Our agenda

CARTO

# Why we are here today
# GIS and spatial thinking
# Welcome spatial indexes
# Workshop

## Data exploration
## Performance face-off
## Visualizing indexed data
## Spatial analysis
## Spatial Autocorrelation
# Outcomes and wrap-up

CART●

# Why we are here today

CART●

**Geospatial data is being produced at a massive scale**

# 1. Size
# 2. Velocity
# 3. Complexity

## What is the largest geospatial dataset?

*Here are some examples from BigQuery Public Data (from July 18, 2022)*

🛰️ NOAA Global Forecast System: 9B rows, 193TB
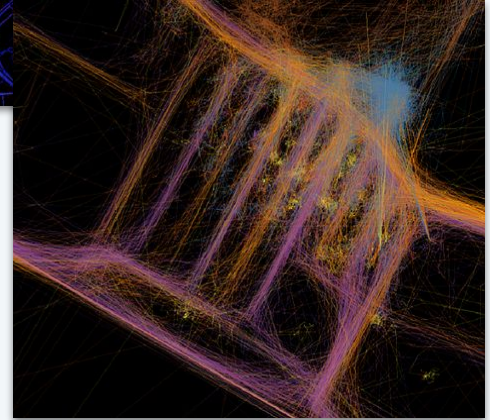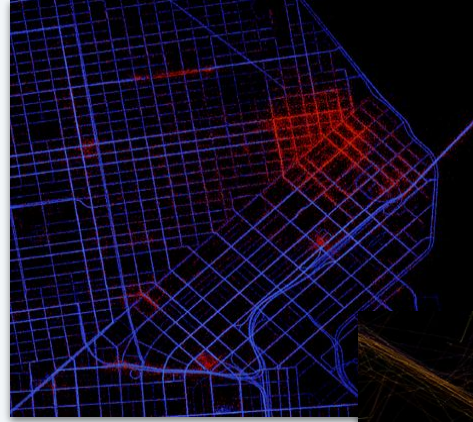
🖼️ OpenStreetMap: 993M rows, 327GB

🚕 NYC Taxis Dataset: ~115M rows, ~15 to 18GB (per year, per taxi type - i.e. Yellow, Green, Uber/Lyft)
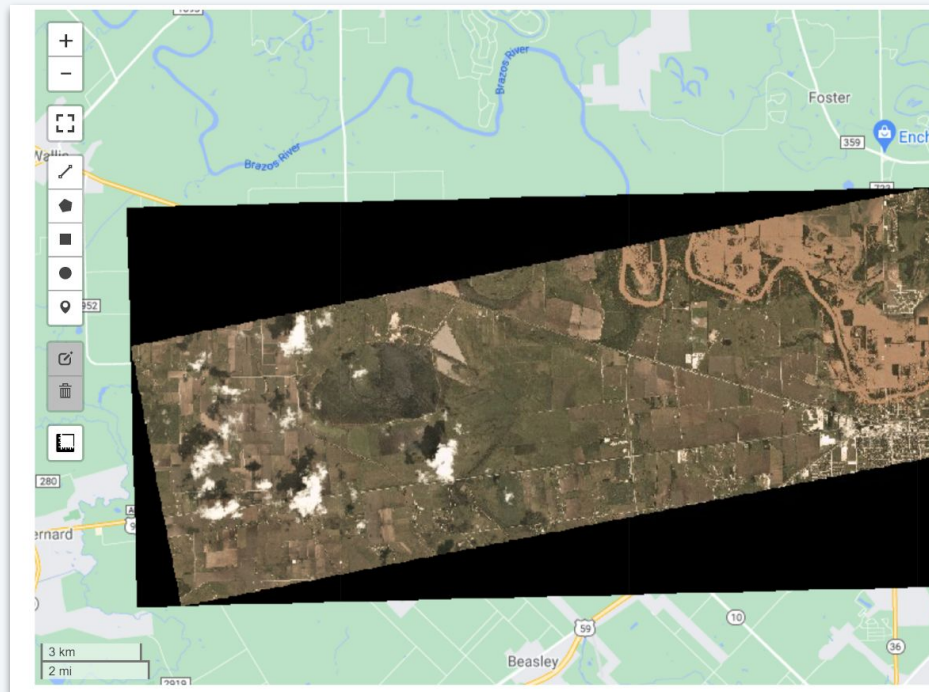
👫 WorldPop 1km Population Grids: 4.6B rows, 858GB

