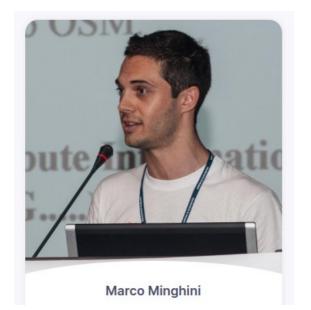
# Using binary data serialization for data storage and sharing

Peter Mooney (IE)
Marco Minghini (IT)









Department of Computer Science, Maynooth University, Co. Kildare, Ireland - peter.mooney@mu.ie

European Commission, Joint Research Centre (JRC), Ispra, Italy marco.minghini@ec.europa.eu



### The practicalities of (geo)data exchange

- Server
  - **Recieves** request
  - Prepares response
  - Send response (JSON, XML, csv, SHP, etc)
- Server
  - Preparation of complex responses (query times)
  - High traffic, networkbandwidth (uncontrollable)









- Client tools
  - Load response
  - Extract or Transform
  - Process: Visualise, analyse, integrate, etc.
- Client tools
  - Delayed response
  - Loading large responses
  - Long times: Extracting, transforming, processing large responses



RESPONSE

REQUEST



**CLIENT(S)** – with software

(Geo)data SERVER - with API

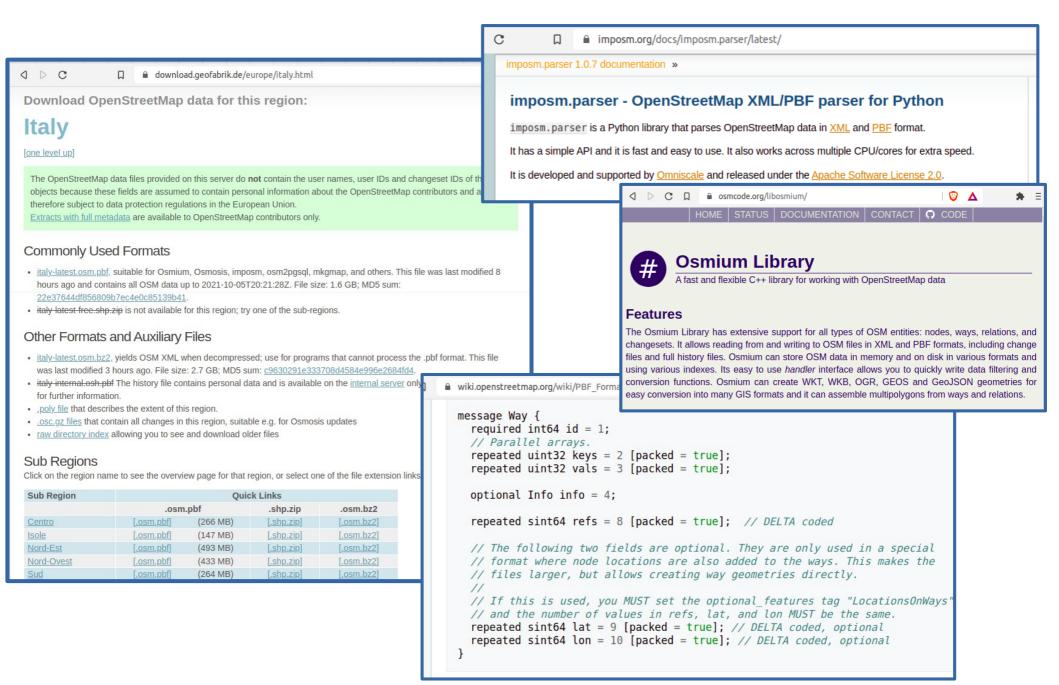
Research Question: Can binary data serialization provide a realistic alternative for geolocation data exchange in 'defacto' standards such as XML or JSON data formats.

