

# Crossing Boundaries Building Tools to Combine Public-Use (and Misaligned) Data

Matthew Graham

LEHD Program

Center for Economic Studies

U.S. Census Bureau

2018 FCSM Research and Policy Conference

# Disclaimer

- Any opinions and conclusions expressed herein are those of the authors and do not necessarily represent the views of the U.S. Census Bureau. All results have been reviewed to ensure that no confidential information is disclosed.
- Additionally, these opinions and conclusions are not representative of other data products or programs within the Census Bureau.

# Background

- LEHD Program
  - Longitudinal data on firms, workers, and jobs
  - >95% coverage of jobs
- Public-use data products (QWI, LODES, J2J) and dissemination tools (e.g. OnTheMap)
- Record of developing user-focused analysis tools with data partners in mind
  - More

# OnTheMap for Emergency Management

- Free public web application that gives emergency preparedness/disaster recovery analysts easy access to economic & demographic summaries for emergency events in near real time.
- Disseminates some statistics from 2010 Decennial, American Community Survey (ACS), and LEHD Origin Destination Employment Statistics (LODES).
- Automatically incorporates real-time weather and hazard event data from various authoritative federal sources into a single searchable archive and map display, for rapid access, visualization, and reporting.
- Building early LEHD tools (OnTheMap) resulted in lots of requests for analysis, particularly during emergency events or disasters.
- But we're not emergency management analysts...
  - So we built OTMEM to automate the mechanical work.

Search:

Filter ▾

U.S. Census Bureau data for disasters, natural hazards, and weather events. Click for more information

on Pacific Storms.

Events as of 02/27/2018

Wildfires

PLEASANT Fire

Zip 93514

Affected Population: 10

Federal Disaster Declarations

DR-4337

Monroe County, FL, Miami-Dade County, FL, Palm Beach County, FL and 64 other Counties

Affected Population: 18,801,310

DR-4223

Edwards County, TX, Duval County, TX, Harris County, TX and 110 other Counties

Affected Population: 17,594,564

EM-3396

Riverside County, CA, Los Angeles County, CA, San Diego County, CA and 2 other Counties

Affected Population: 16,350,775

DR-4332

Harris County, TX, Matagorda County, TX, Brazoria County, TX and 50 other Counties

Affected Population: 14,808,625

DR-4353

Los Angeles County, CA, San Diego County, CA, Santa Barbara County, CA and 1 other Counties

