



# TERRA POPULUS

Integrated Data on Population and Environment

IEEE Big Data Conference, November 1, 2015



# Background

- National Science Foundation: Cyber Infrastructure Grant
- Data Net Initiative
- Enabling research, learning, and policy analysis by providing integrated spatiotemporal data describing people and their environment.



# Big Heterogeneous Data

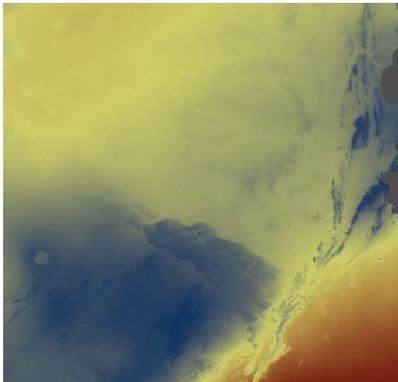
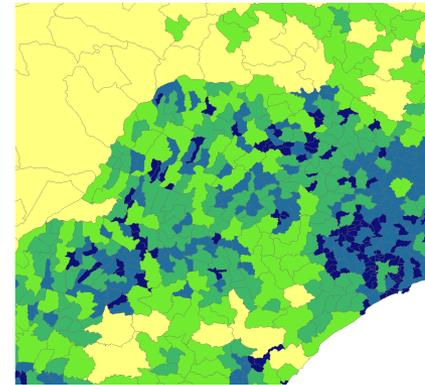


# TerraPop Data Formats



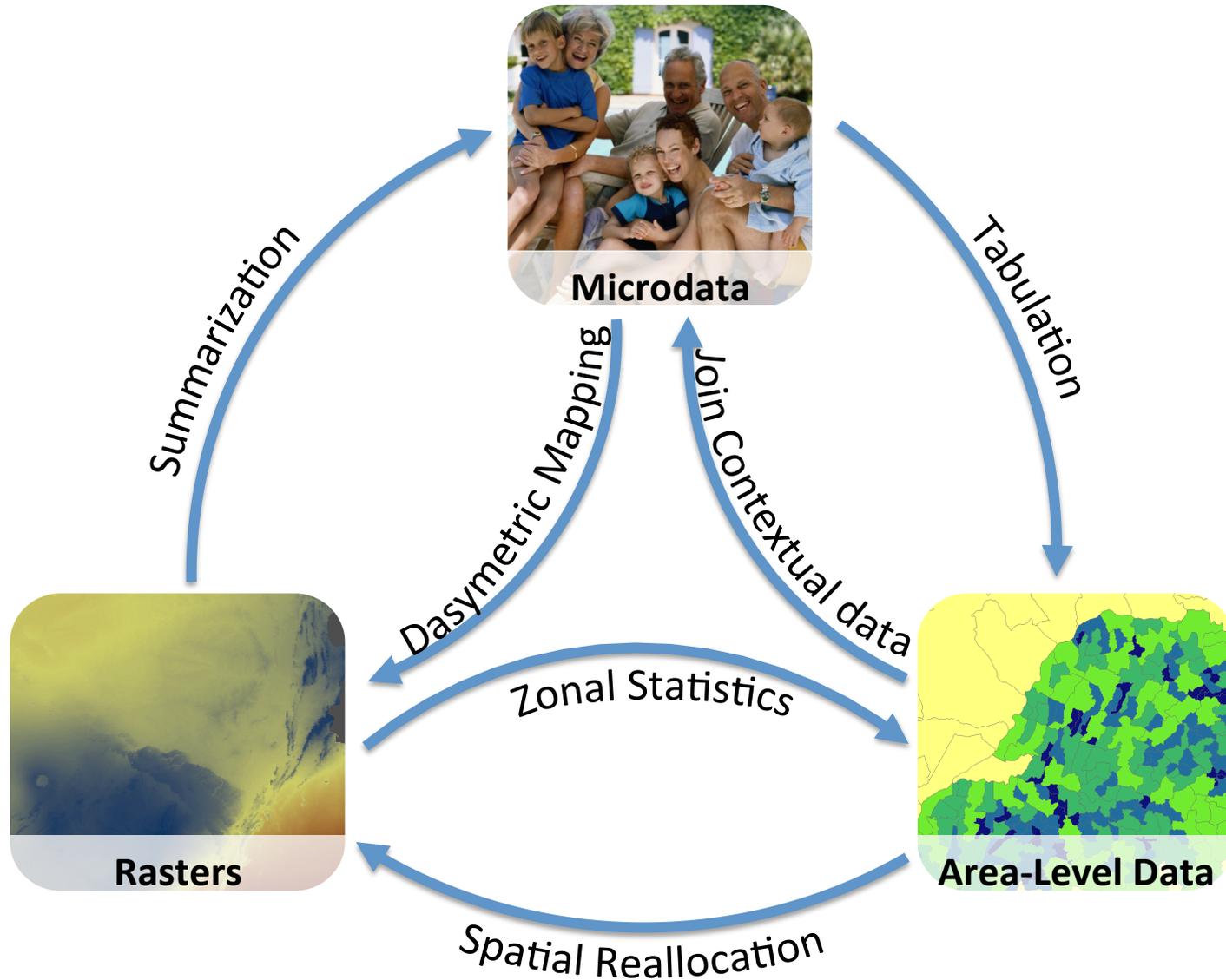
**Microdata:**  
Characteristics of individuals  
and households

