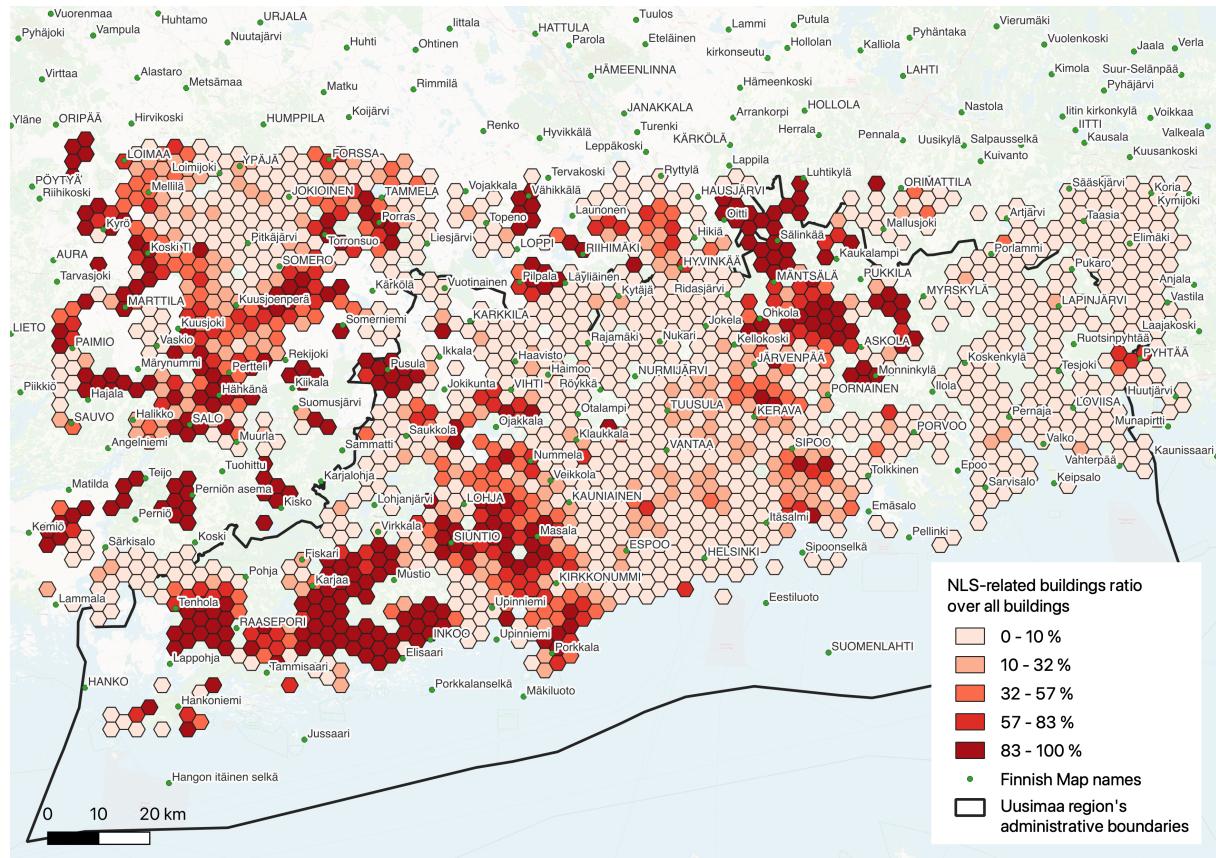
- *the extent of the grid*: covered area by the hexagons
- *the grid interval*: hexagon size

The last parameter is of particular importance: the hexagons shouldn't be too small or too big, with in the former case the result being similar to a simple plot of the data, while in the latter the hexagon size prompts the lose of spatial patterns. In order to take this into account, we'll proceed with an hexagon grid having $2 * \text{apothem} = 2500 \text{ m}$ (ETRS89 / TM35FIN(E,N) [EPSG:3067] CRS Units) as *grid interval*, covering a little more than Uusimaa region's extent.

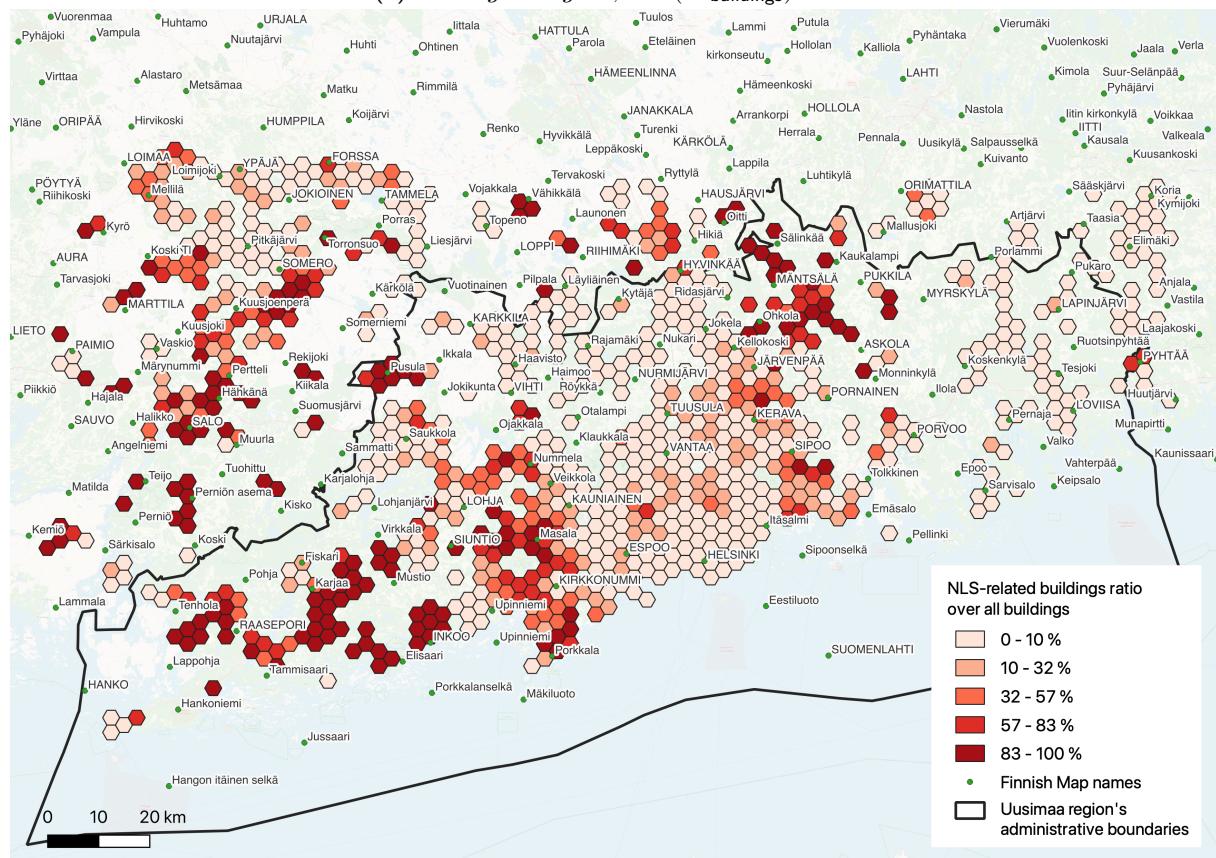
Then, by counting the number of centroids of each polygon from both `nls_buildings_multipolygons` and `buildings_multipolygons` for example inside any hexagon of the grid, the NLS-related buildings ratio over all buildings can be calculated for any hexagon with the following formula:

$$\text{nls_ratio} = \frac{\text{nb of nls_buildings_multipolygons centroids}}{\text{nb of buildings_multipolygons centroids}} * 100$$

One downside of this way of thinking, is that if an hexagon of the grid was to contain only a unique building in its area with this building appearing to be tagged as NLS, then $\text{nls_ratio} = 100\%$, which is completely wrong for our purpose! In fact, the goal here is to identify aggregated groups of features as hexagons. Hence, the *nb of buildings_multipolygons centroids* should be the greatest as possible, as it is an indicator of the size of what we actually consider as a massive contribution. Several map configurations were tested with an ascending minimum number of buildings per hexagon in the Figure 4.5 below.



(a) $\forall \text{hexagon} \in \text{grid}, \min(\text{nbbuildings}) = 20$



(b) $\forall \text{hexagon} \in \text{grid}, \min(\text{nbbuildings}) = 100$

Figure 4.5: Hex maps on NLS-related buildings proportion over all buildings (Map base credits: © OpenStreetMap contributors)

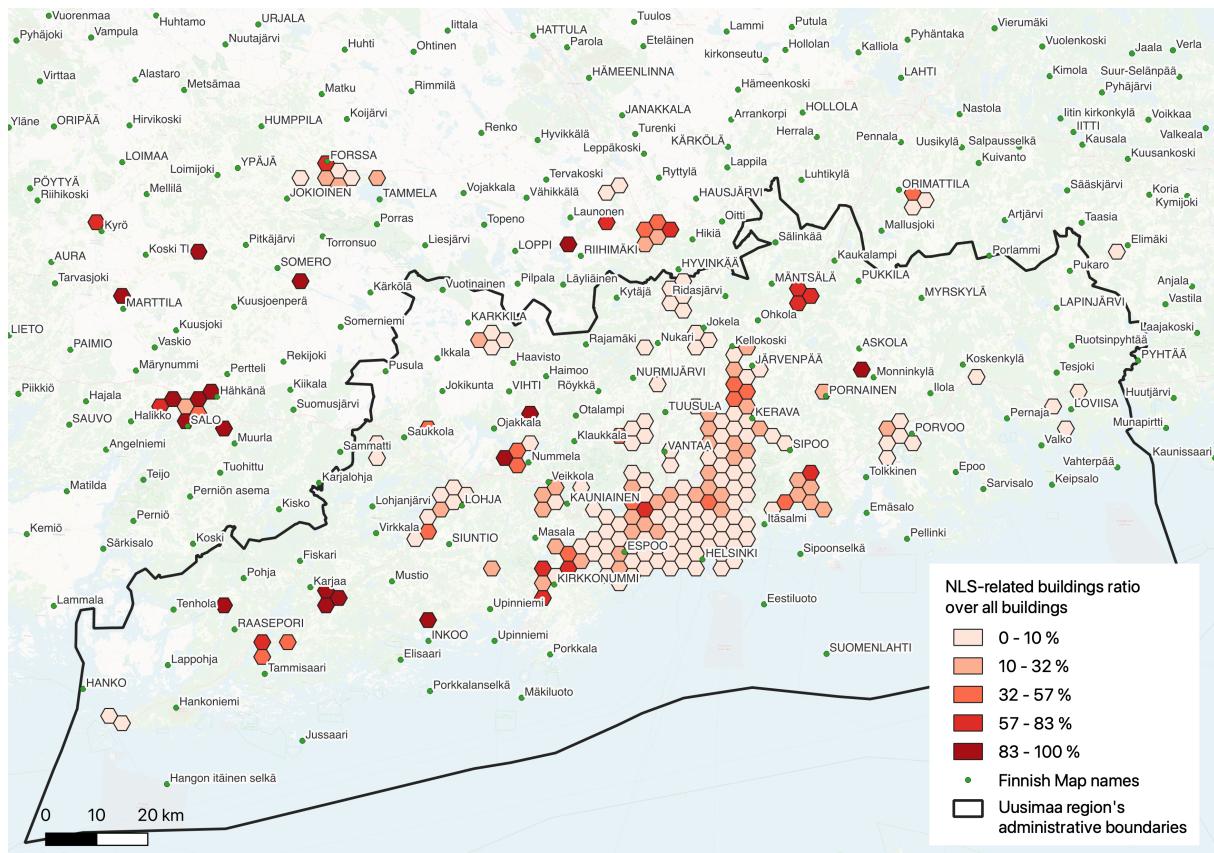


Figure 4.5: Hex maps on NLS-related buildings proportion over all buildings (cont.)

