

# **duckdb-geography: Global vector data in DuckDB**

Dewey Dunnington, Wherobots Inc.

# duckdb-geography

In a nutshell:

```
1 import duckdb
2 con = duckdb.connect()
3
4 con.install_extension("geography", repository="community")
5
6 con.load_extension("geography")
7 con.table_function("s2_data_cities").limit(3)
```

name varchar	population int32	geog geography
Vatican City	832	POINT (12.4533865 41.903282199999999)
San Marino	29000	POINT (12.441770200000002 43.936095799999999)
Vaduz	5342	POINT (9.51666947 47.1337238)

# Data types: Geography

A Geography is a Geometry whose edges follow a the shortest distance on the Earth rather than the shortest distance in Cartesian/planar space

If you want to represent something like Fiji, Russia, or Antarctica without splitting it into pieces, you need Geography!

# Data types: Geography

## Spherely array:

- Order is kept
- In memory data structures
- Requires explicit parallelization

## DuckDB Geography:

- Order is sometimes kept
- Serialized data structures (can be slower)
- Automatic parallelization

# Data types: S2 Cell

```
1 con.sql("""
2 SELECT s2_cellfromlonlat(7.64, 47.53).s2_cell_parent(level)
3 FROM (VALUES (5), (4), (3), (2), (1), (0)) levels(level)
4 """)
```

s2_cell_parent(s2_cellfromlonlat(7.64, 47.53), "level")	s2_cell
---	---------

	2/03302
--	---------

	2/0330
--	--------

	2/033
--	-------

	2/03
--	------

	2/0
--	-----

	2/
--	----



# Data types: S2 Cell

S2 Cells are represented by `uint64_t` and have an important property: **containment can be calculated using a simple comparison**

```
1 con.sql("""
2 SELECT c, s2_cell_range_min(c), s2_cell_range_max(c)
3 FROM (VALUES ('2/0330'::S2_CELL)) cells(c)
4 """)
```

c s2_cell	s2_cell_range_min(c) s2_cell	s2_cell_range_max(c) s2_cell
2/0330	2/03300000000000000000000000000000	2/03303333333333333333333333333333

(...because computers are much better at simple comparisons than spatial ones!)

# Data types: S2 Cell Center

```
1 (  
2     con  
3     .table_function("s2_data_cities")  
4     .select("geog")  
5     .limit(5)  
6 )
```

geog geography
POINT (12.4533865 41.903282199999999)
POINT (12.441770200000002 43.936095799999999)
POINT (9.51666947 47.1337238)
POINT (31.199997100000004 -26.466667500000003)
POINT (6.13000281 49.611660399999999)



# Data types: S2 Cell Center

```
1 (  
2     con  
3     .table_function("s2_data_cities")  
4     .select("geog::S2_CELL_CENTER AS cell")  
5     .limit(5)  
6 )
```

cell s2_cell_center
0/212113230003023131001102313200
0/212112131211123020010233101310
2/033031212023000232111020023022
0/331313212213231033112021020221
2/033022221321121102131101113231

