# Designing a Prototype Environmental Health Data Dashboard for North Carolina

Data + Team: Jerry Fu, Leah Roffman, and Anna Zolotor
Project Manager: Melyssa Minto
Project Leads: Mike Dolan Fliss, Kim Gaetz

## Abstract

Health and wellbeing are largely dependent on both the natural and built environment. More research is published every day which identifies how pollutants in our air or water can affect a number of health outcomes. In recognition of this, the Centers for Disease Control and Prevention (CDC) created the Environmental Public Health Tracking (EPHT) program to visualize population health measures alongside environmental health measures. One of the major goals of the CDC EPHT program is to help states provide a user-friendly portal through which community members, public health professionals, and policy makers can access environmental health data. Currently, North Carolina does not have its own EPHT program (only half of all states in the U.S. have funded EPHT programs). In this project, the team developed a prototype for an EPHT tool for North Carolina in the form of an Environmental Health Data Dashboard (EHDD) which can later be used to apply for a CDC grant for a fully-funded EPHT program.

This pilot includes the integration of data such as air quality, emergency department visits, and demographics data processed in RStudio and visualized as a dashboard in a Tableau (Fig 1). The three main aspects of this project were 1) Development of metadata, 2) Processing of raw datasets, and 3) Visualization in Tableau. In the course of the project, the team filled out descriptive metadata for over 150 required and non-required measures, processed and harmonized 12 raw datasets, and designed six interactive dashboard pages. An analysis of various case studies run on the prototype demonstrated the potential of the EPHT tool.

Fig 1: These icons are representations of the software used throughout the project. a) RStudio was used for data processing, b) Tableau was used for the development of the environmental health data dashboard and the metadata dashboard, and c) git was used for version control management.

## Initial goals

The major goals of this project were to create a prototype data dashboard, in the form of a Tableau story, a metadata document to house various information about each measure, and a package of R scripts to document our data processing workflow and aid future developers. The team developed goals for the prototype dashboard by envisioning what dashboard users would want to be able to get from the dashboard. Accessibility and users' situational needs were considered by compiling user stories as shown in Fig. 2.

Based on these user stories, the team concluded that our dashboard prototype should be able to:
- Show yearly data for each measure in a specific county or tract
- Show comparisons of data trends for two or more specific locations
- Visualize relationships between demographic and socioeconomic data and environmental health data.
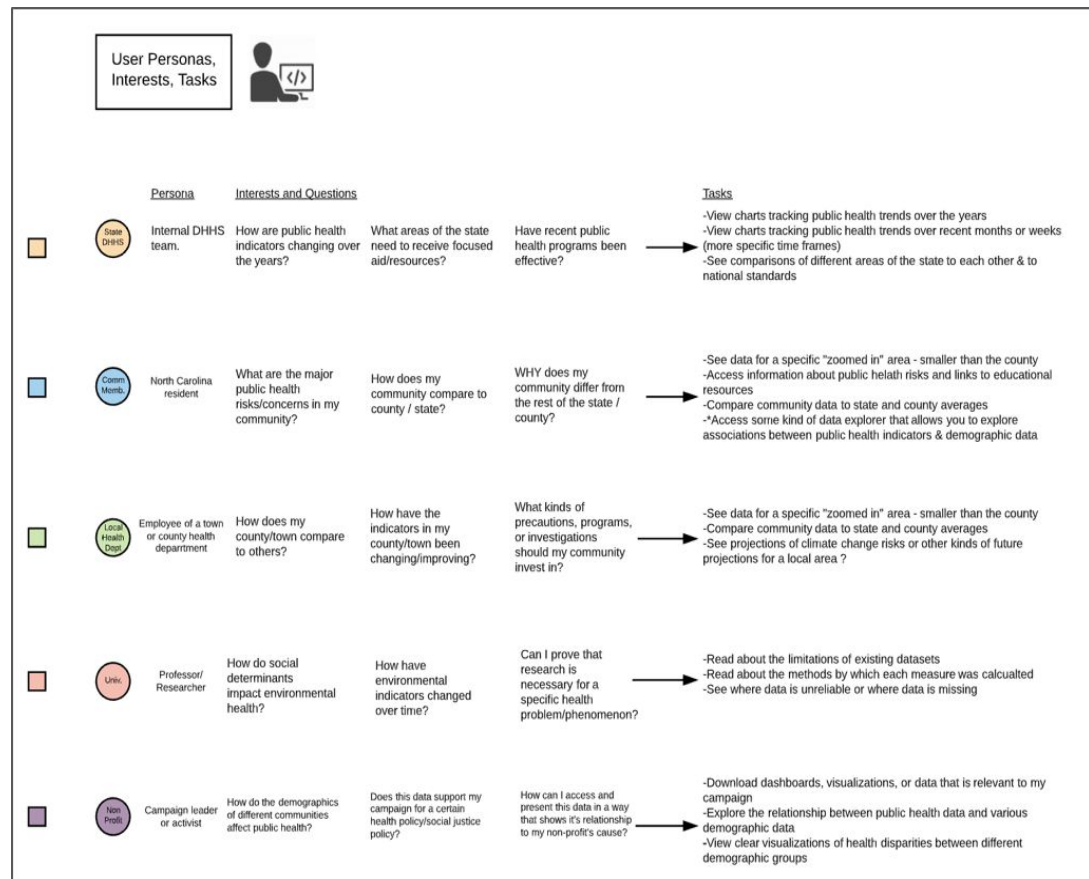