🌐 forrest.nyc     in mbforr

CART○

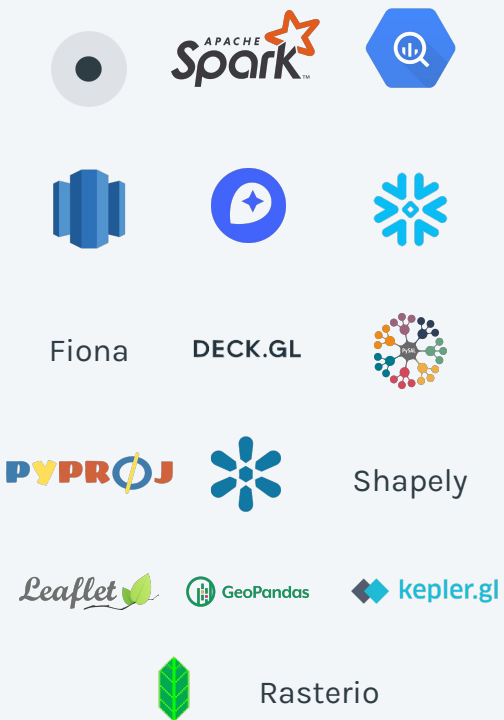# Larger volumes of vector (or vector ready) data

CART●

# Greater reliance on raster data

🇨🇱 **Data Council: Austin**

**2000s**

**2010s**

**2020s**

Fiona

DECK.GL

Shapely

PYPROJ

Leaflet

GeoPandas

kepler.gl

Rasterio

Leafmap

CART⬤

🇨🇱 Data Council: Austin

## Databases

## Desktop

## Web

kepler.gl

## Python

GeoPandas    Leafmap

## Developer

DECK.GL

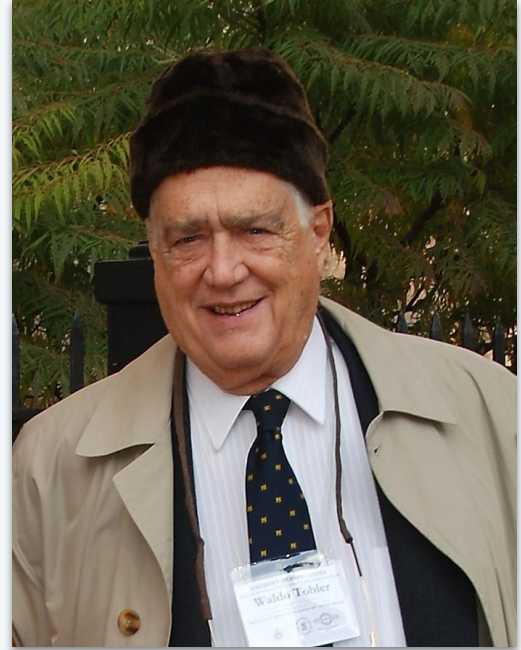## Tools