Affected Population: 14,161,134

EM-3387

Ware County, GA, Burke County, GA, Clinch County, GA and 156 other Counties

Affected Population: 9,687,771

DR-4338

Ware County, GA, Burke County, GA, Clinch County, GA and 156 other Counties

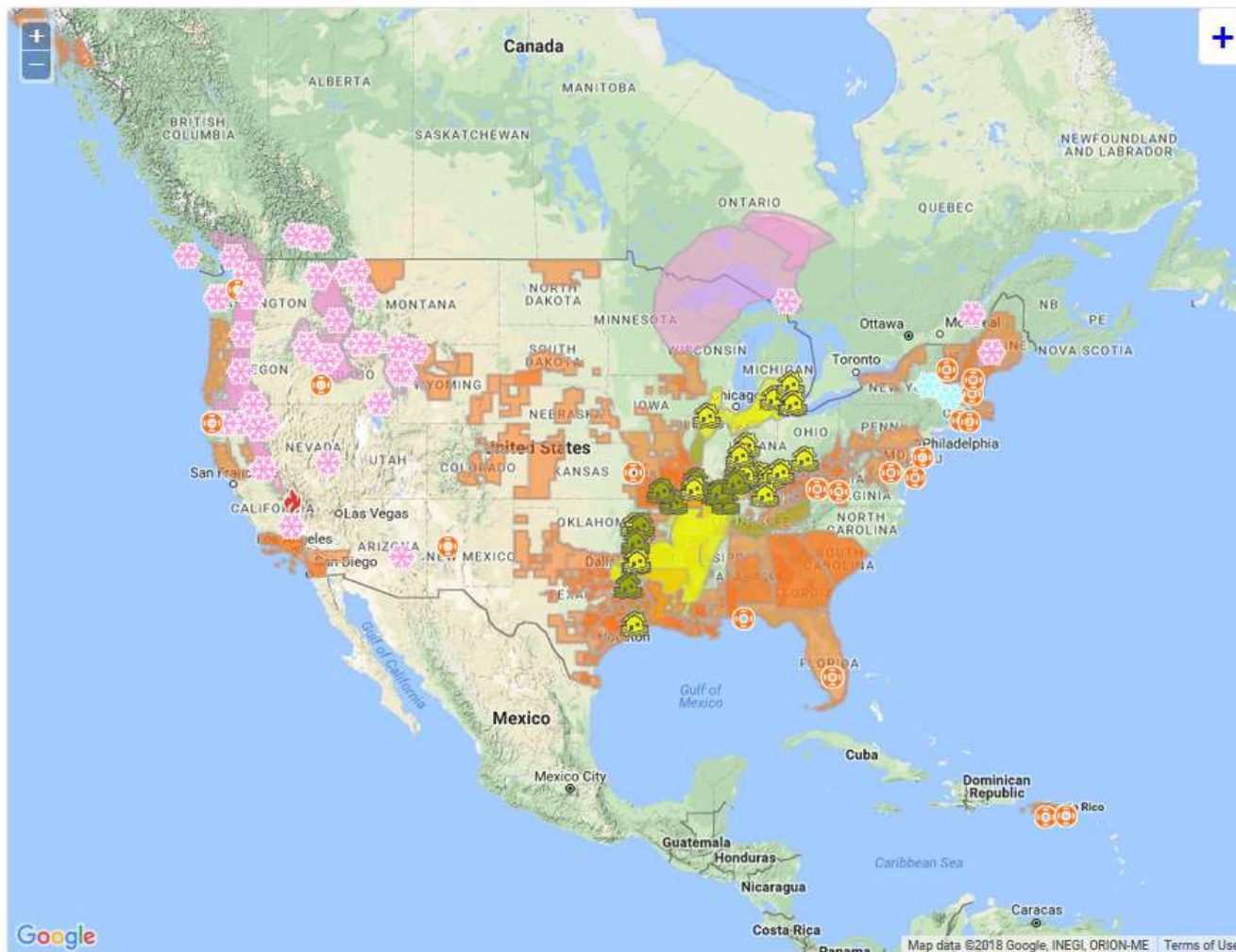
Affected Population: 9,687,771

DR-4272

Harris County, TX, Brazoria County, TX, Hidalgo County, TX and 27 other Counties

Affected Population: 8,286,901

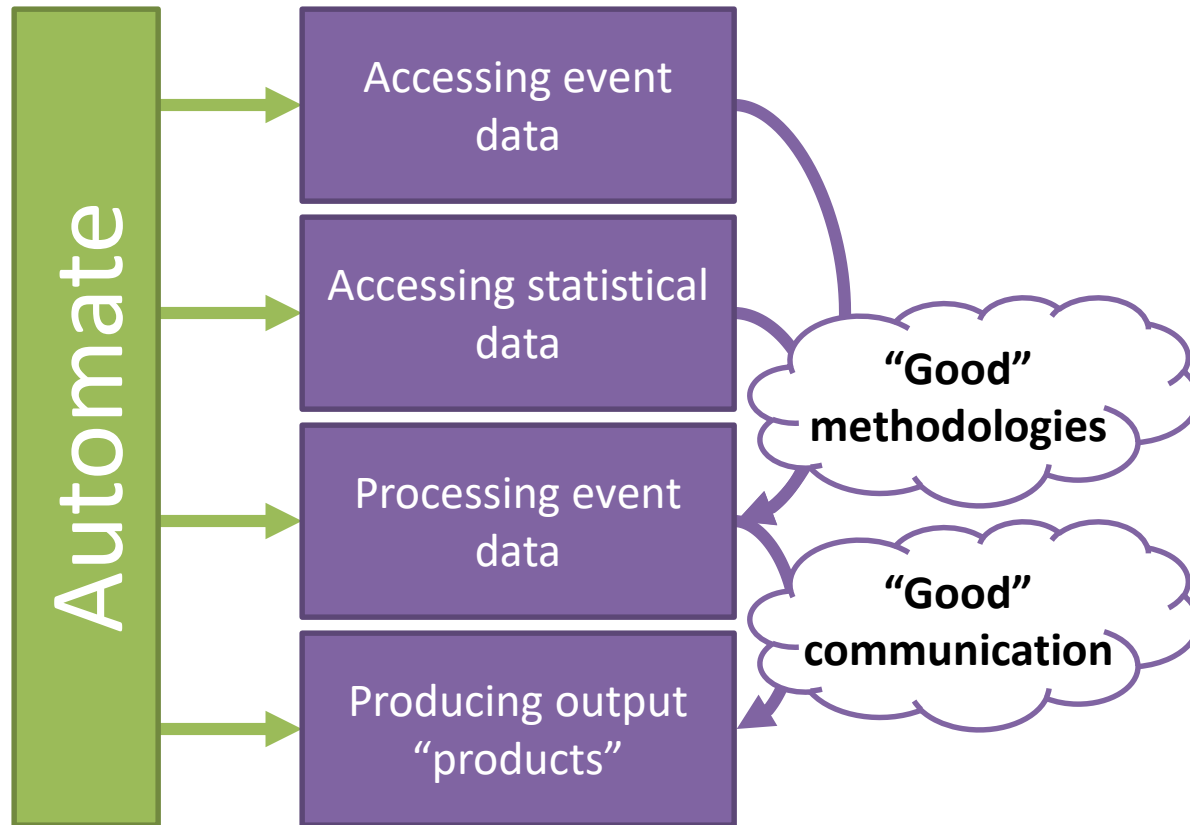
DR-4245



# Partial Timeline

- Late 1990s – LEHD Program started
- 2005 – 1<sup>st</sup> release of OnTheMap (data product and application)
- 2005 – Hurricane Katrina
- 2007 – I-35W Mississippi River bridge collapse
- 2007 – Southern California wildfires
- 2008 – Mississippi River flooding in Iowa
- 2010 – Deepwater Horizon oil spill
- 2010 – 1<sup>st</sup> release of OnTheMap for Emergency Management
- 2014 – Addition of ACS data to OTMEM