Fig 2: The five user stories were designed to reflect the needs of five potential users of the EHDD. The interests of each user were translated into "tasks" the user might seek to complete, which were used to brainstorm features to be included in the dashboard.
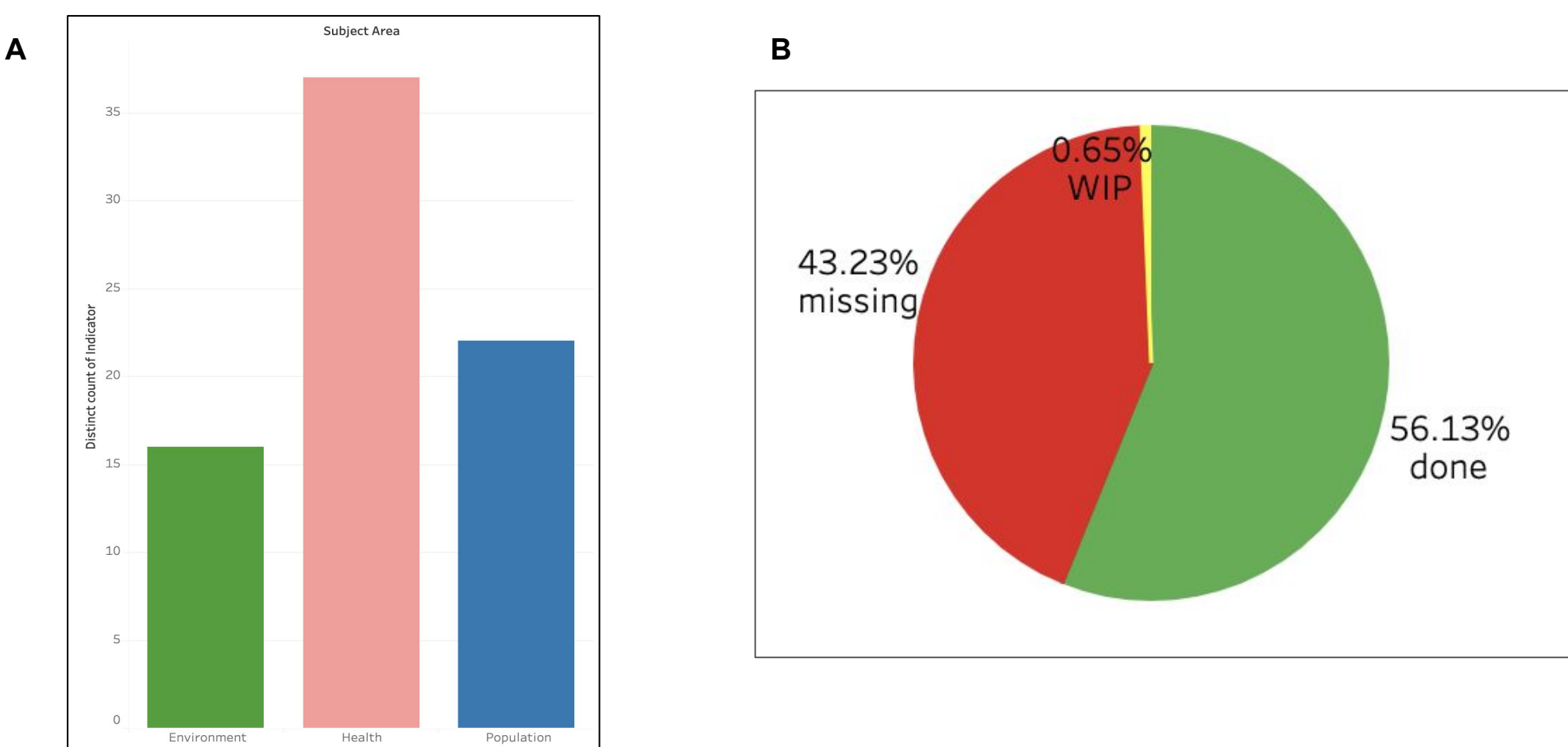
## Data Overview