## Data +
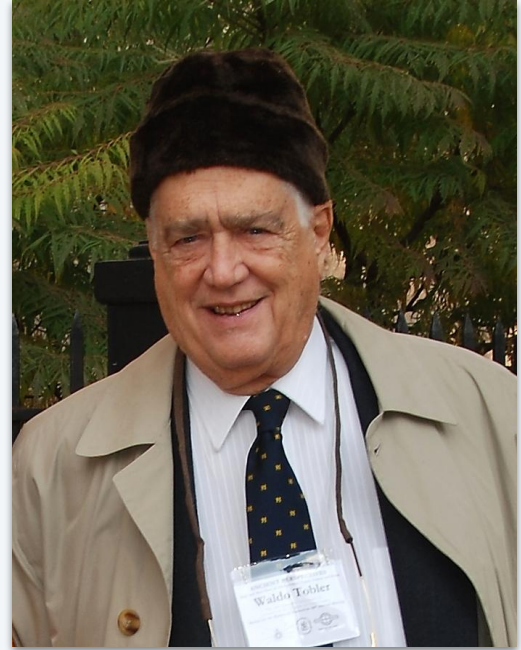
Shapely    Fiona    Rasterio

CARTO
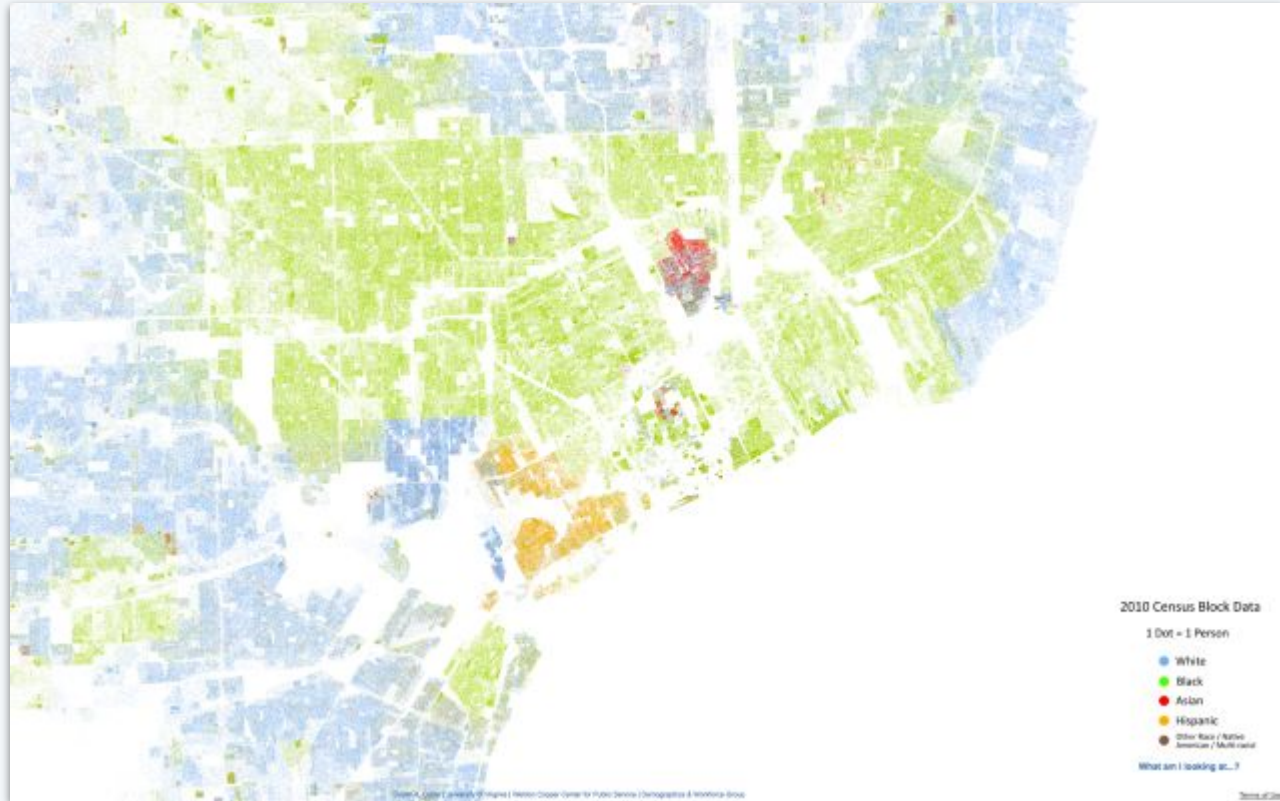
# GIS and spatial thinking

CART●

"Everything is related to everything else, but near things are more related than distant things."

CART●

"The phenomenon external to a geographic area of interest affects what goes on inside."

CART◐

2010 Census Block Data

1 Dot = 1 Person

● White
● Black
● Asian
● Hispanic
● Other Race / Native
   American / Multi-racial

What am I looking at...?