Right from the outset, we note that the density of the grid is inversely proportional to the minimum number of buildings per hexagon. This behavior is consistent, since our selection criterion on the size of massive contributions becomes increasingly restrictive. Then, for a fixed minimum number of buildings per hexagon (i.e. if we focus on one map), hexagons are associated with several percentage categories. The higher the percentage, the higher the rate of NLS buildings in a hexagon, meaning that our massive contributions are concentrated in the upper half of the percentages: the higher the percentage, the more likely it is that the buildings in the hexagon under consideration belong to a massive contribution from NLS.

Indeed, as we've seen it in the previous subsection 4.1.3, it is impossible to determine with certainty the massive contributions from NLS, but this study may support further research in the future. Still, by considering the fifth with the highest ratios from those first results, we can conclude that NLS massive contributions of buildings are likely for the most to be located in the west of the Uusimaa region.

On the other hand, even if the number of buildings per hexagon changes, we can see that the percentages around the Capital Region remain low. One hypothesis for this is that the data in this area comes from the relevant municipalities, from which better quality material than NLS's is available.

4.2 Feature-level modifications underwent by NLS data since its import into OSM

By applying the method described in section 3.3.3, the following numbers in Figure 4.6 were obtained for the NLS-related buildings from Uusimaa.

```
In [3]: runfile('/Users/alexsren/Desktop/test3/nls_buildings_multipolygons_geometry_modif/typology_modif_encoding_copy.py', wdir='/Users/alexsren/Desktop/test3/nls_buildings_multipolygons_geometry_modif')
100%[██████████] 118810/118810 [12:25<00:00, 159.45it/s]

----- Geometry modification statistics -----
Feature count: 118810
Among the 118810 features
- 8240 underwent geometric modifications ( 6.9 % )
- 110555 have not been modified ( 93.1 % )
- 15 could not be analysed ( 0.0 % )
*   Average number of geometry modifications per feature: 0.1
*   Max number of geometry modifications: 24.0

In [4]:
```

Figure 4.6: Statistics over the number of *geometry modification* for nls_buildings_multipolygons

While a negligible number of features encountered some processing errors with the typology_modif_encoding_copy.py code launching, the vast majority of the NLS-related buildings from Uusimaa have not been modified since their import. Overall, the features have also undergone very few changes, but modified features appear to account for a significant proportion in absolute terms, with just over 8,000 modified buildings. It would be worthwhile to work on locating these buildings (related Python code in Appendix K), most of which seem to be clustered around the Capital region as shown in Figure 4.7.

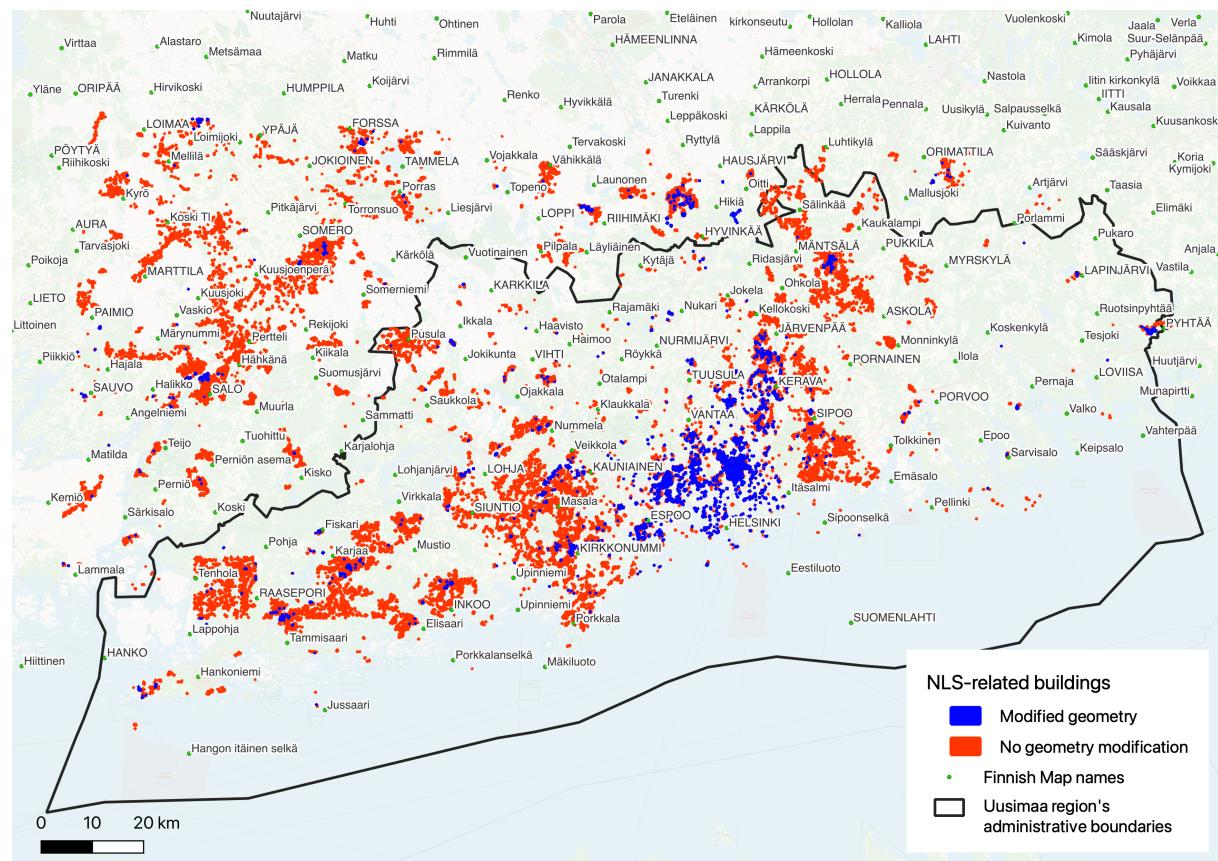


Figure 4.7: Spatial distribution of geometrically modified NLS-related buildings (Map base credits: © OpenStreetMap contributors)

Conclusion

Our study presents an approach to identify massive contributions from a national official agency, the NLS, to the OpenStreetMap database without clear documentation, and to make preliminary assessments of the contributions' intensity on the geometry of the targeted features. We then presented the context of the study, before downloading the data of interest from various data providers. New tools designed for OSM were also used to handle the raw data and transform it to fit the area of Uusimaa, and to retrieve the source key values from both features and changesets, while obtaining the geometric modifications was mainly done by creating our own custom tool using Python, due to the lack of software able to handle historical data from OSM.

Initial results show the importance of including changeset-level source information to obtain the most data of interest, while advanced analysis of the spelling of source values highlights the uncertain nature of our identification process when searching for NLS features. Some simulations were still possible and help to assume that a majority of buildings with massive contributions from the NLS are located in the west of the Uusimaa region. However, considering that in our best case scenario a little less than 50% of the features don't have the source tag filled, further research would be welcome in order to verify if some NLS-related features were not tagged and missed in our research.

Our final discoveries on the geometrically modified features provide insights on the location of the relevant NLS buildings which are of interest to OSM contributors, and can be the start of a deeper investigation on whether the quality of the modified features is acceptable for the NLS to benefit from it. With the emergence of AI-based solutions, these future results could also be compared to those provided by self-sufficient units.

Eventually, the current research should be extended nationwide across the whole of Finland, in order to see any differences or patterns on a larger scale. It may also be interesting to look at more data themes from which different conclusions can be drawn.

References

- [1] GIS StackExchange Community. *Loading .osm.pbf file in QGIS*. StackExchange. 2019. URL: <https://gis.stackexchange.com/questions/39112/loading-osm-pbf-file-in-qgis> (visited on 06/01/2023).
- [2] OpenStreetMap Community. *Buildings*. OpenStreetMap Wiki. 2023. URL: <https://wiki.openstreetmap.org/wiki/Buildings> (visited on 05/30/2023).
- [3] OpenStreetMap Community. *Category:Import from Finland*. OpenStreetMap Wiki. 2013. URL: https://wiki.openstreetmap.org/wiki/Category:Import_from_Finland (visited on 06/28/2023).
- [4] OpenStreetMap Community. *Changeset*. OpenStreetMap Wiki. 2023. URL: <https://wiki.openstreetmap.org/wiki/Changeset> (visited on 06/06/2023).
- [5] OpenStreetMap Community. *Elements*. OpenStreetMap Wiki. 2023. URL: <https://wiki.openstreetmap.org/wiki/Elements> (visited on 06/12/2023).
- [6] OpenStreetMap Community. *Highways*. OpenStreetMap Wiki. 2023. URL: <https://wiki.openstreetmap.org/wiki/Highways> (visited on 05/30/2023).
- [7] OpenStreetMap Community. *Key:converted_by*. OpenStreetMap Wiki. 2022. URL: https://wiki.openstreetmap.org/wiki/Key:converted_by (visited on 06/12/2023).
- [8] OpenStreetMap Community. *Key:created_by*. OpenStreetMap Wiki. 2023. URL: https://wiki.openstreetmap.org/wiki/Key:created_by (visited on 06/12/2023).
- [9] OpenStreetMap Community. *Key:source*. OpenStreetMap Wiki. 2022. URL: <https://wiki.openstreetmap.org/wiki/Key:source> (visited on 06/06/2023).
- [10] OpenStreetMap Community. *OSMCha*. OpenStreetMap Wiki. 2023. URL: <https://wiki.openstreetmap.org/wiki/OSMCha> (visited on 06/06/2023).
- [11] OpenStreetMap Community. *Overpass API*. OpenStreetMap Wiki. 2023. URL: https://wiki.openstreetmap.org/wiki/Overpass_API (visited on 05/29/2023).
- [12] lxml developers. *lxml - XML and HTML with Python*. 2022. URL: <https://lxml.de/> (visited on 08/02/2023).
- [13] gdal.org. *OSM - OpenStreetMap XML and PBF*. GDAL documentation. 2023. URL: <https://gdal.org/drivers/vector/osm.html> (visited on 06/15/2023).
- [14] Geofabrik.de. *Data Extracts - Technical Details*. 2019. URL: <https://osm-internal.download.geofabrik.de/technical.html> (visited on 06/01/2023).
- [15] Andrew Gierth. *hstore*. PostgreSQL 15 Documentation. 2023. URL: <https://www.postgresql.org/docs/15/hstore.html> (visited on 06/20/2023).
- [16] Federico Di Gregorio. *Basic module usage*. Psycopg 2.9.6 documentation. 2021. URL: <https://www.psycopg.org/docs/usage.html> (visited on 06/19/2023).