Fig 3: These figures from the metadata dashboard show the distribution of the measures in our metadata a) by subject area and b) by development status. 56% of the measures in the metadata are finished being processed while others are missing.

Our data comes from a variety of sources, including the CDC and NOAA, and contains environment, health, and population measures. The organizational hierarchy of the data is subject area > content area > indicator group > indicator > measure > measure variant. At a higher level, we have 16 environment indicators, 37 health indicators, and 22 population indicators listed in our metadata. At the smallest level, there may be several measure variants for a single measure, ex) Rate of emergency department visits for COPD by age groups 25-44, 45-64, 65-84, and 85 and older.

## Methods and Workflow

All data processing centered around standardizing data to match the "data skeleton," a set of 11 columns that each set of measure data was standardized to. Those columns include three that can be used to uniquely identify any observation: year, geoid (a numeric code that indicates the place the observation is from), and measure variant identification (a code that refers to the specific way the observation was measured, as well as what was measured).

The metadata includes columns for information about every measure variant, including measure descriptions, development status, whether or not the measures is required by the CDC for the prototype, and who provided the data. Each row represents a unique measure variant. This spreadsheet is joined into the EHDD by measure id and is also used in a separate metadata visualizations dashboard.
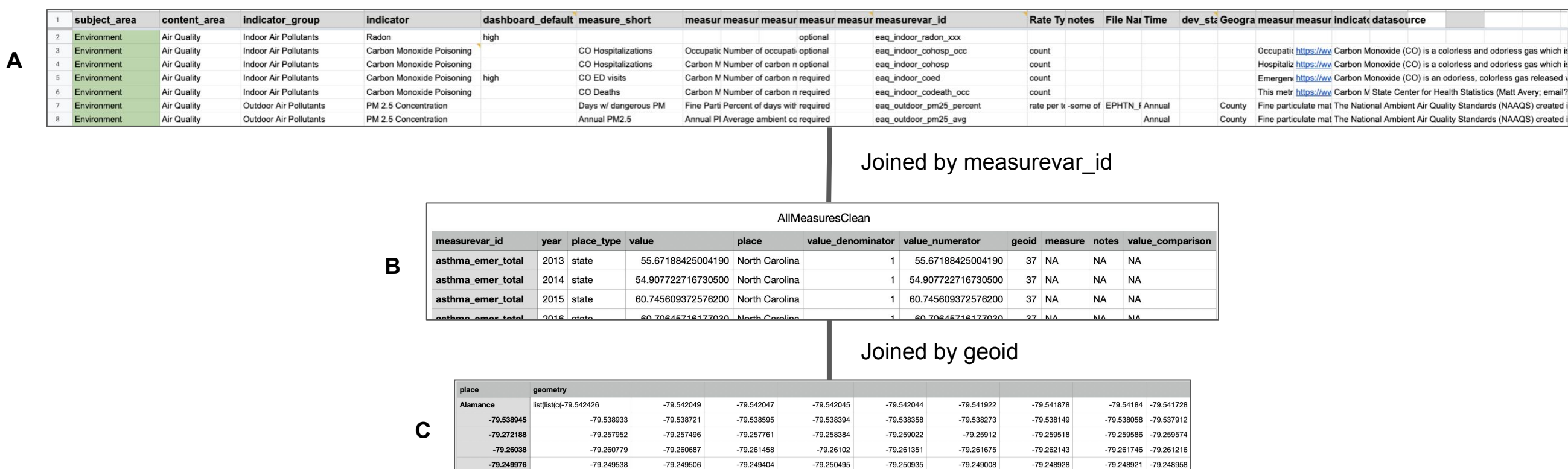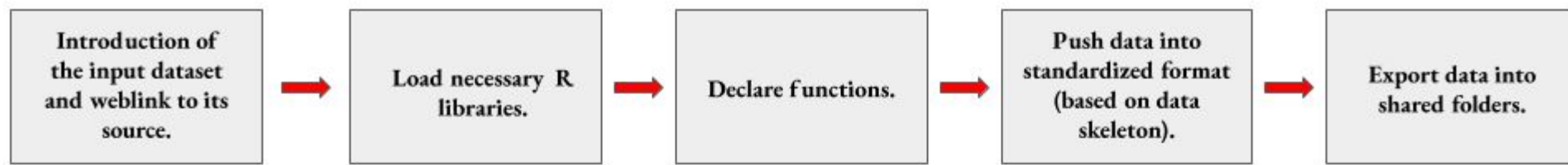


