# Handling Large Spatial Data using PostgreSQL and DuckDB

Kshitij Raj Sharma

OSGeo

Nepal

# Contents

1. Introduction
2. Postgresql
3. DuckDB
4. Handling large spatial data
5. Experiments
6. DuckDB Usage
7. Walkthrough & Demo
8. Contact

Nepal

# PostgreSQL

- Most popular open source database
- Reliable
- Feature Rich
- Powerful indexes
- Large dev support group
- Production ready database support in multiple commercial cloud service provider
- Rich in spatial extensions like postgis , pg_raster

Nepal

# DuckDB

## DuckDB at a glance

### Simple
DuckDB is easy to install and deploy. It has zero external dependencies and runs in-process in its host application or as a single binary.

Read more →

### Portable
DuckDB runs on Linux, macOS, Windows, and all popular hardware architectures. It has idiomatic client APIs for major programming languages.

Read more →

### Feature-rich
DuckDB offers a rich SQL dialect. It can read and write file formats such as CSV, Parquet, and JSON, to and from the local file system and remote endpoints such as S3 buckets.

Read more →

### Fast
DuckDB runs analytical queries at blazing speed thanks to its columnar engine, which supports parallel execution and can process larger-than-memory workloads.

Read more →

### Extensible
DuckDB is extensible by third-party features such as new data types, functions, file formats and new SQL syntax.

Read more →

### Free
DuckDB and its core extensions are open-source under the permissive MIT License.
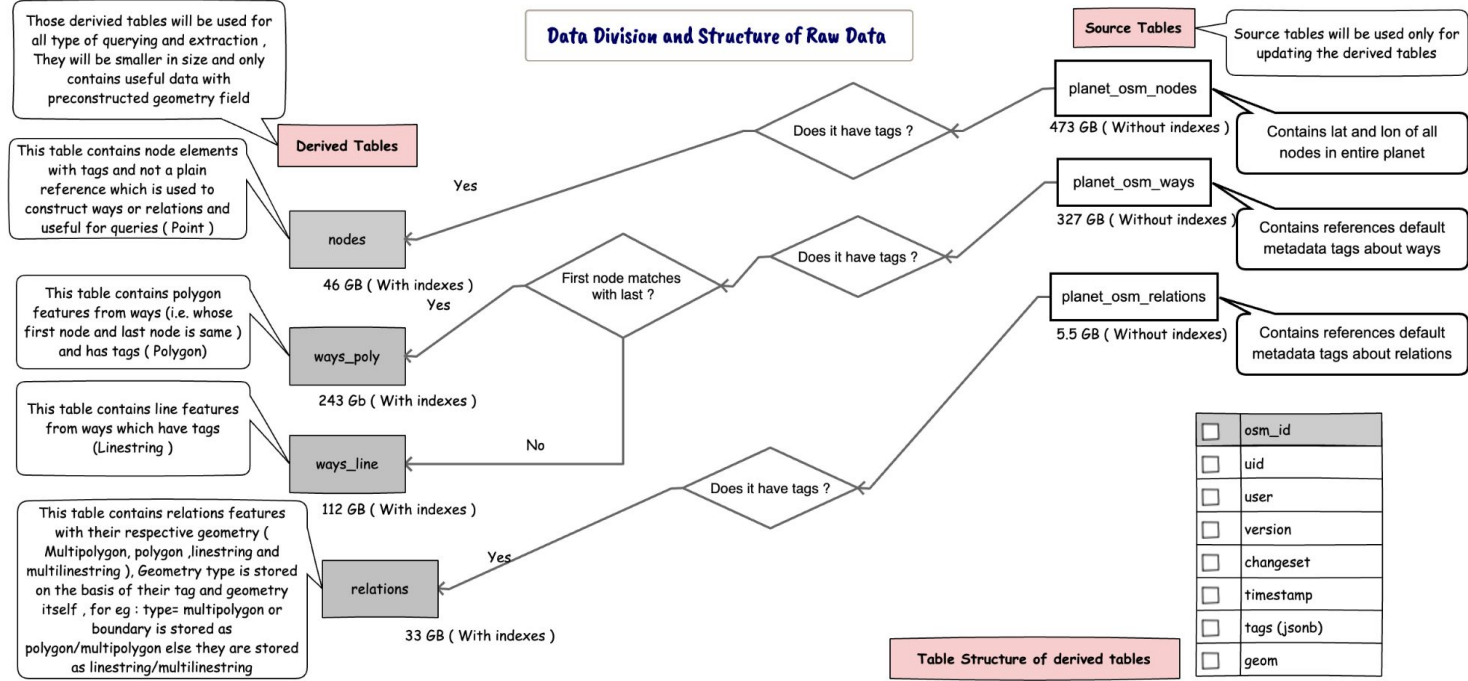
Read more →

- Spatial Support

# Handling Large spatial data

Tips and Tricks :