🇨🇱 **Data Council: Austin**

Buildings data

Vegetation data

Integrated data

CARTO

30m cells
9 billion rows
Multiple bands and
signal strengths

CART●

# 255GB to 28 GB
# 26:05 to 0:03 spatial join

CART●

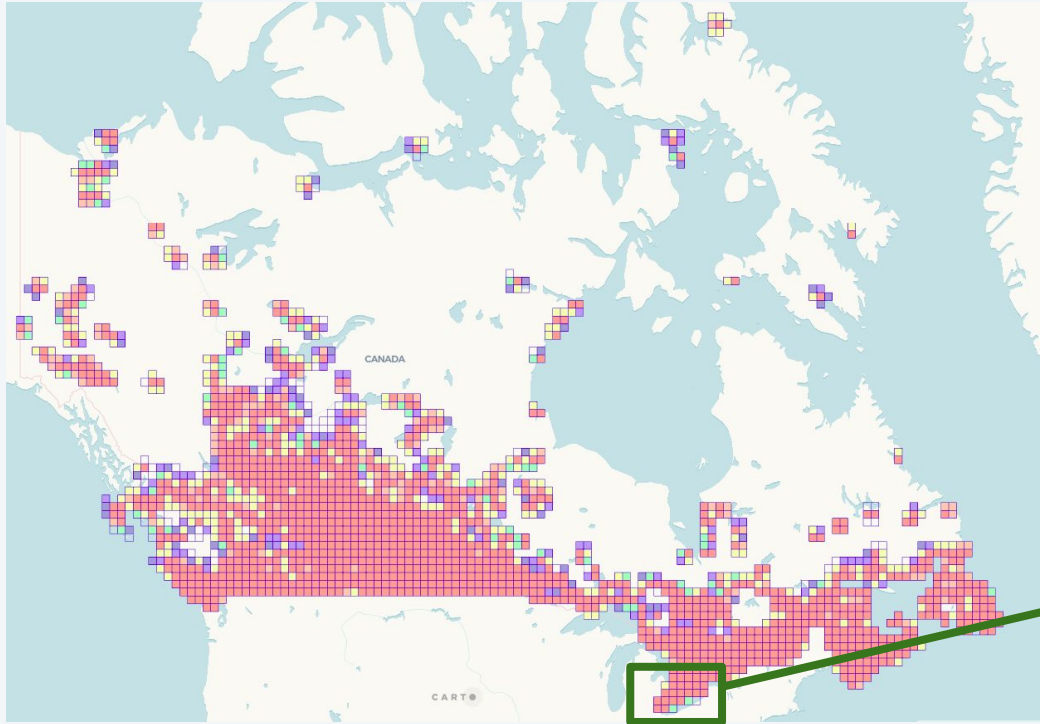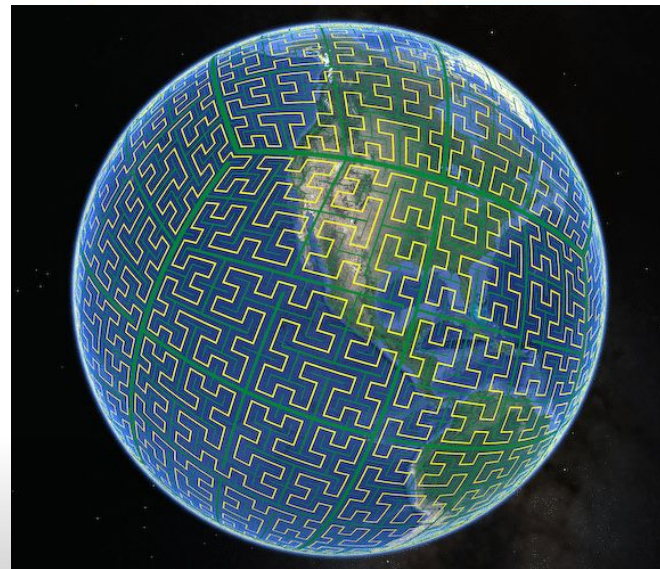# Welcome spatial indexes
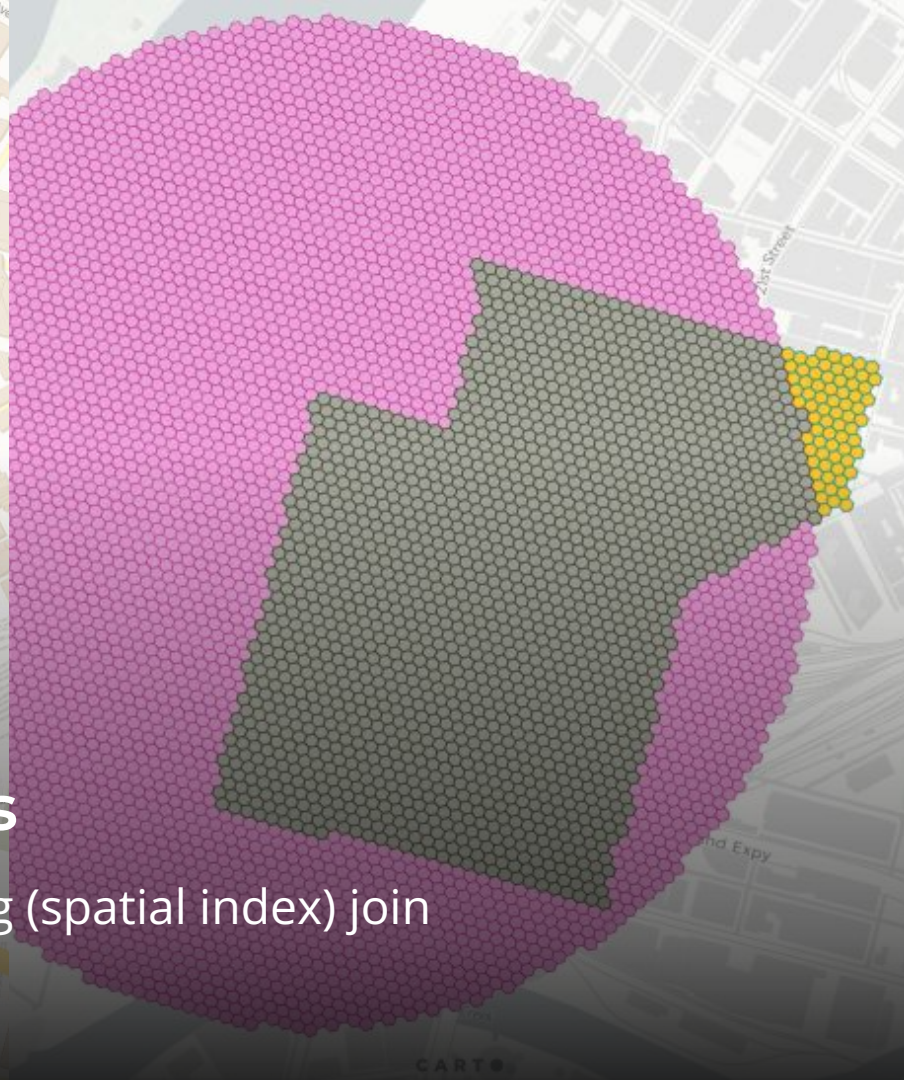
CART⦿