Fig 5: The data that powers the dashboard relies on joins between a) The metadata, created as a google sheet, b) The measure data, which contains vertically joined cleaned data from various data sources, and c) The multipolygon data, which is displayed here as a CSV but is actually a shapefile used for visualizing the data on a map.

## Data processing

All data processing was done in RStudio. As each new dataset was received, the goal was to push it into a standardized form based on the "data skeleton" (as discussed above).

The general order and content of each script are as follows:



Key data processing challenges included variations in the way places were named (for example, some county names were recorded in all capital letters, so converting them into lowercase changed the name of McDowell county to Mcdowell, which resulted in join failures), differences in format from year to year, and the codes for American Community Survey variables changing from year to year.
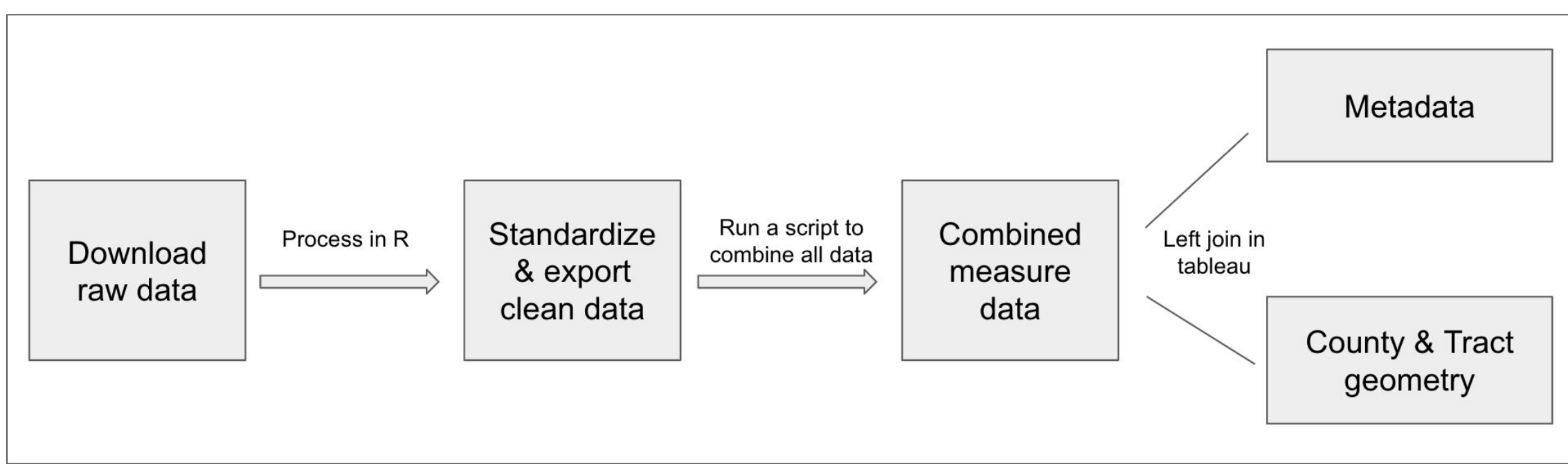


Fig 6: This flowchart shows the overarching data processing workflow. The data is standardized, harmonized, and exported using RStudio and then joined with the metadata in Tableau.

## Dashboard Design

Although the EHDD prototype is referred to as a dashboard, in actuality it consists of six separate Tableau dashboards displayed as tabs in a Tableau story. The process of designing each dashboard page began in week one by reviewing states' existing tools and wireframing pages in Lucidchart. In order to translate our wireframes into Tableau, several parameters, calculated fields and sets were created to help filter the data.