- [17] Ari Hakahuhta. "Turvallisuuskomitean pääsihteeri: Kun naapuri käy sotaa, kuntien ja yritysten olisi syytä rajoittaa karttojensa jakamista". In: *yle.fi* (May 2023). URL: <https://yle.fi/a/74-20033447> (visited on 05/26/2023).
- [18] Sarah Hoffmann. *Welcome to Pyosmium's documentation!* osmcode.org. 2020. URL: <https://docs.osmcode.org/pyosmium/latest/> (visited on 06/07/2023).
- [19] A. Le Guilcher, A.-M. Olteanu-Raimond, and M. B. Balde. "ANALYSIS OF MASSIVE IMPORTS OF OPEN DATA IN OPENSTREETMAP DATABASE: A STUDY CASE FOR FRANCE". In: *ISPRS Annals of the Photogrammetry, Remote Sensing and Spatial Information Sciences* V-4-2022 (2022), pp. 99–106. DOI: [10.5194/isprs-annals-V-4-2022-99-2022](https://doi.org/10.5194/isprs-annals-V-4-2022-99-2022). URL: <https://isprs-annals.copernicus.org/articles/V-4-2022/99/2022/> (visited on 05/22/2023).
- [20] Michael Minn. *MMQGIS*. 2021. URL: <http://michaelminn.com/linux/mmqgis> (visited on 07/26/2023).
- [21] osm2pgsql.org. *Osm2pgsql Manual*. 2023. URL: <https://osm2pgsql.org/doc/manual.html> (visited on 06/14/2023).
- [22] osmcode.org. *Osmium Library*. URL: <https://osmcode.org/libosmium/> (visited on 06/02/2023).
- [23] osmcode.org. *Osmium manual pages*. URL: <https://docs.osmcode.org/osmium/latest/> (visited on 06/02/2023).
- [24] osmcode.org. *Osmium Tool Manual*. URL: <https://osmcode.org/osmium-tool/manual.html> (visited on 06/02/2023).
- [25] python.org. *xml.etree.ElementTree — The ElementTree XML API*. Python Documentation. 2023. URL: <https://docs.python.org/3/library/xml.etree.elementtree.html#parsing-xml> (visited on 06/19/2023).
- [26] A. Sarretta, M. Napolitano, and M. Minghini. "OPENSTREETMAP AS AN INPUT SOURCE FOR PRODUCING GOVERNMENTAL DATASETS: THE CASE OF THE ITALIAN MILITARY GEOGRAPHIC INSTITUTE". In: *The International Archives of the Photogrammetry, Remote Sensing and Spatial Information Sciences* XLVIII-4/W7-2023 (2023), pp. 193–200. DOI: [10.5194/isprs-archives-XLVIII-4-W7-2023-193-2023](https://doi.org/10.5194/isprs-archives-XLVIII-4-W7-2023-193-2023). URL: <https://isprs-archives.copernicus.org/articles/XLVIII-4-W7-2023/193/2023/> (visited on 07/25/2023).

List of Figures

1.1	Finland's Uusimaa region (Map base credits: © OpenStreetMap contributors)	9
1.1	Finland's Uusimaa region (cont.)	10
2.1	QuickOSM usage example in QGIS — <i>roads</i> from Uusimaa retrieval	13
2.2	Finland's Non-Public data Geofabrik downloading webpage — URL: https://osm-internal.download.geofabrik.de/europe/finland.html	14
2.3	Planet OSM website homepage — URL: https://planet.openstreetmap.org/	15
3.1	Bounding box used for Uusimaa's data clipping, Ing/lat coordinates (22.65,59.61,26.65,60.84) — Uusimaa boundaries in <i>light gray</i> inside the blue rectangle (credits: bboxfinder.com)	18
3.2	<i>pgAdmin 4</i> interface showing contents from the <i>uusimaa</i> database	21
3.3	changesets table creation SQL code	22
4.1	First NLS source tag value typology	25
4.2	Feature-level source tag values proportions	26
4.3	Complete source tag values proportions	27
4.4	Source tag value most prominent occurrences from <i>nls_buildings_multipolygons</i>	28
4.5	Hex maps on NLS-related buildings proportion over all buildings (Map base credits: © OpenStreetMap contributors)	30
4.5	Hex maps on NLS-related buildings proportion over all buildings (cont.)	31
4.6	Statistics over the number of <i>geometry modification</i> for <i>nls_buildings_multipolygons</i>	32
4.7	Spatial distribution of geometrically modified NLS-related buildings (Map base credits: © OpenStreetMap contributors)	33
A.1	Source tags fetching Homemade Python code (1/2)	45
A.2	Source tags fetching Homemade Python code (2/2)	46
A.3	Source tags fetching Homemade Python code — Preliminary results for the first 50 features from <i>roads-uusimaa.osm.pbf</i>	47
C.1	<i>uusimaa</i> database UML Class diagram	53
E.1	NLS-related data proportions with additional 'created_by' and 'other_tags' tags	59
F.1	' <i>nls</i> ' tables SQL creation code	61
F.2	' <i>no_source</i> ' tables SQL creation code	62
F.3	' <i>other</i> ' tables SQL creation code	62
G.1	SQL query to select the most prominent occurrences of NLS-like source values	63

List of Tables

Annexes

A	Source tags fetching with Python	45
B	GDAL Configuration file for OSM import	49
C	First UML representation for uusimaa database	53
D	Python script for importing changesets into uusimaa database	55
E	NLS-related data proportions with additional 'created_by' and 'other_tags' tags	59
F	SQL queries for categorised feature table creation	61
G	SQL query to select the most prominent occurrences of NLS-like source values	63
H	Python script for OSM history file import into PostgreSQL	65
I	Naive Python script for OSM history file import into PostgreSQL	73
J	OSM Feature-level modifications analyser Python script	79
K	Python script for retrieving OSM ids from feature which underwent geometric modifications	91

SOURCE TAGS FETCHING WITH PYTHON

APPENDIX A

```

1  #!/usr/bin/env python3
2  # -*- coding: utf-8 -*-
3  """
4  Created on Tue Jun  6 14:48:43 2023
5
6  @author: alexysren
7  """
8
9
10 ### Imports
11
12 import os # to get our current directory path
13 import time # to compute the elapsed time for each function to return something
14
15 from tqdm import tqdm # make loops show a smart progress meter
16
17 from osmcha import changeset # OpenStreetMap Changeset Analyzer
18 import osmium # PyOsmium
19
20 import numpy as np # numpy
21 import pandas as pd # pandas library
22
23
24 ### Global variables
25
26 current_directory_path = os.getcwd()
27
28
29 ### Functions
30
31
32 class DataHandler(osmium.SimpleHandler):
33     def __init__(self):
34         osmium.SimpleHandler.__init__(self)
35         self.elements = []
36
37     def node(self, n):
38         self.elements.append(["node", n.id, n.changeset])
39
40     def way(self, w):
41         self.elements.append(["way", w.id, w.changeset])
42
43     def relation(self, r):
44         self.elements.append(["relation", r.id, r.changeset])
45
46
47 def get_data(filepath):
48     """
49     imports data from an OSM data file and returns features id along its
50     associated changeset id in a Panda Dataframe
51
52     Parameter:
53         filepath (str): path to OSM data file
54
55     Return:
56         (DataFrame): Panda DataFrame containing features id and associated changeset id
57     """
58     dataHandler = DataHandler()
59     dataHandler.apply_file(filepath)
60     col_names = ["type", "id", "changesetId"]
61     return pd.DataFrame(dataHandler.elements, columns=col_names)
62
63
64
65
66
67
68 def get_source(changeset_id):
69     """
70     retrieves the source key value from an input changeset id
71
72     Parameter:
73         changeset_id (int): changeset id
74
75     Return:
76         (str): source key value
77     """
78     ch = changeset.Analyse(changeset_id)
79     return ch.source
80

```

Figure A.1: Source tags fetching Homemade Python code (1/2)

```

81
82
83
84
85 def apply_source(data, out_filename="data_with_source.csv"):
86     """
87         add the source key values column to the input that corresponds to each feature in
88         the input data
89         the result is also saved in a .csv file in root
90
91     Parameter:
92         data (DataFrame): Panda DataFrame containing features id and associated changeset id
93         out_filename (str): filename for the output file (with its extension!)
94
95     Return:
96         (DataFrame): Panda DataFrame with columns [type, id, changesetId, source] (in order)
97         """
98
99     print("\n\n----- Applying source process started... ----- \n\n")
100    start_time = time.time()
101
102
103    # new_data = data[:50].copy()
104    new_data = data.copy()
105    source_column = np.array([])
106
107
108    for i in tqdm(new_data.index):
109        changeset_id = int(new_data.loc[i, "changesetId"])
110        source = get_source(changeset_id)
111
112        if source != "Not reported":
113            source_column = np.append(source_column, source)
114        else:
115            source_column = np.append(source_column, "")
116
117    new_data['source'] = source_column.tolist()
118
119
120
121    new_data.to_csv(out_filename, float_format=".3f", index_label='index', sep=" ")
122
123
124    end_time = time.time()
125
126    print("\n\n----- Job's finished! ----- \nElapsed time: ", round(end_time - start_time, 2), "seconds\n\n")
127
128
129    return new_data
130
131
132
133
134
135    ### Main
136
137
138    if __name__ == '__main__':
139
140        # Data imports
141
142        roads_filename = "roads-uusimaa.osm.pbf"
143        buildings_filename = "buildings-uusimaa.osm.pbf"
144
145        roads_uusimaa = get_data(current_directory_path + "/" + roads_filename)
146
147        # Source values retrieval
148
149        roads_uusimaa_source = apply_source(roads_uusimaa, "roads-uusimaa_source.csv")