# Building OTMEM



# Data Sources



Hurricanes, Floods,  
Winter Storms



Disaster Areas



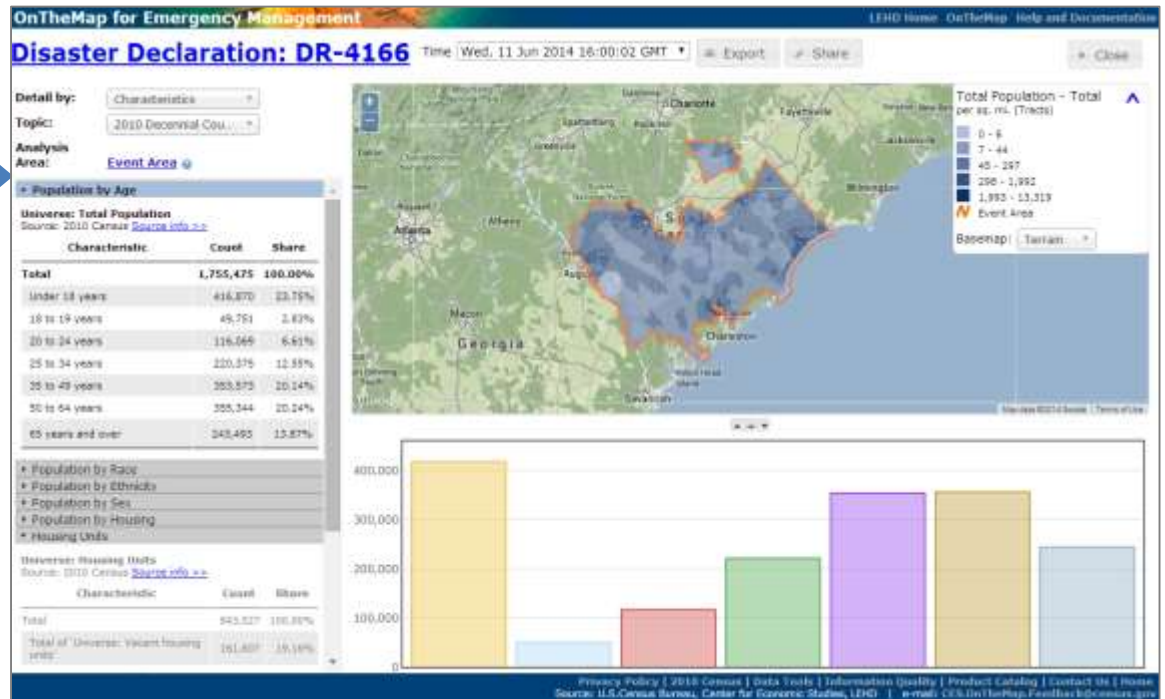
Wildfires



Demographic &  
Economic Data



## OnTheMap for Emergency Management

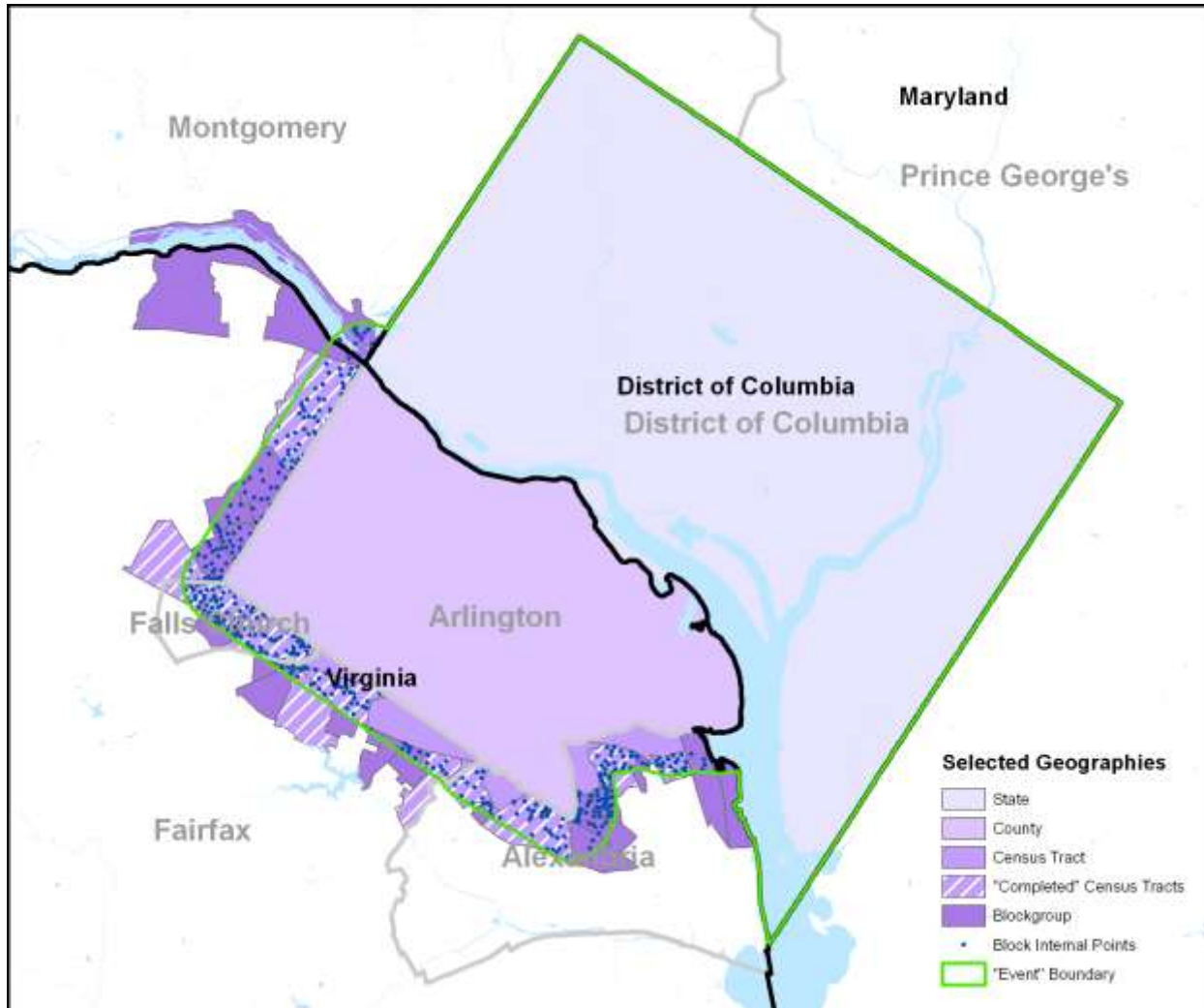




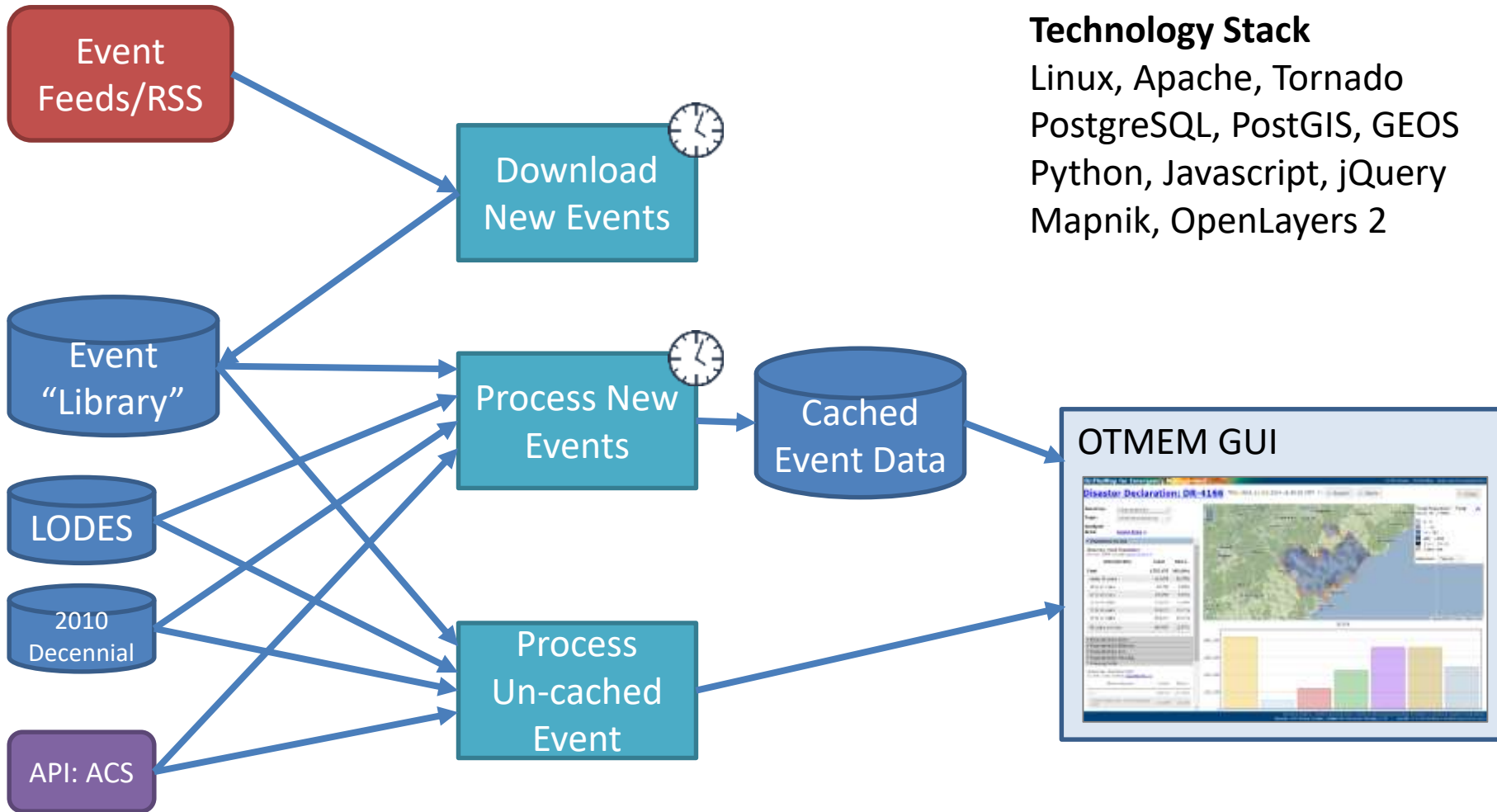
# Geographic Approximation

- Event boundaries do not align with statistical boundaries...
- Many possibilities:
  - Inclusive?/Exclusive?
  - Minimize areal difference? MOE?
  - Data-driven?
  - Deterministic/probabilistic?
- Methodology should be consistent with statistical constraints and be explained to users.

# Visual Example



# Technical Details



## Technology Stack

Linux, Apache, Tornado  
PostgreSQL, PostGIS, GEOS  
Python, Javascript, jQuery  
Mapnik, OpenLayers 2