**Area-level data:**  
Characteristics of places defined  
by boundaries



**Raster data:**  
Values tied to spatial  
coordinates





# Terra Populus



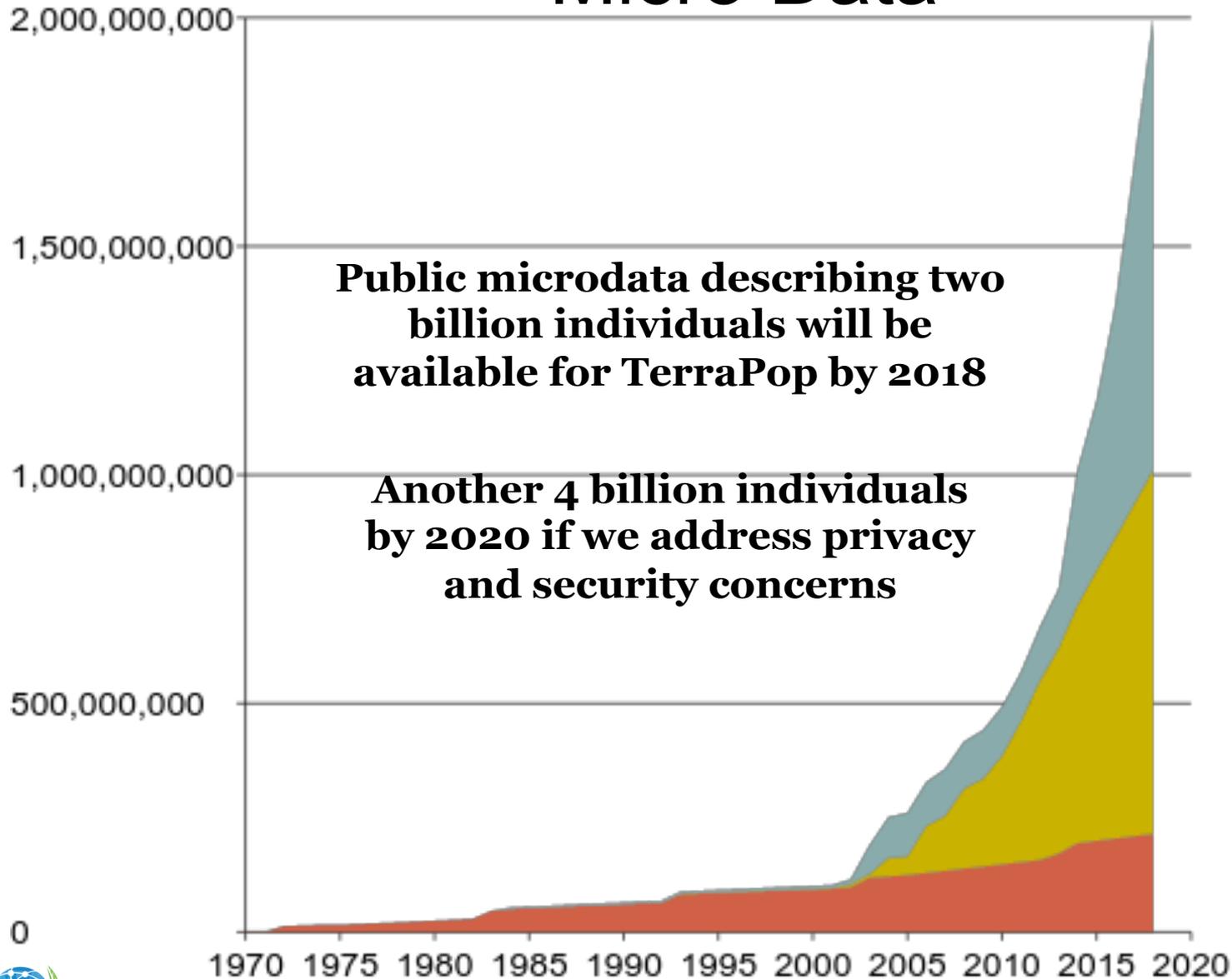
**DATA  
PARAGON  
TABULATOR  
GEOVISUALIZATION**