Quadkey (source)  Uber's H3 (source)  S2 (source)

# Geospatial Hierarchical Indexes

Different strategies to partition the space intro discrete grids

 1. Query performance
2. Storage
3. Visualization performance
4. Intuitive visualization
5. Neighbors and children

CART●

# Spatial joins to string joins

Join on spatial data compared to string (spatial index) join

# Geometries vs Spatial Indexes:

## What do they look like?



POLYGON((-96.196141 41.125515,
-96.195606 41.125514, -96.181864
41.125507, -96.177078 41.125474,
-96.167733 41.125456, -96.160565
41.125456, -96.154682 41.125429,
-96.151094 41.125414, -96.138848
41.125395, -96.138454 41.125394,
-96.138381 41.125394, -96.137158
41.125391, -96.130043 41.125377,
-96.1301...))*

*This represents about a **10th** of the geometric
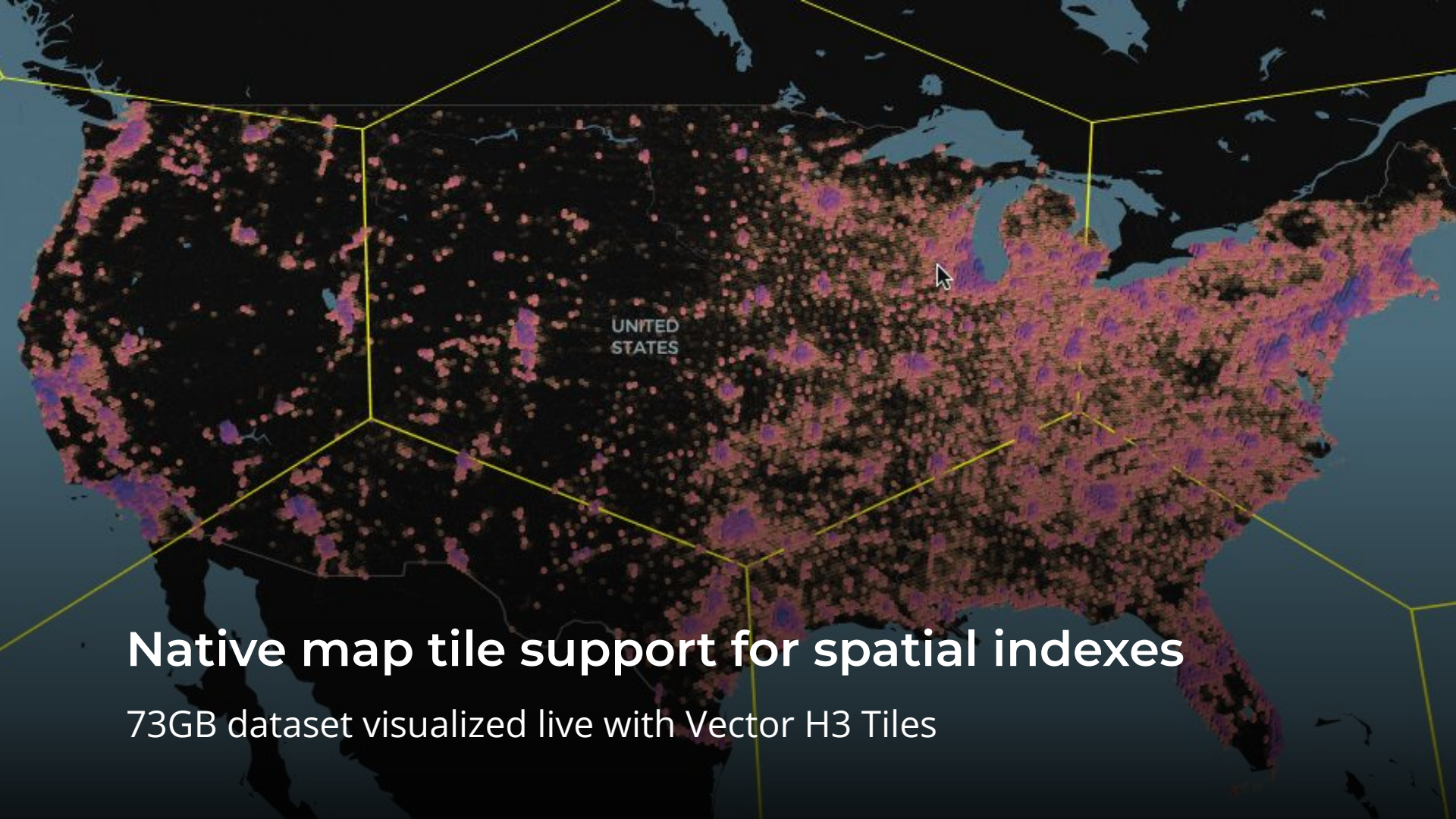description of a census tract, which we have
truncated for readability