```

Figure A.2: Source tags fetching Homemade Python code (2/2)

```

1 index type id changesetId source
2 0 node 130190 9486041
3 1 node 131938 6975707
4 2 node 131939 1188066
5 3 node 131940 2827488
6 4 node 131941 10996384
7 5 node 131954 11897653
8 6 node 131955 11897653
9 7 node 131956 11029829
10 8 node 131957 1370558
11 9 node 131979 3080664
12 10 node 131980 59214754 "Helsinki ortoilmakuva 2017"
13 11 node 131981 59214754 "Helsinki ortoilmakuva 2017"
14 12 node 131982 11614777
15 13 node 131983 11614777
16 14 node 131984 11614777
17 15 node 132101 134857011 "Survey and Helsinki region orthophoto"
18 16 node 132102 22946644
19 17 node 132103 62265520 "Helsingin ortokuva 2017"
20 18 node 132104 62265520 "Helsingin ortokuva 2017"
21 19 node 132105 62265520 "Helsingin ortokuva 2017"
22 20 node 132106 62265520 "Helsingin ortokuva 2017"
23 21 node 132107 17568343
24 22 node 132108 59214754 "Helsinki ortoilmakuva 2017"
25 23 node 132109 59351200 "Helsinki ortoilmakuva 2017, local_knowledge"
26 24 node 132110 428289
27 25 node 132111 5732909
28 26 node 132112 3778762
29 27 node 132113 118535241 survey
30 28 node 132116 59214754 "Helsinki ortoilmakuva 2017"
31 29 node 132118 6474559
32 30 node 134755 804632
33 31 node 134757 124682564 "HSY Orto;Mapillary;survey"
34 32 node 134759 57701713 "Vantaa ortoilmakuva 2017"
35 33 node 134766 645700
36 34 node 134767 124681122 "HSY Orto;Mapillary;survey"
37 35 node 134768 57701713 "Vantaa ortoilmakuva 2017"
38 36 node 134769 57701713 "Vantaa ortoilmakuva 2017"
39 37 node 134771 99992590 "local knowledge, aerial"
40 38 node 134772 6600395
41 39 node 134773 12732517
42 40 node 134774 6771897
43 41 node 134775 3196631
44 42 node 134776 5957783
45 43 node 134777 67983300
46 44 node 134778 67983300
47 45 node 134779 40464087
48 46 node 134783 16173294
49 47 node 134789 75374370
50 48 node 134790 12612790
51 49 node 134791 12612790

```

Figure A.3: Source tags fetching Homemade Python code — Preliminary results for the first 50 features from `roads-uusimaa.osm.pbf`

GDAL CONFIGURATION FILE FOR OSM IMPORT

APPENDIX **B**

Listing B.1: osmconf.ini (original version)

```
1 #
2 # Configuration file for OSM import
3 #
4
5 # put here the name of keys, or key=value, for ways that are
6 # assumed to be polygons if they are closed
7 # see http://wiki.openstreetmap.org/wiki/Map_Features
8 closed_ways_are_polygons=aeroway,amenity,boundary,building,craft,
9 geological,historic,landuse,leisure,military,natural,office,
10 place,shop,sport,tourism,highway=platform,public_transport=
11 platform
12
13 # Uncomment to avoid laundering of keys ( ':' turned into '_' )
14 #attribute_name_laundering=no
15
16 # Some tags, set on ways and when building multipolygons,
17 # multilinestrings or other_relations,
18 # are normally filtered out early, independent of the 'ignore'
19 # configuration below.
20 # Uncomment to disable early filtering. The 'ignore' lines below
21 # remain active.
22 #report_all_tags=yes
23
24 # uncomment to report all nodes, including the ones without any (
25 # significant) tag
26 #report_all_nodes=yes
27
28 # uncomment to report all ways, including the ones without any (
29 # significant) tag
30 #report_all_ways=yes
31
32 [points]
33 # common attributes
34 osm_id=yes
35 osm_version=no
36 osm_timestamp=no
37 osm_uid=no
38 osm_user=no
39 osm_changeset=no
```

```

32 # keys to report as OGR fields
33 attributes=name,barrier,highway,ref,address,is_in,place,man_made
34 # keys that, alone, are not significant enough to report a node
35 # as a OGR point
36 unimportant=created_by,converted_by,source,time,ele,attribution
37 # keys that should NOT be reported in the "other_tags" field
38 ignore=created_by,converted_by,source,time,ele,note,todo,
39 openGeoDB:,fixmeFIXME
40 # uncomment to avoid creation of "other_tags" field
41 #other_tags=no
42 # uncomment to create "all_tags" field. "all_tags" and "
43 # all_tags=yes
44 [lines]
45 # common attributes
46 osm_id=yes
47 osm_version=no
48 osm_timestamp=no
49 osm_uid=no
50 osm_user=no
51 osm_changeset=no
52 # keys to report as OGR fields
53 attributes=name,highway,waterway,aerialway,barrier,man_made
54
55 # type of attribute 'foo' can be changed with something like
56 #foo_type=Integer/Real/String/DateTime
57
58 # keys that should NOT be reported in the "other_tags" field
59 ignore=created_by,converted_by,source,time,ele,note,todo,
60 openGeoDB:,fixmeFIXME
61 # uncomment to avoid creation of "other_tags" field
62 #other_tags=no
63 # uncomment to create "all_tags" field. "all_tags" and "
64 # all_tags=yes
65 #computed_attributes must appear before the keywords _type and
66 #_sql
67 computed_attributes=z_order
68 z_order_type=Integer
69 # Formula based on https://github.com/openstreetmap/osm2pgsql/
70 # blob/master/style.lua#L13
71 # [foo] is substituted by value of tag foo. When substitution is
72 # not wished, the [ character can be escaped with \[ in literals
73 # Note for GDAL developers: if we change the below formula, make
74 # sure to edit ogrosmysql.cpp since it has a hardcoded
75 # optimization for this very precise formula
76 z_order_sql="SELECT (CASE [highway] WHEN 'minor' THEN 3 WHEN '
77 road' THEN 3 WHEN 'unclassified' THEN 3 WHEN 'residential' THEN

```

```

3 WHEN 'tertiary_link' THEN 4 WHEN 'tertiary' THEN 4 WHEN ,
secondary_link' THEN 6 WHEN 'secondary' THEN 6 WHEN ,
primary_link' THEN 7 WHEN 'primary' THEN 7 WHEN 'trunk_link',
THEN 8 WHEN 'trunk' THEN 8 WHEN 'motorway_link' THEN 9 WHEN ,
motorway' THEN 9 ELSE 0 END) + (CASE WHEN [bridge] IN ('yes', ,
true', , '1') THEN 10 ELSE 0 END) + (CASE WHEN [tunnel] IN ('yes
', , 'true', , '1') THEN -10 ELSE 0 END) + (CASE WHEN [railway] IS
NOT NULL THEN 5 ELSE 0 END) + (CASE WHEN [layer] IS NOT NULL
THEN 10 * CAST([layer] AS INTEGER) ELSE 0 END)"

72
73 [multipolygons]
74 # common attributes
75 # note: for multipolygons, osm_id=yes instantiates a osm_id field
    for the id of relations
76 # and a osm_way_id field for the id of closed ways. Both fields
    are exclusively set.
77 osm_id=yes
78 osm_version=no
79 osm_timestamp=no
80 osm_uid=no
81 osm_user=no
82 osm_changeset=no
83
84 # keys to report as OGR fields
85 attributes=name,type,aeroway,amenity,admin_level,barrier,boundary
    ,building,craft,geological,historic,land_area,landuse,leisure,
    man_made,military,natural,office,place,shop,sport,tourism
86 # keys that should NOT be reported in the "other_tags" field
87 ignore=area,created_by,converted_by,source,time,ele,note/todo,
    openGeoDB:,fixme,FIXME
88 # uncomment to avoid creation of "other_tags" field
89 #other_tags=no
90 # uncomment to create "all_tags" field. "all_tags" and "
    other_tags" are exclusive
91 #all_tags=yes
92
93 [multilinestrings]
94 # common attributes
95 osm_id=yes
96 osm_version=no
97 osm_timestamp=no
98 osm_uid=no
99 osm_user=no
100 osm_changeset=no
101
102 # keys to report as OGR fields
103 attributes=name,type
104 # keys that should NOT be reported in the "other_tags" field
105 ignore=area,created_by,converted_by,source,time,ele,note/todo,
    openGeoDB:,fixme,FIXME
106 # uncomment to avoid creation of "other_tags" field

```

```
107 #other_tags=no
108 # uncomment to create "all_tags" field. "all_tags" and "
     other_tags" are exclusive
109 #all_tags=yes
110
111 [other_relations]
112 # common attributes
113 osm_id=yes
114 osm_version=no
115 osm_timestamp=no
116 osm_uid=no
117 osm_user=no
118 osm_changeset=no
119
120 # keys to report as OGR fields
121 attributes=name,type
122 # keys that should NOT be reported in the "other_tags" field
123 ignore=area,created_by,converted_by,source,time,ele,note,todo,
     openGeoDB:,fixmeFIXME
124 # uncomment to avoid creation of "other_tags" field
125 #other_tags=no
126 # uncomment to create "all_tags" field. "all_tags" and "
     other_tags" are exclusive
127 #all_tags=yes
```