# Generalizing this Service (FUTURE!?)

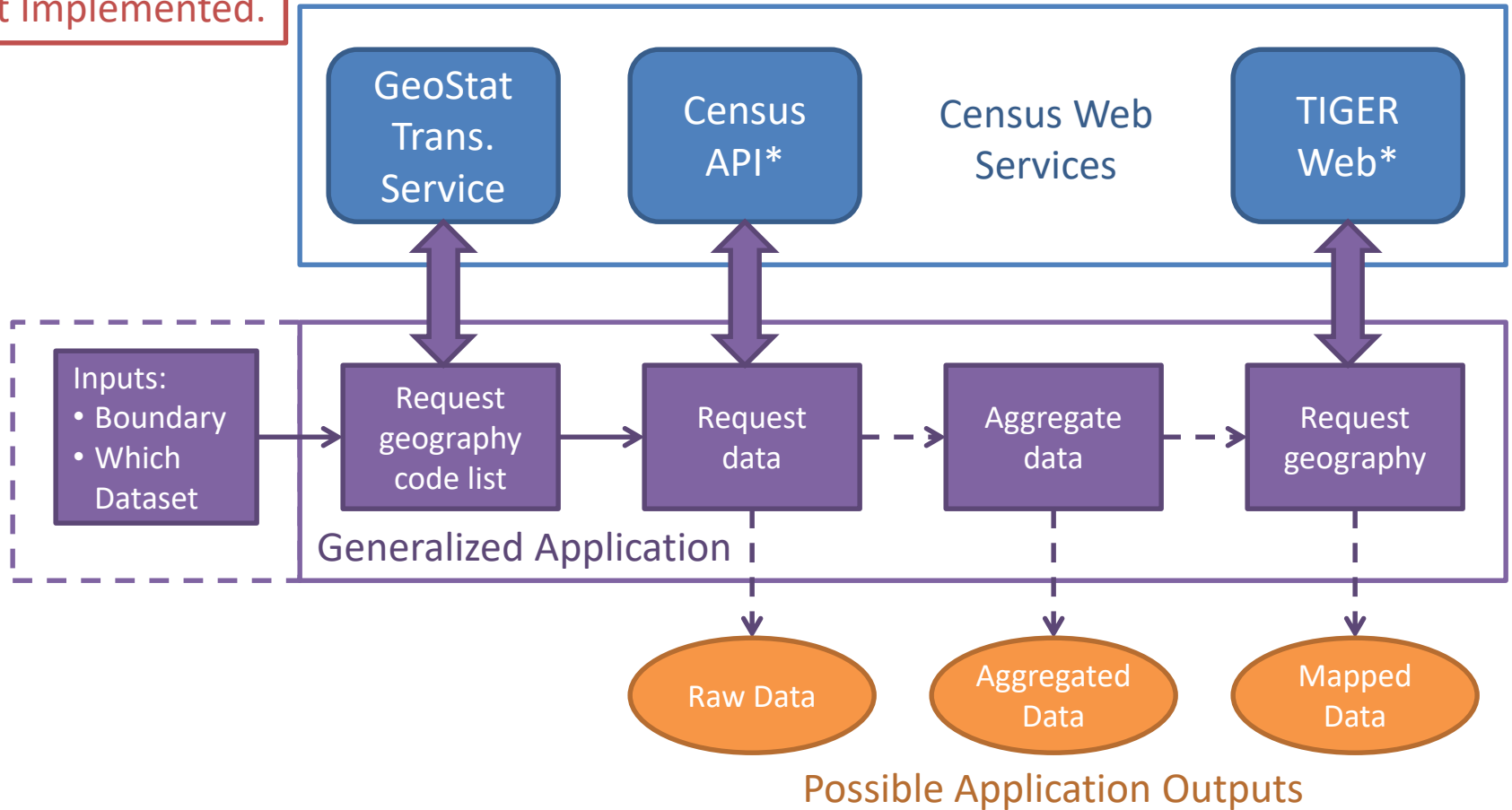
The ultimate goal is to **aggregate data to arbitrary boundaries**.

1. Start with an area of interest (boundary).
2. What data do we want to aggregate?
3. What geography matches that data?
4. What aggregation algorithm should we use?
  - Minimize MOE?
  - Inclusive/exclusive?
  - Areal/population?
  - Other requirements?

Now make it into a web service! What does that look like?

# “GeoStatistical Transformation Service”

Conceptual.  
Not Implemented.



# Potential Challenges

- Lots of datasets → Lots of geography → Large database
- Some algorithms may be slow or scale badly with query size
- Generally non-cacheable (mostly unique queries)
- Need to communicate clearly to the user community:
  - Simple/easy to explain algorithms or defaults for basic users
  - Expanded set of fully documented choices for advanced users
- “Bad” algorithms
  - Though there may be scope for letting users implement their own algorithms

# Thank You

## Questions?

# Links/Contacts

- LEHD
  - <https://lehd.ces.census.gov/>
- OnTheMap for Emergency Management
  - <https://onthemap.ces.census.gov/em>
  - [CES.OnTheMap.Feedback@census.gov](mailto:CES.OnTheMap.Feedback@census.gov)
- Questions
  - [matthew.graham@census.gov](mailto:matthew.graham@census.gov)