# Data

- **Microdata**
  - 70+ countries
  - 2<sup>nd</sup> and 3<sup>rd</sup> level geographies
- **Aggregate census data**
  - Historical data (48 countries)
  - Variables in addition to population by sex (65 countries)
- **Environmental data**
  - CRU monthly time series – precipitation & temperature
  - Elevation and derived characteristics
  - Soils



# Micro Data



**Microdata digitized from historical census manuscripts**

**Microdata from international statistical agencies**

**Microdata from U.S. Census Bureau**



# Areal Data

- US Population
  - 265 billion data points
- 95+ countries whose population < 500,000 people
- 87+ countries with microdata records
- 65 countries have data beyond total population by age and sex
- 54 countries have data at the second geographic level or finer resolution



# Raster Data

- **MODIS Land Data**

- Yearly land cover data derived from the MODIS Terra and Aqua satellites, available for 2001 – 2013
- 5 land cover classifications, 240 Gigabytes

- **Earth Science**

- Aster 30 Meter DEM resolution - 500 Gigabytes
- TAUDem derivatives: slope, solar radiance, wetness index will result in about 6-8 more Terabytes of data

- **Climate Datasets**

- NetCDF Format
- Climate Research Unit – 40 Gigabytes



# What Does TerraPop Do?



# TerraPop

- Access
- Exploration
- Interoperability
- Fusion
- Analysis

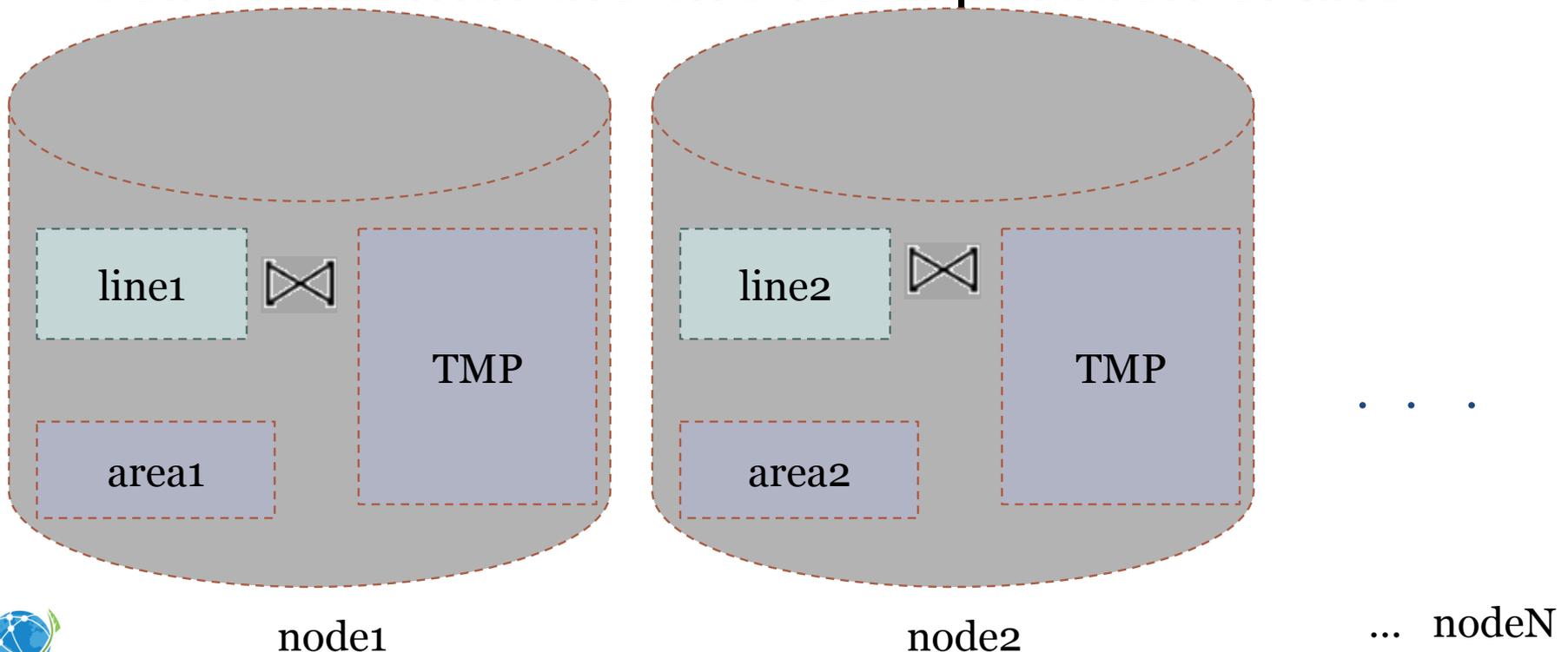


# Paragon



# Joins in Distributed Databases

- Create a temporary table TMP
- Reconstitute area on each node as TMP
- Join TMP with the two local partition of line



# Spatial Join Paragon

Query:

```
select a.gid , b.gid from edges_merge_ca_shall  
as a, arealm_merge_ca_shall  
as b where st_crosses(a.geom, b.geom) ;
```

PostgreSQL (standalone): 463 seconds

Stado-Spatial (2 nodes): 96 seconds



# Tabulator



# Tabulator

- Generates area-level data from microdata using geographic level codes
  - National, First Level (e.g. State), Second Level (e.g. County)
- Parquet on Apache Spark
  - High Compression Ratio
    - ✦ 8 Gigabytes gzip compressed
    - ✦ 3 Gigabytes parquet compressed