FIRST UML REPRESENTATION FOR UUSIMAA DATABASE

APPENDIX C

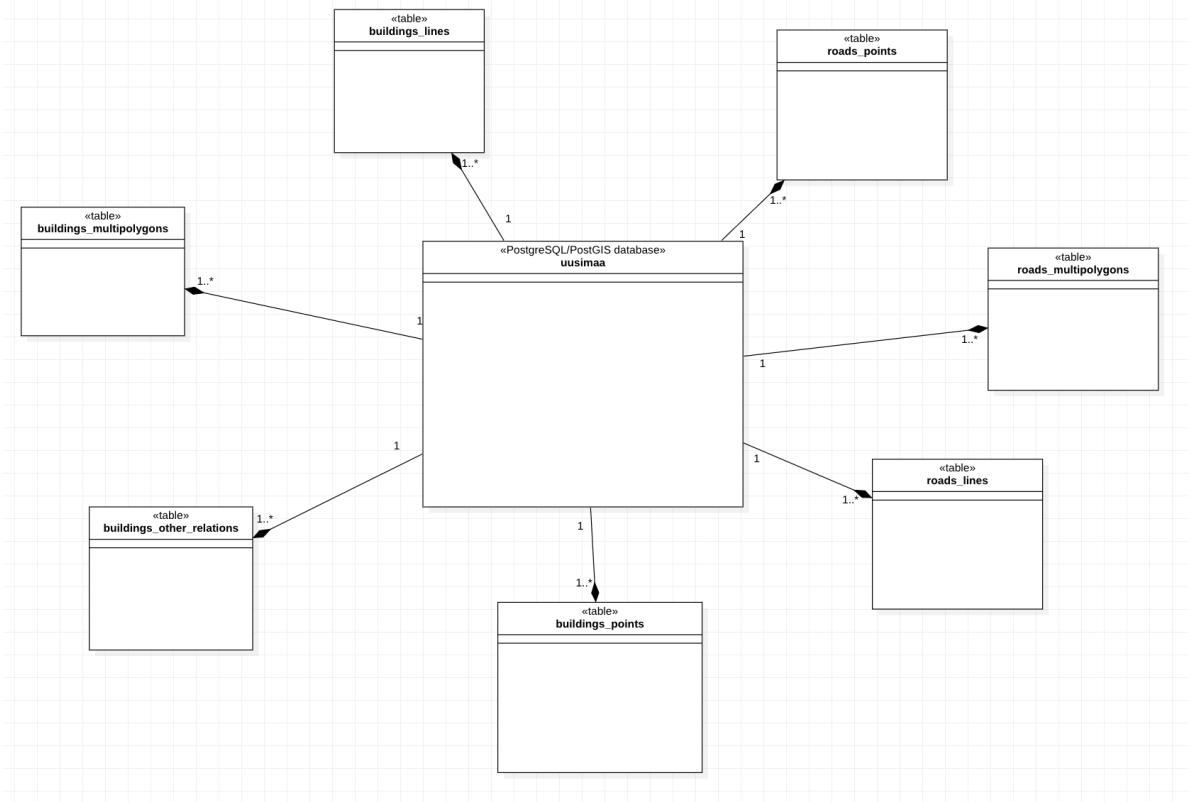


Figure C.1: `uusimaa` database UML Class diagram

PYTHON SCRIPT FOR IMPORTING CHANGESETS INTO UUSIMAA DATABASE

APPENDIX **D**

Listing D.1: chgset_to_pgsql.py Homemade Python script

```
1  #!/usr/bin/env python3
2  # -*- coding: utf-8 -*-
3  """
4  Created on Tue Jun 20 14:31:54 2023
5
6
7  @author: alexysren
8  """
9
10
11 """
12 -----
13
14             OSM CHANGESET FILE IMPORT TO POSTGRESQL
15
16             * custom tool by Alexys Ren - June 2023 *
17
18
19 -----
20 """
21
22
23 """
24 Important prerequisite: if it happens that you have the raw .osm.
25     bz2 file, please extract the .osm from it
26         before launching this code
27
28
29
30 ### Librairies
31
32 from tqdm import tqdm
33
34 import xml.etree.ElementTree as ET
35
```

```
36 import psycopg2
37
38
39
40
41
42
43 ### Changeset file import
44
45 """ /!\ PLEASE MAKE SURE THIS CURRENT PYTHON FILE IS LOCATED IN
46     THE SAME DIRECTORY AS YOUR
47             CHANGESET FILE /!\
48 """
49
50
51
52
53
54 ### Connecting to the PostgreSQL Database
55 # feel free again to change the connection details to fit your
56 # needs
57
58 conn = psycopg2.connect(database="uusimaa", user="postgres",
59                         password="postgres", host="localhost", port="5432")
60 cur = conn.cursor() # to perform database operations
61
62
63
64 ### Code
65
66 root = tree.getroot()
67
68
69 changesets = []
70
71 print("\n\n Parsing changesets... ")
72
73 for child in tqdm(root):
74     # every try/except below are meant to test if the related
75     # info exists!
76     try:
77         changeset_id = child.attrib['id']
78     except KeyError:
79         changeset_id = None
80     try:
81         created_at = child.attrib['created_at']
82     except KeyError:
```

```

82         created_at = None
83     try:
84         closed_at = child.attrib['closed_at']
85     except KeyError:
86         closed_at = None
87     try:
88         Open = child.attrib['open']
89     except KeyError:
90         Open = None
91     try:
92         user = child.attrib['user']
93     except KeyError:
94         user = None
95     try:
96         uid = child.attrib['uid']
97     except KeyError:
98         uid = None
99     try:
100        min_lat = child.attrib['min_lat']
101    except KeyError:
102        min_lat = None
103    try:
104        min_lon = child.attrib['min_lon']
105    except KeyError:
106        min_lon = None
107    try:
108        max_lat = child.attrib['max_lat']
109    except KeyError:
110        max_lat = None
111    try:
112        max_lon = child.attrib['max_lon']
113    except KeyError:
114        max_lon = None
115    try:
116        num_changes = child.attrib['num_changes']
117    except KeyError:
118        num_changes = None
119    try:
120        comments_count = child.attrib['comments_count']
121    except KeyError:
122        comments_count = None
123
124
125
126 # Preparing the additional tags for storing in hstore format
127 # (if they exist!)
128
129 additional_tags = []
130
131 try:
132     temp = child[0]

```

```
132     except IndexError:
133         additional_tags = None
134     else:
135         for tag in child:
136             key = tag.attrib['k']
137             value = tag.attrib['v']
138             additional_tags.append([key, value])
139
140     changesets.append((changeset_id, created_at, closed_at, Open,
141                       user, uid, min_lat, min_lon, max_lat, max_lon, num_changes,
142                       comments_count, additional_tags))
143
144 #### Inserting the output data into PostgreSQL
145
146 # the targeted table structure has to match with the columns
147 # pattern below, and also the table name must be the same!
148 # feel free to change if needed!
149 print("\n\n Inserting changesets into database... ")
150
151 for changeset in tqdm(changesets):
152     cur.execute("""INSERT INTO changesets VALUES (%s, %s, %s, %s,
153               %s, %s, %s, %s, %s, %s, %s, %s, hstore(%s));""", changeset
154 )
155
156 conn.commit() # Make the changes to the database persistent
157
158 cur.close()
159 conn.close()
```

NLS-RELATED DATA PROPORTIONS WITH ADDITIONAL 'CREATED_BY' AND 'OTHER_TAGS' TAGS

APPENDIX **E**

Tagging / Type of features	Points (49,353 features)	Lines (85 features)	Multipolygons (502,786 features)	All features (total)
NLS	0.008% (4 features)	0% (0 features)	6.6% (33,159 features)	6 %
Other	99.992% (49,349 features)	100% (85 features)	93.4% (469,627 features)	94 %

Table 3 — NLS-related features proportions for the “buildings” dataset (552,224 features)

Tagging / Type of features	Points (105,452 features)	Lines (468,982 features)	Multipolygons (4,214 features)	All features (total)
NLS	0.2% (201 features)	0.7% (3,458 features)	0.2% (9 features)	0.6%
Other	99.8% (105,251 features)	99.3% (465,524 features)	99.8% (4,205 features)	99.4%

Table 4 — NLS-related features proportions for the “roads” dataset (578,648 features)

Figure E.1: NLS-related data proportions with additional ‘created_by’ and ‘other_tags’ tags

SQL QUERIES FOR CATEGORISED FEATURE TABLE CREATION

APPENDIX **F**

Those are examples for the `lines` layer from `roads` data set.

```
1 -- Making data output permanent (for QGIS visualisation)
2 CREATE TABLE nls_roads_lines AS
3
4 -- SELECT operation on features coming from the NLS (see WHERE condition below)
5 SELECT * FROM
6
7 -- LEFT JOIN table between roads_lines & changesets
8 (SELECT r.osm_id AS osm_feature_id, c.changesetId,
9 r.source AS feature_source,
10 c.other_tags -> 'source' AS changeset_source,
11 r.wkb_geometry
12
13 FROM roads_lines AS r LEFT JOIN changesets AS c
14 ON r.osm_changeset = c.changesetId) AS temp
15
16 -- one feature is selected if NLS is mentioned either in feature_source
17 -- or in changeset_source
18 WHERE
19 temp.feature_source ILIKE '%mml%' OR
20 temp.feature_source ILIKE '%nls%' OR
21 temp.feature_source ILIKE '%maanmittauslaitos%' OR
22 temp.feature_source ILIKE '%national land survey of finland%' OR
23 temp.feature_source ILIKE '%finnish geospatial research institute%' OR
24 temp.feature_source ILIKE '%fgi%' OR
25 temp.feature_source ILIKE '%Paikkatietokeskus%' OR
26 temp.feature_source ILIKE '%lmv%' OR
27 temp.feature_source ILIKE '%Lantmäteriverket%' OR
28 temp.feature_source ILIKE '%Geodatacentralen%' OR
29
30 temp.changeset_source ILIKE '%mml%' OR
31 temp.changeset_source ILIKE '%nls%' OR
32 temp.changeset_source ILIKE '%maanmittauslaitos%' OR
33 temp.changeset_source ILIKE '%national land survey of finland%' OR
34 temp.changeset_source ILIKE '%finnish geospatial research institute%' OR
35 temp.changeset_source ILIKE '%fgi%' OR
36 temp.changeset_source ILIKE '%Paikkatietokeskus%' OR
37 temp.changeset_source ILIKE '%lmv%' OR
38 temp.changeset_source ILIKE '%Lantmäteriverket%' OR
39 temp.changeset_source ILIKE '%Geodatacentralen%';
```

Figure F.1: '`nls`' tables SQL creation code

```

1 -- Making data output permanent (for QGIS visualisation)
2 CREATE TABLE no_source_roads_lines AS
3
4 -- SELECT operation on features with no source (see WHERE condition below)
5 SELECT * FROM
6
7 -- LEFT JOIN table between roads_lines & changesets
8 (SELECT r.osm_id AS osm_feature_id, c.changeset_id,
9 r.source AS feature_source,
10 c.other_tags -> 'source' AS changeset_source,
11 r.wkb_geometry
12
13 FROM roads_lines AS r LEFT JOIN changesets AS c
14 ON r.osm_changeset = c.changeset_id) AS temp
15
16
17 WHERE temp.feature_source IS NULL
18 AND
19 temp.changeset_source IS NULL;

```

Figure F.2: 'no_source' tables SQL creation code

```

1 -- Making data output permanent (for QGIS visualisation)
2 CREATE TABLE other_roads_lines AS
3
4 -- SELECT operation on features that are neither from the NLS and neither with no source (see WHERE condition below)
5 SELECT * FROM roads_lines
6
7 WHERE NOT EXISTS (
8     SELECT 1
9     FROM nls_roads_lines
10    WHERE nls_roads_lines.osm_feature_id = roads_lines.osm_id
11 )
12 AND NOT EXISTS (
13     SELECT 1
14     FROM no_source_roads_lines
15    WHERE no_source_roads_lines.osm_feature_id = roads_lines.osm_id
16 );

```

Figure F.3: 'other' tables SQL creation code

SQL QUERY TO SELECT THE MOST PROMINENT OCCURRENCES OF NLS-LIKE SOURCE VALUES

```
1  SELECT * FROM (
2    SELECT temp.source, SUM(temp.occurrences) AS occurrences FROM (
3      SELECT feature_source AS source, COUNT(*) AS occurrences FROM nls_buildings_multipolygons GROUP by source
4      UNION
5      SELECT changeset_source, COUNT(*) FROM nls_buildings_multipolygons GROUP BY changeset_source
6    ) AS temp
7    GROUP BY temp.source
8  ) AS tmp
9
10 WHERE
11 tmp.source ILIKE '%mml%' OR
12 tmp.source ILIKE '%nls%' OR
13 tmp.source ILIKE '%maanmittauslaitos%' OR
14 tmp.source ILIKE '%national land survey of finland%' OR
15 tmp.source ILIKE '%finnish geospatial research institute%' OR
16 tmp.source ILIKE '%fgi%' OR
17 tmp.source ILIKE '%Paikkatietokeskus%' OR
18 tmp.source ILIKE '%lmv%' OR
19 tmp.source ILIKE '%Lantmäteriverket%' OR
20 tmp.source ILIKE '%Geodatacentralen%'
21
22 ORDER BY tmp.occurrences DESC;
```