Several challenges were involved in designing the dashboards:
- Creating the map legends required making dynamic value bins to reflect the range of values for each measure.
- Displaying the units for various types of measures in the dashboard required using calculations to selectively display both rate and raw count values for certain measures
- The existence of duplicate census tract names in North Carolina required creating several calculated fields to display each place name uniquely
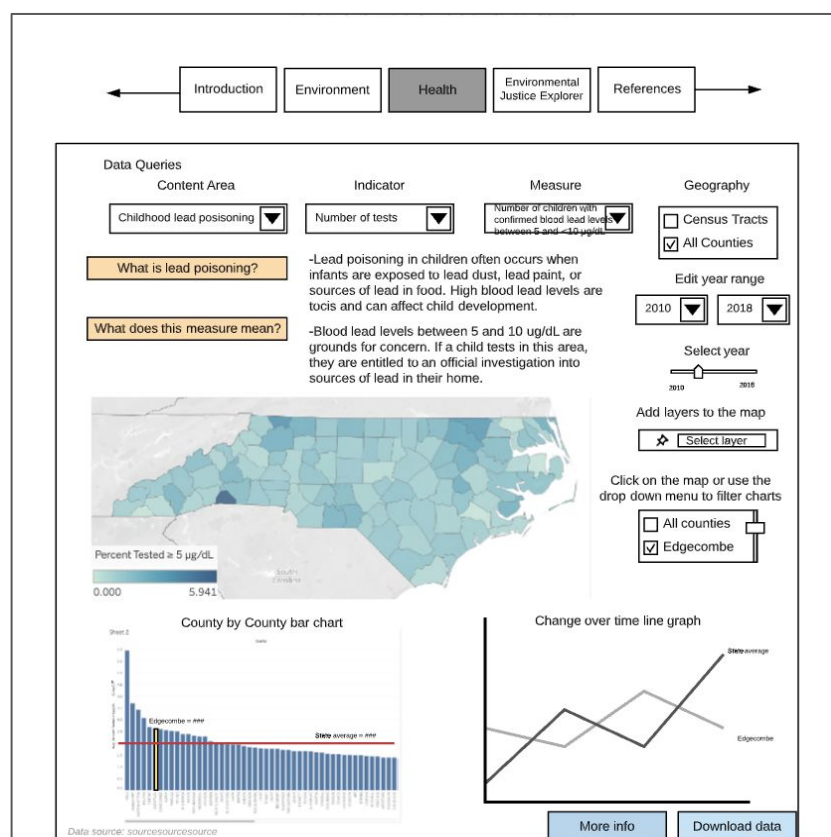


Fig 7: In the initial wireframes completed in Lucidchart, the map dashboard included map layers such as hospitals, flood zones, and highways.
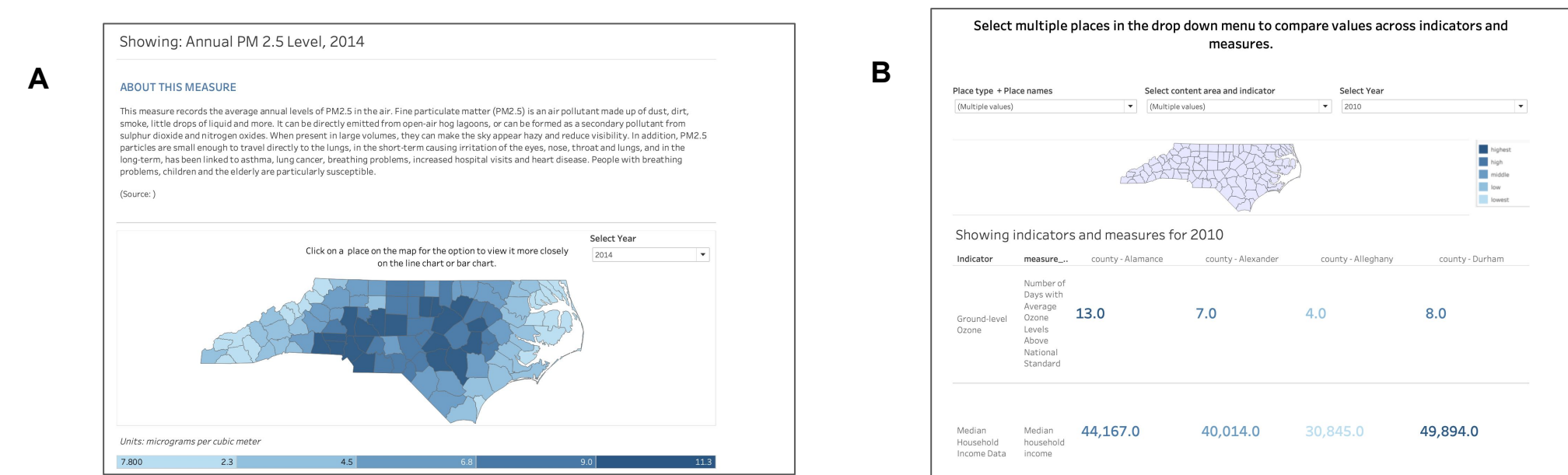
## Results and Deliverables