8a2aa84ec307fff

**CARTO**

## Performance comparisons

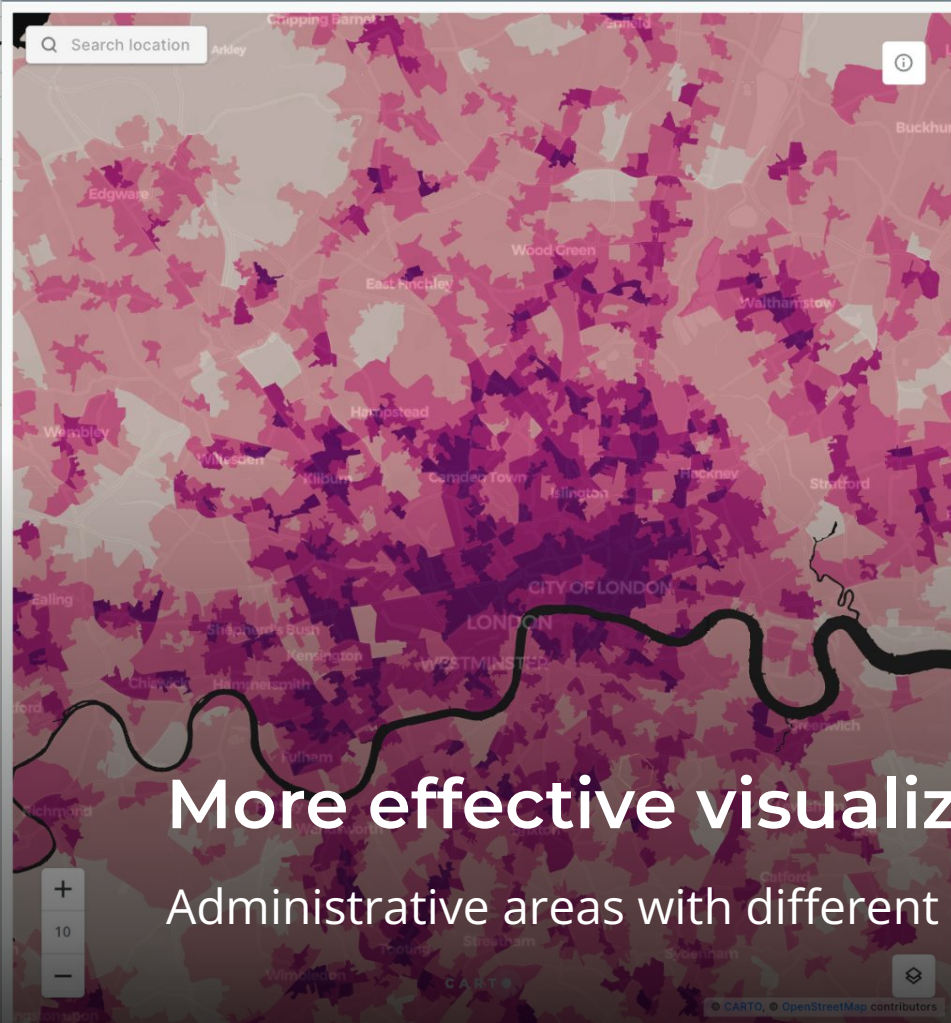| Example ETL use case | Geometries | Spatial Indexes | Gain of spatial index over geometries |
|---|---|---|---|
| Processing time | 12 days | 7 hours | **98% time saved** |
| Data transferred into the Database | 4 TB | 1.5 TB | **62% less data transferred** |
| RAM to process the largest file | 256 GB | 28 GB | **89% reduction in RAM** |
| Time to process a spatial join with population | 26 minutes | 3 seconds | **99% less time** |
| Time to generate a tileset | 23 minutes | 1.5 minutes | **94% time saved** |
| Population coverage | 15.48% | 15.48% | **0% coverage lost** |