  - Columnar Storage (3,000+ columns)
  - 3 Nodes with each allocated 8 gigs of memory



# Query Performance

Number of datasets	Number of records	Time to aggregate by 1 column	Time to aggregate by 2 columns	Time to aggregate by 3 columns
1	12 million	5 seconds	9 seconds	10 seconds
10	25 million	7 seconds	10 seconds	18 seconds
20	82 million	10 seconds	12 seconds	27 seconds
56	128 million	7.5 seconds	20 seconds	30 seconds



# Visualization



# Landing Page

A Web Page

http://

**Select Variable** *You have not selected a variable*

Area Level Data | Raster Data

- Birthplace and Nativity
- Demographic
- Education
- Employment
- Household Amenities
- Household: Dwelling Characteristics
- Household Economic
- Household Utilities
- Urban

**Select Geographic Region** *You have not selected a geographic region*

Indicating Area Level Data Availability 1960 - 2013

Chart | Graph | Table

The image shows a web browser window titled "A Web Page" with a search bar containing "http://". Below the browser window is a landing page for a data visualization tool. The page is divided into two main sections: "Select Variable" and "Select Geographic Region". The "Select Variable" section has two tabs: "Area Level Data" (selected) and "Raster Data". Under "Area Level Data", there is a list of variables: Birthplace and Nativity, Demographic, Education, Employment, Household Amenities, Household: Dwelling Characteristics, Household Economic, Household Utilities, and Urban. The "Select Geographic Region" section features a world map where green areas represent data availability from 1960 to 2013. A legend below the map states "Indicating Area Level Data Availability 1960 - 2013". At the bottom of the page, there are three buttons: "Chart", "Graph", and "Table".



# Terra Populus Software Stack

## Web Application



## Data Processing



## Geospatial Data Repository

PostgreSQL



# Websites



**EXTRACT BUILDER**  
**[HTTPS://DATA.TERRAPOP.ORG/](https://data.terrapop.org/)**

**TERRACLIP**  
**[HTTPS://DATA.TERRAPOP.ORG/TERRACLIP](https://data.terrapop.org/terraclip)**

# Questions



?