Fig 8: These figures show a) a partial screenshot of the prototyped Map Dashboard in Tableau, and b) a partial screenshot of the prototyped Comparison Dashboard in Tableau.

Our final dashboard prototype contained six dashboard pages in total: a welcome page, a map dashboard, a content area explorer, an indicator zoom, a fact sheet, and a comparison dashboard. Each page allows for user interactivity through drop down menus that will filter visualizations by content area, indicator and explorer. In addition, the team produced a metadata dashboard to help with data management, and a write-up paper to document challenges and suggestions for future development.

## Analysis and Limitations

In order to test the usefulness of our prototype, the team ran several test cases on it by attempting to use the dashboards to answer questions about environmental health concerns in North Carolina.
- One case study aimed to detect disparities in asthma emergency department visits around low-income housing stock like mobile home parks. The prototype was helpful in identifying counties in which there were noticeably high rates of asthma ED visits, as well as counties with a high percent of the population spending more than 30% of income on housing. However, neither of these measures were available at a census tract level of aggregation.
- In this test case and others, one of the greatest limitations of the prototype is the lack of tract-level and neighborhood-level data.
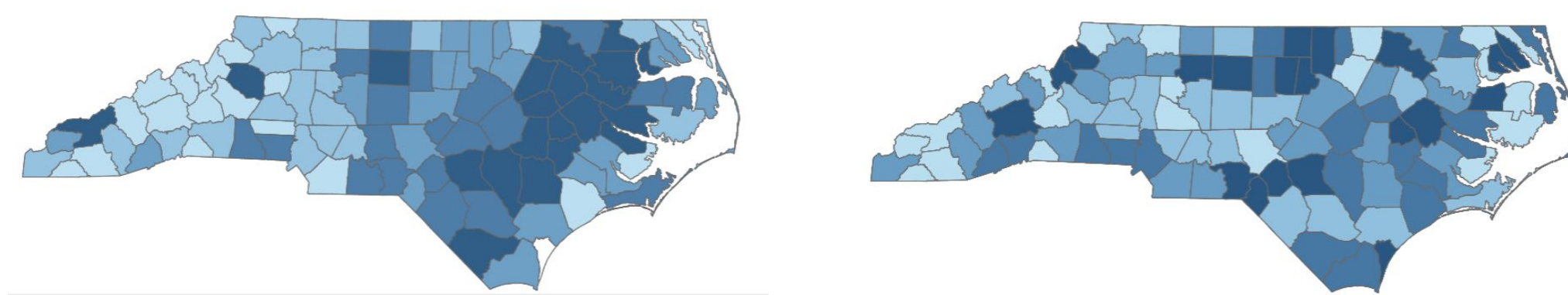


Fig. 9 and 10, a) Rate of asthma ED visits in 2019, b) Percent of population spending 30 percent or more of income on housing, 2018