Figure G.1: SQL query to select the most prominent occurrences of NLS-like source values

PYTHON SCRIPT FOR OSM HISTORY FILE IMPORT INTO POSTGRESQL

APPENDIX **H**

Listing H.1: history_to_pgsql_enhanced.py Homemade Python script

```
1 #!/usr/bin/env python3
2 # -*- coding: utf-8 -*-
3 """
4 Created on Wed Jul 19 11:42:19 2023
5
6 @author: alexysren
7 """
8
9
10
11 """
12 -----
13
14         OSM HISTORY FILE IMPORT INTO POSTGRESQL
15             [ENHANCED VERSION]
16
17             * custom tool by Alexys Ren      August 2023 *
18
19
20 -----
21 """
22
23 """
24 User guideline: Please take a look at the 'MAIN' part at the end
25     of this script,
26             so to make the necessary changes to import your
27                 own OSM history file !
28
29
30
31 ### Librairies import
32
33 from tqdm import tqdm      # for displaying the elapsed time after
34                             launching the code, and the number of iterations per second
```

```

34 import psycopg2      # PostgreSQL driver for Python support
35 from psycopg2 import sql # for generating dynamically SQL queries
   (for choosing dynamically a table name)
36 import lxml.etree    # for processing XML data, using C-based
   libraries libxml2 and libxslt      superior to the native
   ElementTree API in terms of processing time and functionalities
37
38
39
40
41 ### Code
42
43
44 def history_importer(osm_history_filename, pgsql_tablename='
  OSM_history'):
45     """
46         imports an OSM history file (XML) as a new table in the
        PostgreSQL database
47     provided by the user in the 'MAIN'
48
49     Parameters
50     -----
51     osm_history_filename : string
        OSM history filename with its extension.
53     pgsql_tablename : string
        Tablename for output PostgreSQL table.
54
55     Returns
56     -----
57     None.
58
59     """
60
61
62     # Creating a new PostgreSQL table for hosting history data
63     cur.execute(sql.SQL("CREATE TABLE public.{}"
64 (
65         id bigint,
66         osm_id bigint,
67         osm_type text,
68         version integer,
69         "timestamp" timestamp without time zone,
70         uid bigint,
71         "user" text,
72         changeset_id bigint,
73         visible boolean,
74         lat double precision,
75         lon double precision,
76         members_refs bigint[],
77         other_tags hstore,
78         PRIMARY KEY (id)
79 );""").format(sql.Identifier(pgsql_tablename)))

```

```

80
81     # Making the previous change to the database persistent
82     conn.commit()
83
84
85     with open(osm_history_filename, "rb") as f:
86
87         id_primary_key = 0
88
89
90         print("\n\n Parsing history from XML, on-the-fly
91             importing... ")
92         print("\n N.B.: For reference, importing a 7.73GB history
93             file took me 1:30h (around 22 million feature versions
94             ). Please also consider your memory capacity!")
95         print("\n\n ** Ongoing process, please make sure the
96             screen saver mode is disabled on your computer! (
97                 otherwise the kernel will be interrupted!) **")
98
99
100        for event, child in tqdm(lxml.etree.iterparse(f)):
101
102
103        osm_type = child.tag
104
105        if osm_type == 'osm':
106            # Ignore 'OSM' XML tag (metadata)
107            continue
108
109        if osm_type in ['node', 'way', 'relation']:
110
111            try:
112                if osm_id == child.attrib['id'] and version
113                    == child.attrib['version']:
114
115                    """ The closing XML tag for a OSM feature
116                        can either be '/>' or either like '</
117                            feature_type>'.
118
119                    In the case it is '</feature_type>', we
120                        jump directly to the next iteration, as
121                        the iterparse function considers this
122                        closing XML tag as a OSM feature in its
123                        own right
124
125                    with all the information related to the
126                        very previous feature """
127
128                    continue
129
130            except UnboundLocalError: # triggered only once
131                at the first iteration since both 'osm_id' and
132                'version' are not initialised yet!
133
134                pass
135
136
137            osm_id = child.attrib['id']
138            version = child.attrib['version']

```

```

116
117     id_primary_key += 1 # inserting a new id column
118             so every record has a unique identifier
119             # (one feature can have
120             several versions using the
121             same feature ID!)
122
123 # every try/except below is meant to test if the
124 # related info exists!
125
126 try:
127     timestamp = child.attrib['timestamp']
128 except KeyError:
129     timestamp = None
130
131 try:
132     uid = child.attrib['uid']
133 except KeyError:
134     uid = None
135
136 try:
137     user = child.attrib['user']
138 except KeyError:
139     user = None
140
141
142 if osm_type == 'node': # if the OSM object is a
143     node
144     try:
145         lat = child.attrib['lat']
146     except KeyError:
147         lat = None
148     try:
149         lon = child.attrib['lon']
150     except KeyError:
151         lon = None
152 else:
153     lat = None
154     lon = None
155
156
157 # Saving currently available data on the current
158 # OSM feature in a list
159 history_record = [id_primary_key,
160                   osm_id,
161                   osm_type,

```

```

161                     version,
162                     timestamp,
163                     uid,
164                     user,
165                     changeset_id,
166                     visible,
167                     lat,
168                     lon]
169
170         """
171         (A): The 'etree.iterparse' method from lxml
172             read all tags and members related to any
173             OSM feature before the feature itself.
174             Below, we retrieve any additional
175             information
176             about the current OSM feature, and append it
177             to the end of 'history_record' list.
178
179         """
180
181     try:
182         if len(next_members_refs) > 0:
183             # If member references were read
184             history_record.append(next_members_refs)
185         else:
186             history_record.append(None)
187     except UnboundLocalError: # triggered only once
188         at the first iteration since 'next_members_refs'
189         ' is not initialised yet!
190         history_record.append(None)
191
192     try:
193         if len(next_additional_tags) > 0:
194             # If tags were read
195             history_record.append(
196                 next_additional_tags)
197         else:
198             history_record.append(None)
199     except UnboundLocalError: # triggered only once
200         at the first iteration since '
201         next_additional_tags' is not initialised yet!
202         history_record.append(None)
203
204     """
205         No more information about the current OSM
206         feature is contained in the XML file, we
207         proceed with the insertion phase.
208
209         """
210
211     # SQL query for inserting the data

```

```

200     # N.B.: the 'INSERT' operation is repeated over
201     # all the OSM feature versions in the OSM history
202     # file
203     cur.execute(sql.SQL("""INSERT INTO {} VALUES (%s,
204         %s, %s, %s, %s, %s, %s, %s, %s, %s, %s,
205         hstore(%s));""").format(sql.Identifier(
206             pgsql_tablename)), history_record)
207
208     # Make the changes to the database persistent
209     conn.commit()
210
211     # Freeing the memory
212     del history_record
213
214     # Freeing memory and variable initialisation for
215     # the next iteration
216     try:
217         del next_members_refs
218     except NameError: # triggered only once at the
219         # first iteration since 'next_members_refs' is
220         # not initialised yet!
221         pass
222     next_members_refs = []
223     try:
224         del next_additional_tags
225     except NameError: # triggered only once at the
226         # first iteration since 'next_additional_tags' is
227         # not initialised yet!
228         pass
229     next_additional_tags = []
230
231
232
233     # if the XML element is actually a OSM tag
234     """
235         See (A).
236     """
237
238
239     # if the tag is a referenced member of a way or a
240     # relation
241     if osm_type == "nd" or osm_type == "member":
242
243         member_id = int(child.attrib['ref'])
244         next_members_refs.append(member_id)

```

```

240             # if the tag is just a common attribute
241         else:
242
243             key = child.attrib['k']
244             value = child.attrib['v']
245
246             tag = [key,value]
247
248             next_additional_tags.append(tag)
249
250
251             # Freeing the memory for the next iteration
252             child.clear()
253
254     print("\n\n Process completed! \n Your OSM history file
255           content should now be appearing in your '", pgsql_tablename
256           , "' PostgreSQL table!")
257
258
259     ### Main      Code execution
260
261
262     if __name__ == '__main__':
263
264         # Connecting to the PostgreSQL Database
265         conn = psycopg2.connect(database="uusimaa", user="postgres",
266                               password="postgres", host="localhost", port="5432")
267
268         # Opening a cursor to perform database operations
269         cur = conn.cursor()
270
271         # Selecting the OSM history file for import
272         """
273             Important prerequisite: if it happens that you have the raw .
274                 osh.pbf file,
275                         please convert it to .osh format (XML
276                             ) before launching this code
277                         ('osmium cat' command can be used for
278                             that purpose)
279             """
280
281         #filename = 'otaniemi-history.osh'
282         filename = 'uusimaa-history.osh'
283
284         """ /!\ PLEASE MAKE SURE THIS CURRENT PYTHON FILE IS LOCATED
285             IN THE SAME DIRECTORY AS YOUR
286                         HISTORY FILE /!\
287             """
288
289         # Choosing a name for the output PostgreSQL table (optional ,

```

```
    default='OSM_history')
284 tablename = 'history'
285
286 # Importing history file content
287 history_importer(filename, tablename)
288
289
290 # Close communication with the database
291 cur.close()
292 conn.close()
```