- Not just a "time vrs space" analysis of binary vrs JSON/XML serialization approaches
- Focus on interoperability, usability, scalability, Open source, open approaches, ease of use ....
- Investigate conditions where binary serialization could replace or compliment the 'de-facto' standards (JSON, XML, and so on...)
- Google Protocol Buffers, Apache Avro

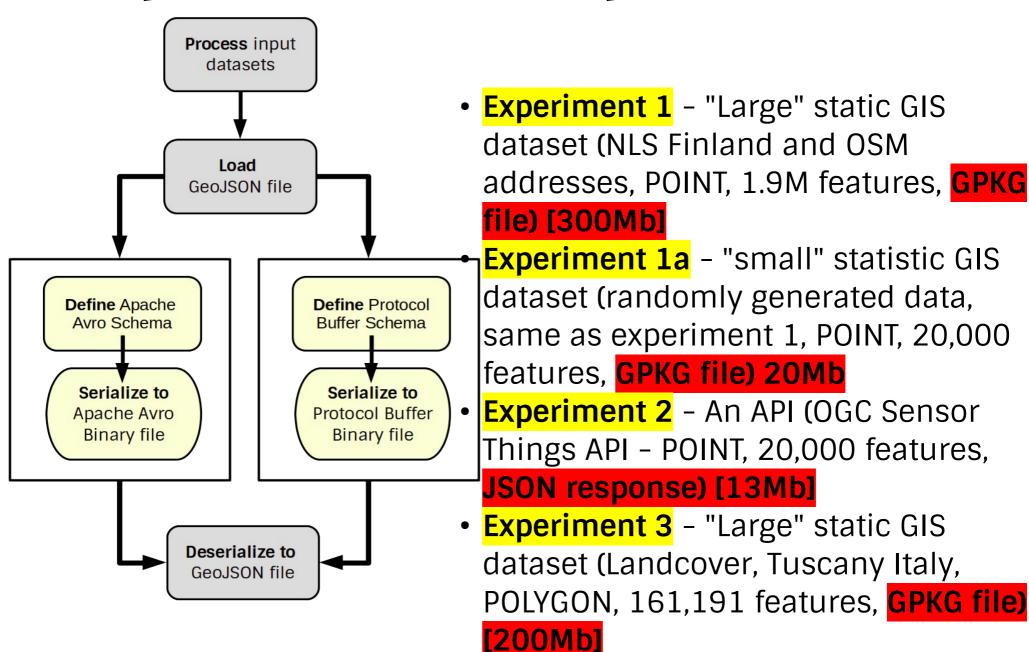




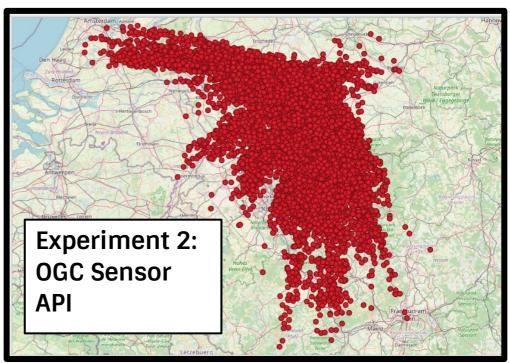
## Real world Example: OpenStreetMap – dissemination of data in PBF format