## Directions for Further Development



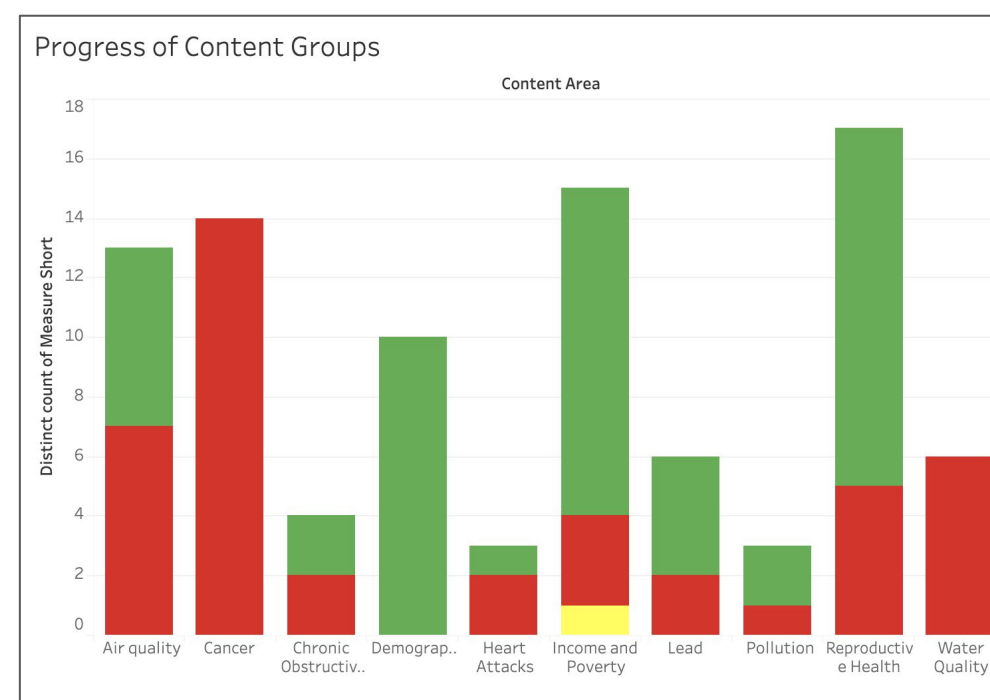Fig 11: In this figure from the metadata dashboard, red bars indicate the number of measures within each content area that have yet to be processed.

One of the first and most important steps for future developers will be to acquire and process the rest of the CDC-required indicators, including hospitalizations, cancer, and pesticide prevalence.

There are also several other areas for development:
- Creating a vulnerability index specific to environmental health.
- Continuing to perfect language of measure descriptions and improve accessibility.
- Acquiring data at more granular geographic levels, including community data if possible
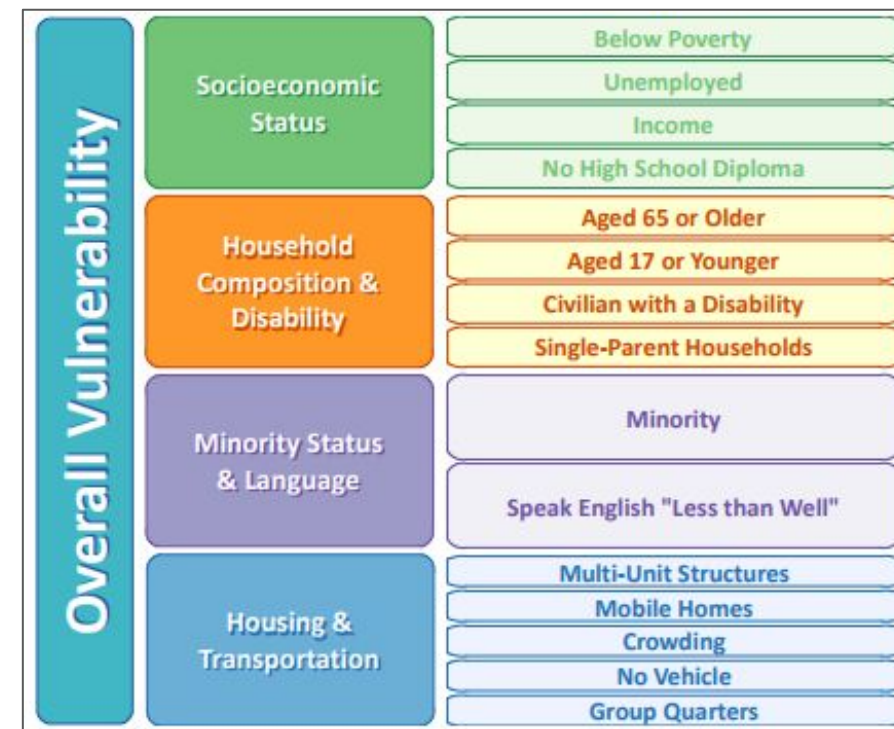


Fig 12: The CDC currently structures their social vulnerability index by equally weighing four different categories of vulnerability measures.

## Acknowledgements