- Understanding and organizing data
- Processing raw data
- Choice of database based on need
- Schema
- Managing Indexes
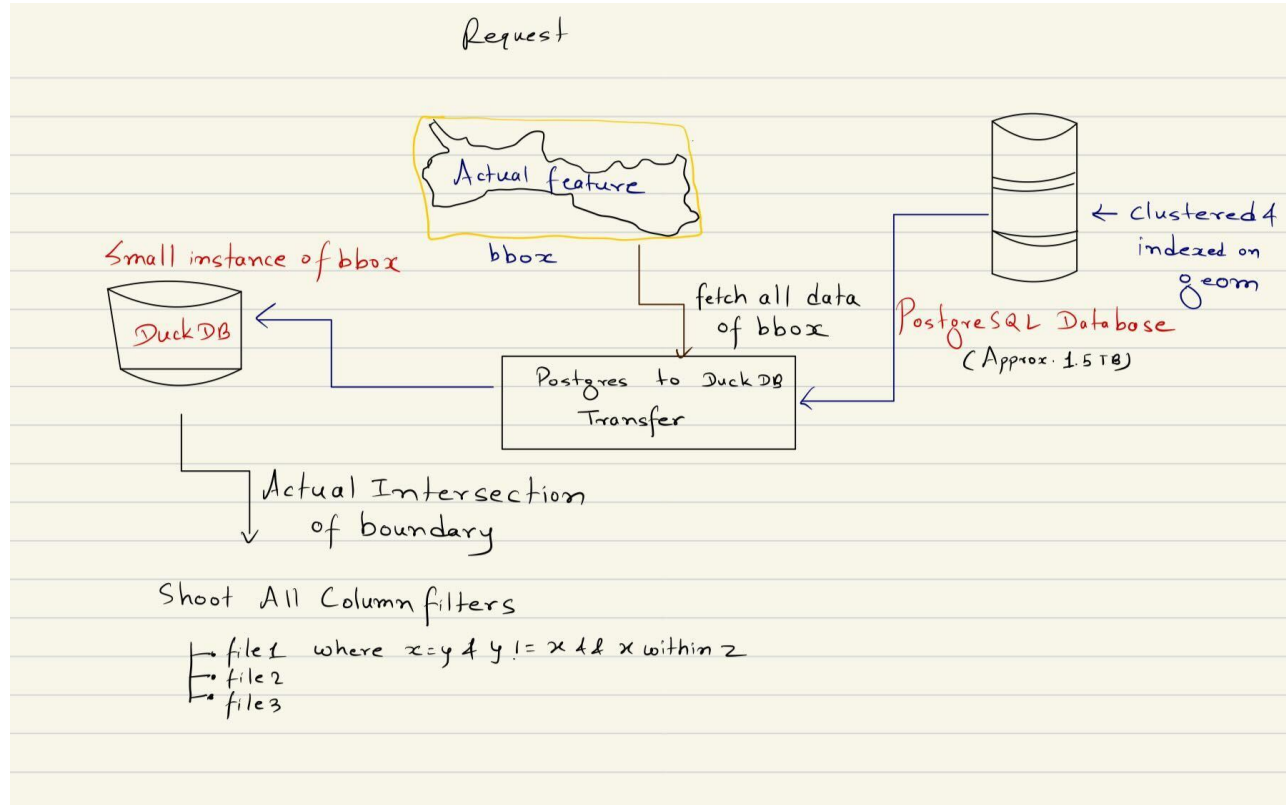- Query plan
- Scalability

Nepal

# Experiments



Data Division and Structure of Raw Data

Those derivied tables will be used for all type of querying and extraction , They will be smaller in size and only contains useful data with preconstructed geometry field

Derived Tables

This table contains node elements with tags and not a plain reference which is used to construct ways or relations and useful for queries ( Point )

nodes

46 GB ( With indexes )

This table contains polygon features from ways (i.e. whose first node and last node is same ) and has tags ( Polygon)

ways_poly

243 Gb ( With indexes )

This table contains line features from ways which have tags (Linestring )

ways_line

112 GB ( With indexes )

This table contains relations features with their respective geometry ( Multipolygon, polygon ,linestring and multilinestring ), Geometry type is stored on the basis of their tag and geometry itself , for eg : type= multipolygon or boundary is stored as polygon/multipolygon else they are stored as linestring/multilinestring

relations

33 GB ( With indexes )

Does it have tags ?

First node matches with last ?

Does it have tags ?

Does it have tags ?

Yes

Yes

No

Yes

Source Tables

Source tables will be used only for updating the derived tables

planet_osm_nodes

473 GB ( Without indexes )

Contains lat and lon of all nodes in entire planet

planet_osm_ways

327 GB ( Without indexes )

Contains references default metadata tags about ways

planet_osm_relations

5.5 GB ( Without indexes)

Contains references default metadata tags about relations

Table Structure of derived tables

| | |
|---|---|
| ☐ | osm_id |
| ☐ | uid |
| ☐ | user |
| ☐ | version |
| ☐ | changeset |
| ☐ | timestamp |
| ☐ | tags (jsonb) |
| ☐ | geom |

Nepal

# DuckDB Usage

- DuckDB is really useful when you want to do quick spatial operations in large files in your disk without having to setup whole database
- Faster columnar query as compared to postgresql without indexes
- Relatively new and doesn't contain lots of spatial datasets
- Currently focused on vector datasets
- Supports both spatial and non spatial query
- Has support of OGR formats enabling larger GIS format export support directly from the database
- Can be useful on google buildings , overture datasets / your custom csv datasets

Nepal

# PostgreSQL + DuckDB



Request

Actual feature
bbox

Small instance of bbox
DuckDB

← Clustered &
indexed on
geom

PostgreSQL Database
(Approx. 1.5 TB)

fetch all data
of bbox

Postgres to DuckDB Transfer

Actual Intersection
of boundary

Shoot All Column filters

- file 1  where x=y & y != x && x within z
- file 2
- file 3

# Walkthrough & Demo

# Contact Me :



Kshitij Raj Sharma

**LinkedIn / GitHub / Mastodon** :
@kshitijrajsharma

**Email** : skshitizraj@gmail.com

Nepal