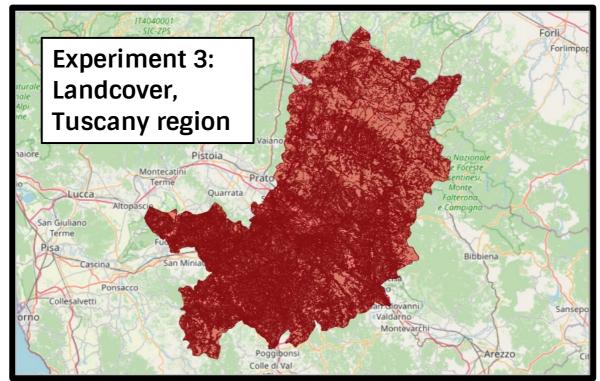


### **Experimental Setup - Workflow**



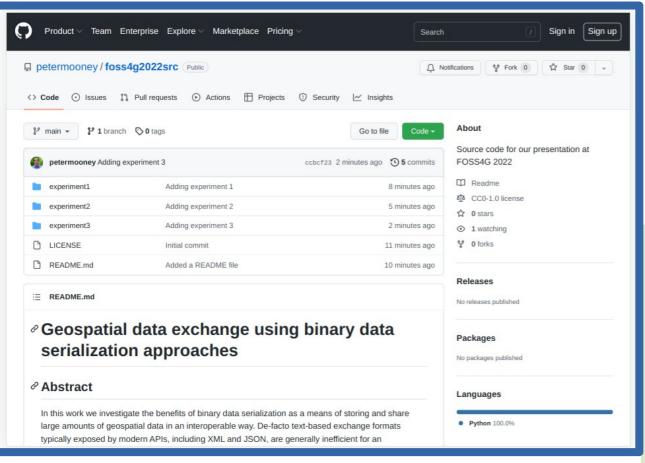






### **Experimental Setup - Software**

https://github.com/petermooney/foss4g2022src









#### Software setup

- Fully reproducible code (GitHub)
- · No "hacks"
- Use open source and widely supported Python libraries only
- Interoperable

### Creating Avro and Protobuf Schemas

#### **Original GPKG**

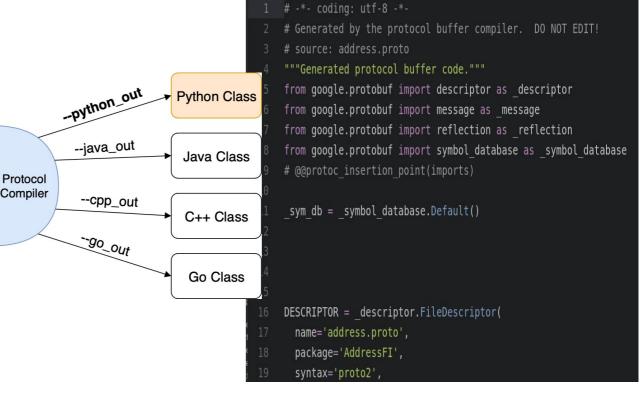
Layer Properties — original-dataset-experiment1 test-geopackage — Fields								
ld ▼	Name	Alias	Туре	Type name	Length	Precision	Comment	
123 0	fid		qlonglong	Integer64	0	0		
abc <b>1</b>	addr:housenumber		QString	String	0	0		
abc 2	addr:street		QString	String	0	0		
abc 3	addr:country		QString	String	0	0		
abc 4	addr:city		QString	String	0	0		
abc 5	source		QString	String	0	0		
abc 6	fullAddress		QString	String	0	0		
abc 7	addr:unit		QString	String	0	0		

```
package AddressFI;

package AddressFI;

message Address_prop {
    required string addrHousenumber = 1;
    required string addrCity = 3;
    required int32 fid = 4;
    required string source = 5;
    required string addrUnit = 6;
    required string fullAddress = 7;
    required string geometry = 8;

message Address {
    repeated Address_prop address = 1;
}
```



## Experiment 1 Results: Conflated addresses dataset

Process input datasets

Load GeoJSON file

Deserialize to GeoJSON file **Define Protocol** 

Buffer Schema

Serialize to

Protocol Buffer

Binary file

**Define** Apache

Avro Schema

Serialize to

Apache Avro

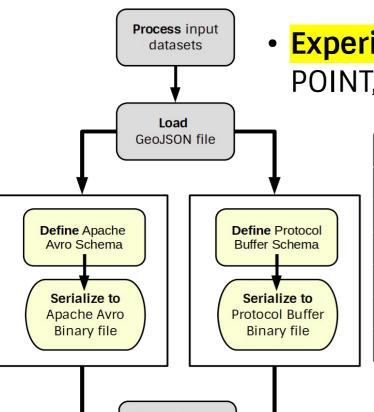
Binary file

 Experiment 1 - "Large" GIS dataset (NLS Finland and OSM addresses, POINT, 1.9M features,
 GPKG file) [300Mb]