# Data types: S2 Cell Union

```
1 con.sql("SELECT s2_data_country('Switzerland').s2_covering()")
```

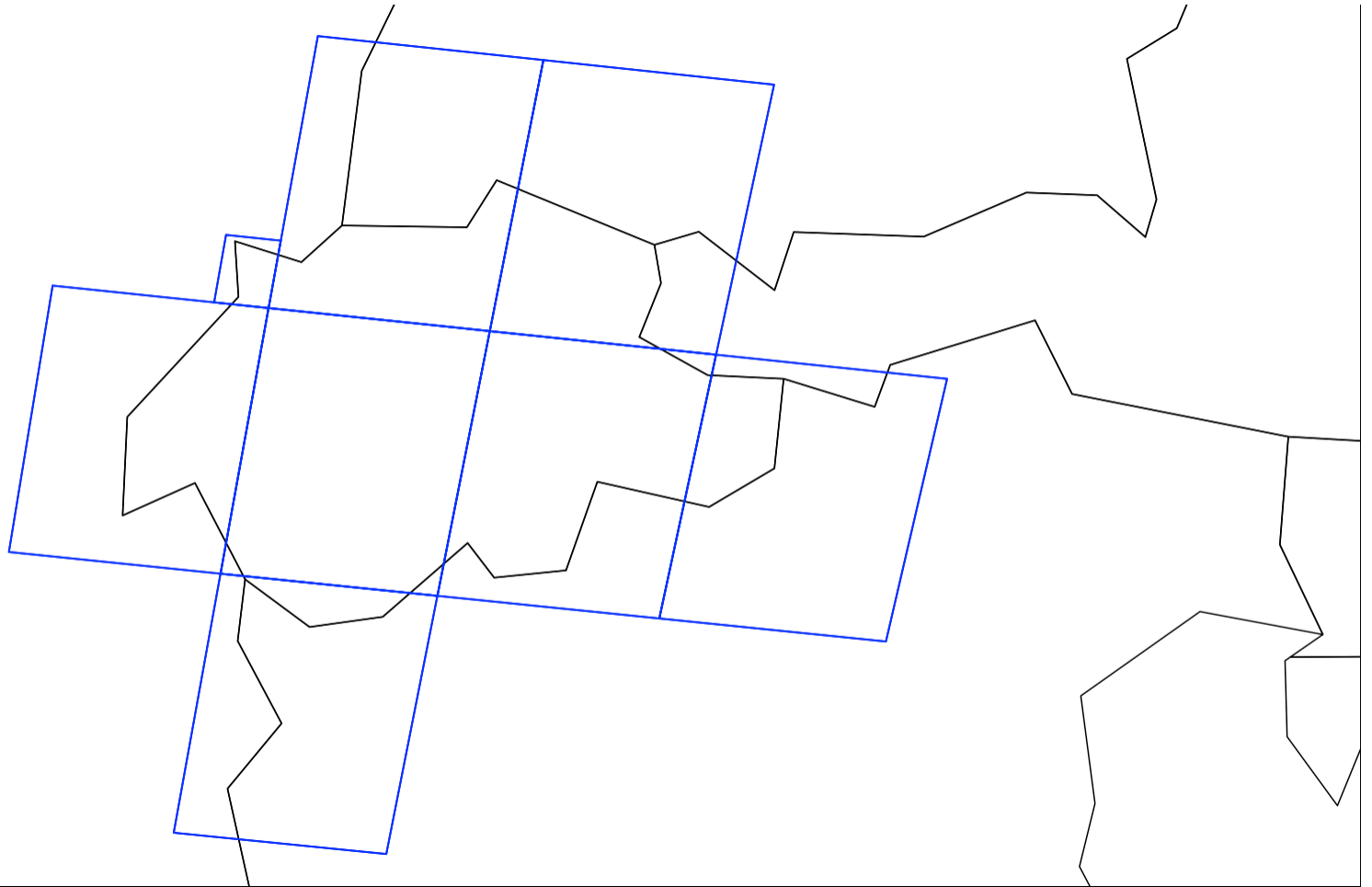
s2_covering(s2_data_country('Switzerland'))
[2/033001, 2/033002, 2/033010, 2/033012, 2/033013, 2/033020, 2/03302100, 2/033031]

s2\_cell\_union

[2/033001, 2/033002, 2/033010, 2/033012, 2/033013, 2/033020, 2/03302100, 2/033031]

# Data types: S2 Cell Union



```
1 con.sql("SELECT s2_data_country('Switzerland').s2_covering()")
```



# Practical example: using S2 cells to do a join

- <https://dewey.dunnington.ca/post/2024/wrangling-and-joining-130m-points-with-duckdb-the-open-source-spatial-stack/#duckdb-geography>
- <https://dewey.dunnington.ca/post/2024/partitioning-strategies-for-bigger-than-memory-spatial-data/#using-s2-cells-fixed-cell-level>

# duckdb-geography

- [dewey.dunnington.ca/slides/geopython2025](https://dewey.dunnington.ca/slides/geopython2025)
-  [paleolimbot/duckdb-geography](https://github.com/paleolimbot/duckdb-geography)
-  [duckdb.org/community\\_extensions](https://duckdb.org/community_extensions)

```
1 import duckdb
2 duckdb.install_extension("geography", repository="community")
```

```
1 INSTALL geography FROM community;
```