NAIVE PYTHON SCRIPT FOR OSM HISTORY FILE IMPORT INTO POSTGRESQL

APPENDIX I

Listing I.1: history_to_pgsql.py Homemade Python script

```
34
35 import xml.etree.ElementTree as ET
36
37 import psycopg2
38
39
40
41
42
43
44 ### History file import
45
46 """ /!\ PLEASE MAKE SURE THIS CURRENT PYTHON FILE IS LOCATED IN
47     THE SAME DIRECTORY AS YOUR
48             HISTORY FILE /!\
49 """
50 print("\n\n Parsing XML... ")
51 tree = ET.parse('otaniemi-history.osh') # feel free to change the
52     filename here to fit your needs
53
54
55
56 ### Connecting to the PostgreSQL Database
57 # feel free again to change the connection details to fit your
58     needs
59 conn = psycopg2.connect(database="uusimaa", user="postgres",
60     password="postgres", host="localhost", port="5432")
61 cur = conn.cursor() # to perform database operations
62
63
64
65
66 ### Code
67
68 root = tree.getroot()
69
70
71 history_records = []
72 id_primary_key = 0
73
74
75 print("\n\n Parsing history... ")
76
77 for child in tqdm(root):
78     id_primary_key += 1 # inserting a new id column so every
79         record has a unique identifier
80             # (one feature can have several versions
```

```

                                using the same feature ID!)
80    # every try/except below are meant to test if the related
     info exists!
81    try:
82        osm_id = child.attrib['id']
83    except KeyError:
84        osm_id = None
85    osm_type = child.tag
86    try:
87        version = child.attrib['version']
88    except KeyError:
89        version = None
90    try:
91        timestamp = child.attrib['timestamp']
92    except KeyError:
93        timestamp = None
94    try:
95        uid = child.attrib['uid']
96    except KeyError:
97        uid = None
98    try:
99        user = child.attrib['user']
100   except KeyError:
101       user = None
102   try:
103       changeset_id = child.attrib['changeset']
104   except KeyError:
105       changeset_id = None
106   try:
107       visible = child.attrib['visible']
108   except KeyError:
109       visible = None
110
111
112 if child.tag == 'node': # if the OSM object is a node
113     try:
114         lat = child.attrib['lat']
115     except KeyError:
116         lat = None
117     try:
118         lon = child.attrib['lon']
119     except KeyError:
120         lon = None
121 else:
122     lat = None
123     lon = None
124
125
126
127
128
```

```

129
130
131     # Preparing the additional tags for storing in hstore format
132     # (if they exist!)
133     # and Saving the node references for the ways, and the member
134     # (node, way, or relation) references for the relations
135     #      --> In fact, geometry information is only contained in
136     #      nodes. So to access the geometry
137     #          of a way, you need the geometries of the nodes that
138     #          are part of this way. The same goes
139     #          for the geometries of relations.
140
141     additional_tags = None
142     members_refs = None # references to nodes for ways and to
143                         # members for relations
144
145     # For initialising the storage list (see below)
146     stop1 = True
147     stop2 = True
148
149     # Checking if the OSM object has additional tags
150     try:
151         temp = child[0]
152     except IndexError:
153         pass
154     else:
155         for tag in child:
156
157             # if the tag is a referenced member of a way or a
158             # relation
159             if (child.tag == "way" and tag.tag == "nd") or (child
160                 .tag == "relation" and tag.tag == "member"):
161
162                 if stop1: # if it's the first referenced member
163                     # encountered, initialise an empty list for
164                     # storing
165                     members_refs = []
166                     stop1 = False
167                     # NB: this initialisation won't happen again,
168                     # since 'stop1' variable is not used anymore
169                     # from now on
170
171                     members_refs.append(int(tag.attrib['ref'])) #
172                         # adding the member id to the list
173
174                 # if the tag is just a common attribute
175                 else:
176
177                     if stop2: # if it's the first additional tag
178                         # encountered, initialise an empty list for
179                         # storing

```

```

166         additional_tags = []
167         stop2 = False
168         # NB: same remarks as for 'stop1' variable
169
170         key = tag.attrib['k']
171         value = tag.attrib['v']
172         additional_tags.append([key,value])
173
174     history_records.append((id_primary_key,
175                             osm_id,
176                             osm_type,
177                             version,
178                             timestamp,
179                             uid,
180                             user,
181                             changeset_id,
182                             visible,
183                             lat,
184                             lon,
185                             members_refs,
186                             additional_tags))
187
188
189
190 ### Inserting the output data into PostgreSQL
191
192 # the targeted table structure has to match with the columns
193 # pattern below, and also the table name must be the same!
194 # feel free to change if needed!
195 print("\n\n Inserting history into database... ")
196
197 for history_record in tqdm(history_records):
198     cur.execute("""INSERT INTO otaniemi_history VALUES (%s, %s, %
199                 %s, hstore(%s));""",
200                 history_record)
201
202 conn.commit() # Make the changes to the database persistent
203
204 cur.close()
205 conn.close()

```


OSM FEATURE-LEVEL MODIFICATIONS ANALYSER PYTHON SCRIPT

Listing J.1: typology_modif_encoding_copy.py Homemade Python script

```
1 #!/usr/bin/env python3
2 # -*- coding: utf-8 -*-
3 """
4 Created on Thu Jul 20 17:01:22 2023
5
6 @author: alexysren
7 """
8
9 """
10 -----
11
12         OSM FEATURE-LEVEL MODIFICATIONS ANALYSER
13             [BETA VERSION]
14
15             * custom tool by Alexys Ren      July 2023 *
16
17
18 -----
19 """
20 """
21 """
22 ['BETA VERSION']: Since some challenges were encountered in our
23 investigation, the initial goal was
24 rerouted toward statistics on geometry modification. However, all
25 the functions below were implemented so to be easily developed
26 further to achieve the initial goal.
27
28 """
29 [Quick start]: If you would like to use the code as it is, just
30     modify the 'MAIN' part at the end of the code to fit your needs
31     !
32 """
33 """
```

```

32
33 ##### Libraries import
34
35 from tqdm import tqdm # for displaying a progress bar on loops
36 import os
37
38 import psycopg2      # PostgreSQL driver for Python support
39 import pandas as pd
40 import numpy as np
41
42
43
44 ##### Functions
45
46
47 def osm_versions(osm_id):
48     """
49         Returns all the versions of a OSM feature based on its id
50
51     Parameters
52     -----
53     osm_id : int
54         OSM FEATURE ID.
55
56     Returns
57     -----
58     pandas DataFrame
59         SQL QUERY RESULT.
60
61     """
62     cur.execute("SELECT * FROM history WHERE osm_id = %s ORDER BY
63                 version;", (osm_id,))
64     return pd.DataFrame(cur.fetchall(), columns=[ 'id',
65                                         'osm_id',
66                                         'osm_type',
67                                         'version',
68                                         'timestamp',
69                                         'uid',
70                                         'user',
71                                         'changeset_id',
72                                         'visible',
73                                         'lat',
74                                         'lon',
75                                         'members_refs',
76                                         'other_tags'])
77
78
79
80 def geometry_references_resolver(member_versions,
81                                 feature_timestamp):

```

```

81     """
82     When resolving references for accessing ways and relations
83         geometries,
84     referenced members may also have several versions to be taken
85         into account.
86
87     This function takes all the versions of a referenced member
88         from a way or a relation,
89     and return a unique version from it.
90
91     Parameters
92     -----
93     member_versions : pandas DataFrame
94         CONTAINS ALL THE VERSIONS OF A REFERENCED MEMBER.
95     feature_timestamp : pandas Timestamp
96         TIMESTAMP FOR THE OSM FEATURE TO WHICH THE REFERENCE IS
97             LINKED.
98
99     Returns
100    -----
101    pandas Series
102        UNIQUE REFERENCED MEMBER FOR THE WAY/RELATION.
103
104
105     # Calculating time difference
106     member_versions['time_difference'] = member_versions['
107         timestamp'] - feature_timestamp
108
109     # Filtering versions older than 'feature_timestamp'
110     filtered_df = member_versions[member_versions['timestamp'] <=
111         feature_timestamp]
112
113     n, m = filtered_df.shape
114     if n == 0:
115         # tests revealed that some referenced members can
116         # be younger (in time) than the ways or relations they
117             are referenced to, which
118             # is theoretically impossible: DATA error (Error handling
119                 in 'geometry_modification')
120             return filtered_df
121
122
123     # Find the version with minimal time difference from the
124         filtered DataFrame
125     min_time_difference_index = filtered_df['time_difference'].idxmax() # relative time!
126
127     return member_versions.iloc[min_time_difference_index]
```

```

121
122
123
124
125
126
127 ## Topology of modifications
128
129 # Geometry modification
130
131 def geometry_modification(t1,t2,osm_type):
132     """
133         [RECURSIVE FUNCTION]
134         returns 1 if there has been a geometry modification, None if
135             not,
136         or 2 if there has been an error in the process
137
138     Parameters
139     -----
140         t1 : pandas Series
141             PREVIOUS VERSION.
142         t2 : pandas Series
143             CURRENT VERSION.
144         osm_type : string
145             FEATURE TYPE.
146
147     Returns
148     -----
149         int
150             see main description.
151
152     """
153
154     if osm_type == 'node':
155         if (t1.lat == t2.lat and t1.lon == t2.lon) or (np.isnan(
156             t2.lat) and np.isnan(t2.lon)):
157             # if coordinates are equals or if the geometry has
158                 been erased
159             return None
160         return 1
161
162
163     # So frow now on, only ways and relations are considered
164     t1_members_refs = t1['members_refs']
165     t2_members_refs = t2['members_refs']
166
167     if t2_members_refs == None:
168         # if the way/relation geometry has been erased
169         return None
170
171     if t1_members_refs == None:
172         # if the way/relation has no referenced members: DATA

```

```

    error
169     # -> it's theoretically impossible since way and relation
170     # existences relies
171     #       on their references!
172     return 2
173
174     n, m = len(t1_members_refs), len(t2_members_refs)
175
176     if n != m:
177         # if the numbers of nodes/members are not equal, there
178         # has been a modification!
179         return 1
180
181
182     t1_timestamp = t1['timestamp']
183     t2_timestamp = t2['timestamp']
184
185
186     if osm_type == 'way':
187         # Iterating over referenced nodes
188         for i in range(n):
189
190             # Referenced nodes from ways may also have several
191             # versions needed to be taken into account
192             t1_i_node_versions = osm_versions(t1_members_refs[i])
193             t2_i_node_versions = osm_versions(t2_members_refs[i])
194
195
196             # If the referenced node doesn't exist in our
197             # database
198             # -> cut during the geographical extraction of
199             # Otaniemi: EDGE EFFECT error
200             if t1_i_node_versions.shape[0] == 0 or
201                 t2_i_node_versions.shape[0] == 0:
202                 return 2
203
204
205             # Selecting the right version of the referenced node
206             # based on timestamps
207             #print("t=", t1_i_node_versions['timestamp'], "
208             #      timestamp=", t1_timestamp)
209             t1_i_node = geometry_references_resolver(
210                 t1_i_node_versions, t1_timestamp)
211             t2_i_node = geometry_references_resolver(
212                 t2_i_node_versions, t2_timestamp)
213
214             if t1_i_node.shape[0] == 0 or t2_i_node.shape[0] ==
215                 0:
216                 # if there has been an error (see '
217                 #      geometry_references_resolver' inline comments)
218                 return 2
219
220
221     # Comparing nodes geometry

```

```

207         if geometry_modification(t1_i_node, t2_i_node, 'node')  
208             != None:  
209                 return 1  
210     return None  
211  
212 else: # if osm_type == 'relation':  
213     # Iterating over referenced nodes/ways/relations  
214     for i in range(n):  
215  
216         # Referenced members from relations may also have  
217         # several versions needed to be taken into account  
218         t1_i_member_versions = osm_versions(t1_members_refs[i]  
219             ])  
220         t2_i_member_versions = osm_versions(t2_members_refs[i]  
221             ])  
222  
223         # If the referenced member doesn't exist in our  
224         # database  
225         # -> cut during the geographical extraction of  
226         # Otaniemi: EDGE EFFECT error  
227         if t1_i_member_versions.shape[0] == 0 or  
228             t2_i_member_versions.shape[0] == 0:  
229             return 2  
230  
231         # Selecting the right version of the referenced  
232         # member based on timestamps  
233         t1_i_member = geometry_references_resolver(  
234             t1_i_member_versions, t1_timestamp)  
235         t2_i_member = geometry_references_resolver(  
236             t2_i_member_versions, t2_timestamp)  
237  
238         if t1_i_member.shape[0] == 0 or t2_i_member.shape[0]  
239             == 0:  
240             # if there has been an error (see '  
241             # geometry_references_resolver' inline comments)  
242             return 2  
243  
244         # Comparing ways or relations geometry  
245  
246         t1_i_member_type = t1_i_member['osm_type']  
247         t2_i_member_type = t2_i_member['osm_type']  
248  
249             # if the feature types are different between the  
250             # compared members, there has been a modification  
251             !  
252         if t1_i_member_type != t2_i_member_type:  
253             return 1  
254  
255         if geometry_modification(t1_i_member, t2_i_member,