<b>Processing steps from Figure 1</b>	Time (s)	Size (Kb)
Convert (GPKG $\rightarrow$ GeoJSON)	327, 11.3	614,859
Load (GeoJSON)	81, 3.2	614,859
Serialize (GeoJSON $\rightarrow$ Avro)	301, 3.4	228,102
Serialize (GeoJSON $\rightarrow$ PBF)	306, 2.9	235,821
Deserialize (PBF $\rightarrow$ GeoJSON)	378, 2.8	542,878
Deserialize (Avro → GeoJSON)	389, 3.1	546,616

"Long (de)serialization processing times. The difference in GeoJSON file sizes (614MB vrs 546MB) is due to coordinate precision representation in Python libraries. Binary files (Avro and PBF) are ~ 21% smaller"

## Experiment 2 Results: OGC Sensor Things API



Deserialize to

GeoJSON file

 Experiment 2 – An API (OGC Sensor Things API – POINT, 20,000 features, JSON response) [13Mb]

<b>Processing steps from Figure 1</b>	Time (s)	Size (Kb)
JSON Response download	n/a	12,926
Load (GeoJSON)	1.23, 0.07	11,539
Serialize (GeoJSON → Avro)	0.34, 0.04	7,001
Serialize (GeoJSON $\rightarrow$ PBF)	0.32, 0.04	7,109
Deserialize (PBF → GeoJSON)	1.14, 0.07	11,515
Deserialize (Avro $\rightarrow$ GeoJSON)	1.10, 0.03	11,554



"Very positive results. Binary files (Avro and PBF) ~ 46% smaller with serialization very fast. Deserialization is 3.5 times slower and this would need to be improved, perhaps with a more efficient implementation."

# Experiment 3 Results: Polygons - Tuscany

• Experiment 3 – "Large" static GIS dataset

(Landcover, Tuscany Italy, POLYGON, 161,191 features, GPKG file) [200Mb]

<b>Processing steps from Figure 1</b>	Time (s)	Size (Kb)
Convert (GPKG $\rightarrow$ GeoJSON)	128.7, 2.30	528,386
Reduce GeoJSON attributes	148.3, 2.2	320,472
Load (GeoJSON - reduced)	48.7, 0.31	320,472
Serialize (GeoJSON → Avro)	61.3, 0.18	264,555
Serialize (GeoJSON $\rightarrow$ PBF)	58.0, 0.15	264,937
Deserialize (PBF → GeoJSON)	80.4, 1.53	320,472
Deserialize (Avro $\rightarrow$ GeoJSON)	82.6, 1.93	320,472



**Define** Protocol

Buffer Schema

Serialize to

Protocol Buffer

Binary file

**Define** Apache

Avro Schema

Serialize to

Apache Avro

Binary file

Deserialize to GeoJSON file

"Many redundant attributes in original dataset which we removed. Original dataset expresses coordinates with 15 decimal places whereas Python GeoJSON allows 10. Avro and PBF ~ 17% smaller than the *reduced* GeoJSON file."

## Experimental Results – Technical considerations

- File sizes: PBF and Avro were, on average, at least
   20% smaller than the GeoJSON files
- Deserialization was always slower (futher investigation needed) than serialization
- Deserialization back to GeoJSON was performed as binary files lack basic query capabilities
- Data model considerations: Precision of geographic coordinates, inclusion of redundant attributes, and so on must be carefully considered in the source datasets.

### Experimental Results – Practical considerations

- Positive results but many overheads remain (nonautomated schema generation, specialist programmer knowledge, user community support, wider (geo) tool support)
- Binary data serialization could work very well for specific scenarios, services and applications (OSM is the leading example)
- It still remains a challenge to measure and understand "success" in regards to the possible replacement of existing 'de-facto' standards with binary data serialization

### Opportunities for future work

- Further computational experimentation on different types of (geo)data and (geo)services
- Making binary data approaches user-friendly like their non-binary data alternatives. This includes more integrated software tool support
- Schema definitions: automated semantic interoperabilty using linked geodata



### Publications related to this work



The International Archives of the Photogrammetry, Remote Sensing and Spatial Information Sciences, Volume XLVIII-4/W1-2022
Free and Open Source Software for Geospatial (FOSS4G) 2022 – Academic Track, 22–28 August 2022, Florence, Italy

#### GEOSPATIAL DATA EXCHANGE USING BINARY DATA SERIALIZATION APPROACHES

P. Mooney 01, M. Minghini 02

Department of Computer Science, Maynooth University, Co. Kildare, Ireland - peter.mooney@mu.ie
European Commission, Joint Research Centre (JRC), Ispra, Italy - marco.minghini@ec.europa.eu

Commission IV, WG IV/4

KEY WORDS: Binary data serialization, interoperability, geospatial, GeoJSON, Protocol Buffers, Apache Avro.