**Estimated reduction in cloud data warehouse bill by 85%**
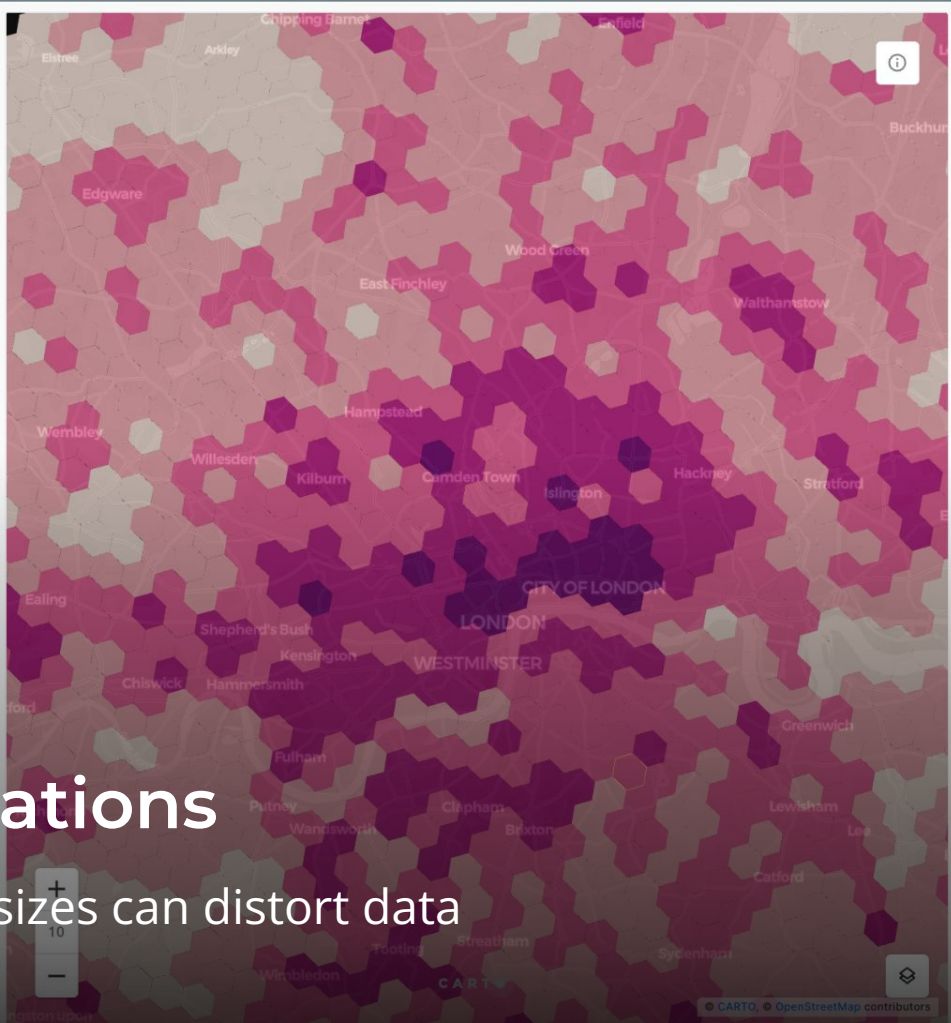
CART●

**Native map tile support for spatial indexes**

73GB dataset visualized live with Vector H3 Tiles

# More effective visualizations

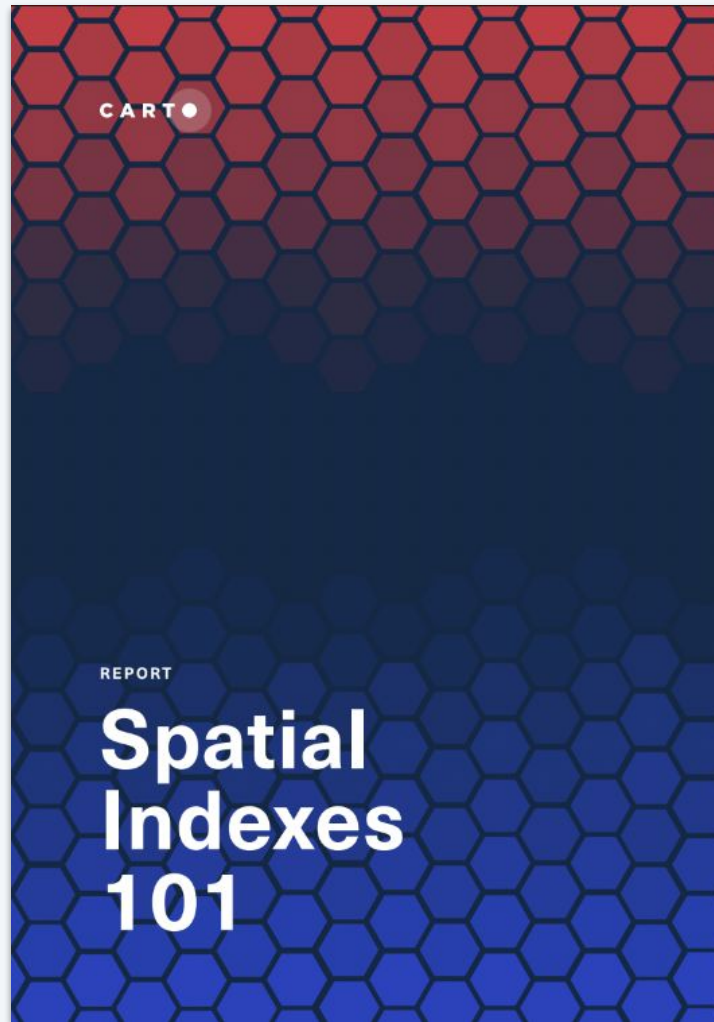Administrative areas with different sizes can distort data

# Analyze parents, children, and neighbors

Easy to scale up and down and analyze relationships

1. Loss of raw data
2. Precise spatial coverage
3. Original data quality and precision
4. Boundary effects

CART●

🇨🇱 **Data Council: Austin**

**Workshop**

CART○

1. Data exploration
2. Performance face-off
3. Visualizing indexed data
4. Enriching grids
5. Spatial analysis
6. Spatial Autocorrelation

CART●

[**Workshop docs here**](#)

CART●

**Q&A**

CART●

🇨🇱 Data Council: Austin

# Thank you!

**Matt Forrest**
**matt@carto.com**
**in mbforr**

CARTO