```

```

244         t1_i_member_type) != None:
245             return 1
246     return None
247
248 # Semantic modifications (not our case study, to be filled if
249 # needed!)
250
251 def tag_enrichment(t1,t2):
252     # same etc.
253     return None
254
255 def tag_suppression(t1,t2):
256     return None
257
258 def tag_modification(t1,t2):
259     return None
260
261
262 ## Encoding modification types into sequence
263 """
264
265 [IMPORTANT NOTE]: The below function was originally meant for
266 preparing the clustering analysis
267 based on encoded sequences of modifications with semantic ones
268 included. It is now used
269 for geometry-focused studies, but can be developed further on to
270 fit your needs.
271
272 """
273
274 def sequence_modifications(osm_id):
275     """
276     returns the encoded sequence of modifications for one OSM
277     feature based on its id
278
279     Parameters
280     -----
281     osm_id : int
282         OSM FEATURE ID.
283
284     Returns
285     -----
286     seq : int list
287         SEQUENCE OF MODIFICATIONS.
288
289 """
290
291
292 # Retrieving all the feature versions
293 versions = osm_versions(osm_id)

```

```

289     n, m = versions.shape # n = number of versions, m = number of
290     # fields
291
292     if n < 1:
293         # tests has shown that it may happen a OSM feature is not
294         # recorded in the history
295         return [2]
296
297
298     # Retrieving the feature type (is it a node, a way, or a
299     # relation?)
300     feature_type = versions['osm_type'][0]
301
302     # Iterating over versions (Pairwise similarity comparing)
303     for i in range(n-1):
304         geometryModification = geometry_modification(versions.
305             iloc[i], versions.iloc[i+1], feature_type)
306         if geometryModification != None:
307             seq.append(geometryModification)
308
309         # tag_enrichment =
310         # if tag_enrichment is not None:
311             # seq.append(tag_enrichment(versions.iloc[i], versions
312                 .iloc[i+1]))
313
314         # tag_suppression =
315         # tag_modification =
316
317         # .
318         # .
319         # . and so on...
320
321
322
323
324     def nb_geometry_modification(osm_id):
325         """
326             based on the current state of the 'sequence_modifications'
327             function,
328             nb_geometry_modification returns the number of geometry
329             modifications
330             underwent by an input OSM feature based on its ID.
331
332             Parameters
333             -----
334             osm_id : int

```

```

333     OSM FEATURE ID.
334
335     Returns
336     -----
337     int
338         NUMBER OF GEOMETRY MODIFICATIONS FOR THE INPUT FEATURE.
339
340     """
341     seq = sequence_modifications(osm_id)
342     if 2 in seq:
343         return None
344     return len(seq)
345
346
347
348 def get_geometry_modification():
349     """
350         saves locally and returns a pandas DataFrame containing one
351             line per feature with its
352             associated number of geometry modifications.
353
354         Returns
355         -----
356         df : pandas DataFrame
357             see main description.
358
359     """
360     # Retrieving OSM feature ids from '
361     # nls_buildings_multipolygons'
362     cur.execute("""SELECT * FROM
363     (SELECT CAST(osm_id AS bigint) FROM nls_buildings_multipolygons
364     UNION
365     SELECT CAST(osm_way_id AS bigint) FROM
366     nls_buildings_multipolygons) AS temp
367
368     WHERE temp.osm_id IS NOT NULL;""")
369
370     df = pd.DataFrame(cur.fetchall(), columns=['osm_id'])
371
372     nbGeometryModification = []
373
374     for osm_id in tqdm(df.osm_id):
375         nbGeometryModification.append(nb_geometry_modification(
376             osm_id))
377
378     df['nb_geometry_modification'] = nbGeometryModification
379
380     df.to_csv(
381         "nls_buildings_multipolygons_nb_geometry_modification.csv",
382         float_format=".3f", index_label="index", sep=" ")

```

```

378     return df
379
380
381 ## Typology of modifications: Statistics
382
383
384 def massive_contributions_extract(filename):
385     # /!\ TO BE MODIFIED:
386         # additional input: feature id list containing ids of NLS
387             # massive contributions
388         # filter the below dataframe to return only features
389             # coming from NLS massive contributions
390         # thus allowing proper statistics to answer our initial
391             # problem
392
393         """ [IN REALITY]: this function is just about loading a .csv
394             file, saved previously to avoid launching 'get_geometry_modification' at each code running
395                 (let's say you want to display the statistics
396                 again later...).
397
398             Indeed, 'get_geometry_modification' is the
399                 most time consuming in the whole process! (the
400                     whole history is scanned!) """
401
402     df = pd.read_csv(os.path.join(os.getcwd(), filename),
403                     delimiter=" ")
404     return df
405
406
407
408
409 def modifications_analyser(dataframe):
410     """
411     displays detailed statistics about modifications underwent by
412         the input features
413     (can be developed further on to include semantic
414         modifications)
415
416     Parameters
417     -----
418     dataframe : pandas DataFrame
419         OSM FEATURES WITH NUMBER OF GEOMETRY MODIFICATIONS.
420
421     Returns
422     -----
423     None.
424
425     """
426
427     nb_features = dataframe.shape[0]
428     nb_geometry_modification = dataframe['nb_geometry_modification']
429     arr = np.array(nb_geometry_modification)

```

```

417
418     nb_nan_values = nb_features - nb_geometry_modification.count()
419
420     # Creating boolean masks for non-zero and non-'NaN' values
421     non_zero_mask = arr != 0
422     non_nan_mask = ~np.isnan(arr)
423     # Counting the number of values that satisfy both conditions
424     nb_non_zero_and_non_nan_values = np.count_nonzero(np.logical_and(non_zero_mask, non_nan_mask))
425
426     # Creating a boolean mask for zero values
427     zero_mask = arr == 0
428     # Counting the number of zero values using the boolean mask
429     nb_zero_values = np.count_nonzero(zero_mask)
430
431
432     print("\n\n----- Geometry modification statistics\n-----\n")
433     print("Feature count:", nb_features)
434     print("Among the ", nb_features, "features")
435     print(" -", nb_non_zero_and_non_nan_values, "underwent geometric modifications (",
436           round(100 * nb_non_zero_and_non_nan_values/nb_features, 1), "% )")
437     print(" -", nb_zero_values, "have not been modified (",
438           round(100*nb_zero_values/nb_features, 1), "% )")
439     print(" -", nb_nan_values, "could not be analysed (",
440           round(100*nb_nan_values/nb_features, 1), "% )")
441     print("* Average number of geometry modifications per feature:",
442           round(nb_geometry_modification.mean(), 1))
443     print("* Max number of geometry modifications:",
444           nb_geometry_modification.max(), "\n")
445
446     ### Code execution      'MAIN'
447
448     if __name__ == '__main__':
449
450         # # Connecting to PostgreSQL database (feel free to change
451             # the details to fit your needs)
452         # conn = psycopg2.connect(database="uusimaa", user="postgres",
453             # password="postgres", host="localhost", port="5432")
454
455         # # Opening a cursor to perform database operations
456         # cur = conn.cursor()
457
458
459         """ /!\ This code assumes the existence of a PostgreSQL table
460             containing your OSM history data, which also need to
461             respect a particular table structure.

```

```
455     If you haven't already done so, please import your
456         history data using the 'history_to_pgsql_enhanced.
457             py' Python script,
458     then ensure that the SQL query in the 'cur.execute',
459         method of 'osm_versions' function matches the name
460             of your history table! /!\'
461
462     N.B.: Modifying the SQL query in the code was a coding
463         mistake,
464             as I should have asked the user for a tablename
465                 input in this 'MAIN' directly to make this easier
466     """
467
468     # # Saving the number of geometry modifications for each
469         feature belonging to 'nls_buildings_multipolygons' table on
470             our computer
471     # get_geometry_modification()
472
473     # Importing the saved file (from 'get_geometry_modification')
474     nls_buildings_multipolygons = massive_contributions_extract("nls_buildings_multipolygons_nb_geometry_modification.csv")
475
476     # Displaying statistics on geometry modification
477     modifications_analyser(nls_buildings_multipolygons)
478
479     # # Closing communication with the database
480     # cur.close()
481     # conn.close()
```

PYTHON SCRIPT FOR RETRIEVING OSM IDS FROM FEATURE WHICH UNDERWENT GEOMETRIC MODIFICATIONS

APPENDIX **K**

Listing K.1: `get_geom_modif_ids.py` Homemade Python script

```
1 #!/usr/bin/env python3
2 # -*- coding: utf-8 -*-
3 """
4 Created on Tue Aug  8 10:52:14 2023
5
6 @author: alexysren
7 """
8
9 ### Libraries import
10
11 import os
12
13 import pandas as pd
14 import numpy as np
15
16
17 ### Code
18
19
20 def get_geom_modif_ids(filename):
21     """
22         create a new .csv file containing the OSM feature ids of
23             those that have
24             underwent geometry modification from the input file
25
26     Parameters
27     -----
28     filename : string
29             CSV filename with its extension.
30
31     Returns
32     -----
33     None.
34
35     """
```

92 Python script for retrieving OSM ids from feature which underwent geometric modifications

```
36     df = pd.read_csv(os.path.join(os.getcwd(), filename),
37                       delimiter=" ")
38
39     osm_id = df['osm_id']
40     nb_geometry_modification = df['nb_geometry_modification']
41     arr = np.array(nb_geometry_modification)
42
43     # Creating boolean masks for non-zero and non-'NaN' values
44     non_zero_mask = arr != 0
45     non_nan_mask = ~np.isnan(arr)
46     # Retrieving index of values that satisfy both conditions
47     non_zero_and_non_nan_values = np.nonzero(np.logical_and(
48         non_zero_mask, non_nan_mask))
49
50     osm_ids = pd.DataFrame(np.take(osm_id,
51                                    non_zero_and_non_nan_values[0]), columns=["osm_id"])
52
53     output_filename = filename[:-4] + "_geom_modif_ids.csv"
54     osm_ids.to_csv(output_filename, float_format=".3f",
55                     index_label="index", sep=" ")
56
57     print("\n\nThe result has been successfully saved in '" + os.
58           getcwd() + "/" + output_filename + "'!")
59
60
61     #### Main
62
63     if __name__ == '__main__':
64         filename = 'nls_buildings_multipolygons_nb_geometry_modification.csv'
65         get_geom_modif_ids(filename)
```