#### ARSTRACT:

In this paper we investigate the benefits of binary data serialization as a means of storing and sharing large amounts of geospatial data in an interoperable way. De-facto text-based exchange encodings typically exposed by modern Application Programming Interfaces (APIS), including exkensible Markup Language (KML) and JavaScript Object Notation (ISON), are generally inefficient for an increasingly higher number of applications due to their inflated volumes of data, low speed and the high computational cost for parsing and processing. In this work we consider comparisons of JSON/Geospatial JSON (GeolSON) and two popular binary data encodings (Protocol Buffers and Apache Avro) for storing and sharing geospatial data. Using a number of experiments, we illustrate the advantages and disadvantages of both approaches for common workflows that make use of geospatial data encodings such as GeoPackage and GeoJSON. The paper contributes a number of practical recommendations around the potential for binary data serialization for interoperable (geospatial) data storiage and sharing in the future.

#### 1. INTRODUCTION AND MOTIVATION

Data-driven innovation has seen recent advances enabled by a dynamic technological context characterised by the continuous influx of data, miniaturization and massive deployment of sensing technology, data-driven algorithms, and the Internet of Things (IoT) (Granell et al., 2022). Data-driven innovation is considered a key part of several policy actions worldwide. The recently published European strategy for data (European Commission, 2020) envisions Europe's digital future in the data economy through the establishment of dedicated data spaces. These would allow the efficient flow of data between actors and sectors so that data-driven innovation is translated into concrete benefits for the economy and society. To exploit Europe's potential in our data-driven society, technologies currently used for the management, exchange and consumption of data, including geospatial data, must be evaluated in terms of their suitability to efficiently scale and adapt to streams of larger data and datasets. The increasing number of users accessing data services through mobile devices is putting pressure on service providers to make larger volumes of data available efficiently to these particular users. For many years, encodings such as JavaScript Object Notation (JSON), Geospatial JSON (GeoJSON), Comma-Separated Values (CSV) and eXtensible Markup Language (XML) have been considered as the de facto interoperable standards for data serialisation. The majority of Application Programming Interfaces (APIs) available today facilitate data sharing and exchange using these encodings (Vaccari et al., 2020). Whilst such encodings, mainly JSON and XML, have an almost universal support and many other advantages (e.g. they are human and machine-readable as well as open and standards-based), they also have many limitations. These limitations are particularly pronounced when the volume of data is large, resulting in reduced computational performance when exchanging or managing large data volumes.

\* Corresponding author

As an alternative, binary data serialization approaches also allow for the interoperable exchange of large volumes of data (Vanura and Kriz, 2018). Binary data serialization is a process to transform data structures of an object into a binary stream that can be transmitted (and/or stored) and reconstructured later. Popular distributors of geospatial data have also begun mak-



perimental analysis for our three common GIS workflow scen-

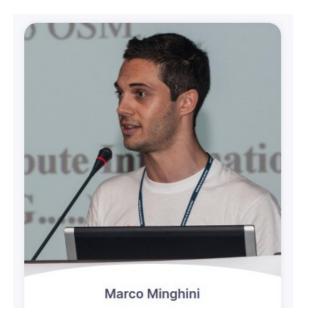
### Grazie a tutti – Thanks everyone!

## Peter Mooney (IE) Marco Minghini (IT)









Department of Computer Science, Maynooth University, Co. Kildare, Ireland - peter.mooney@mu.ie

European Commission, Joint Research Centre (JRC), Ispra, Italy marco.minghini@ec.europa.eu

