

A Need-Finding Study with Users of Geospatial Data



Parker Ziegler

peziegler@cs.berkeley.edu

<https://parkie-doo.sh/>



Sarah E. Chasins

schasins@cs.berkeley.edu



Check out
the paper

CHI '23 • Working with Data • April 25, 2023

Berkeley
UNIVERSITY OF CALIFORNIA



Ok, but hold up, Parker. What is geospatial data?

(And why should we study how domain experts work with it?)

Background

Geospatial Data

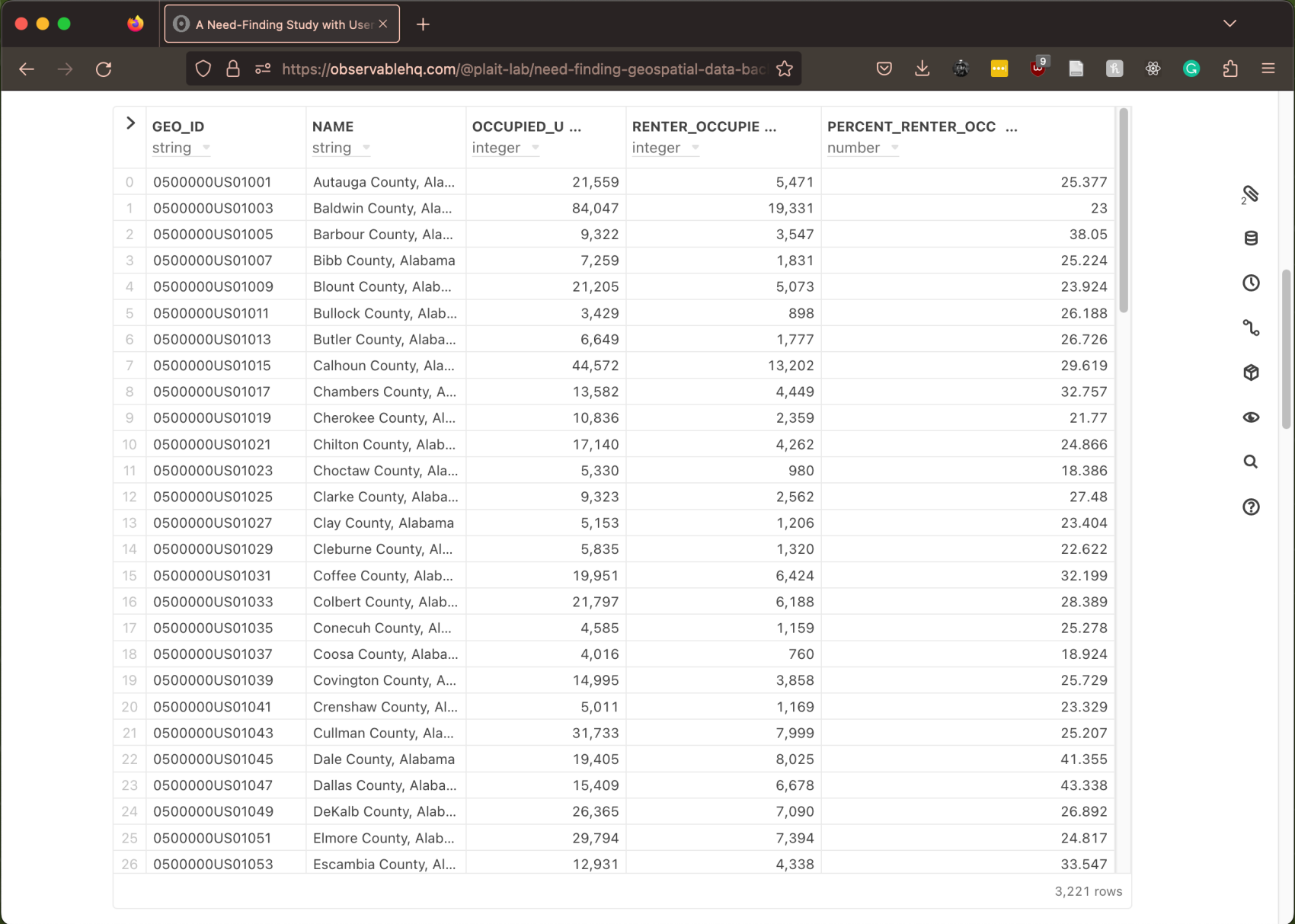
Geospatial data describes the **location** and **attributes** of phenomena on the Earth's surface.

Background

Geospatial Data

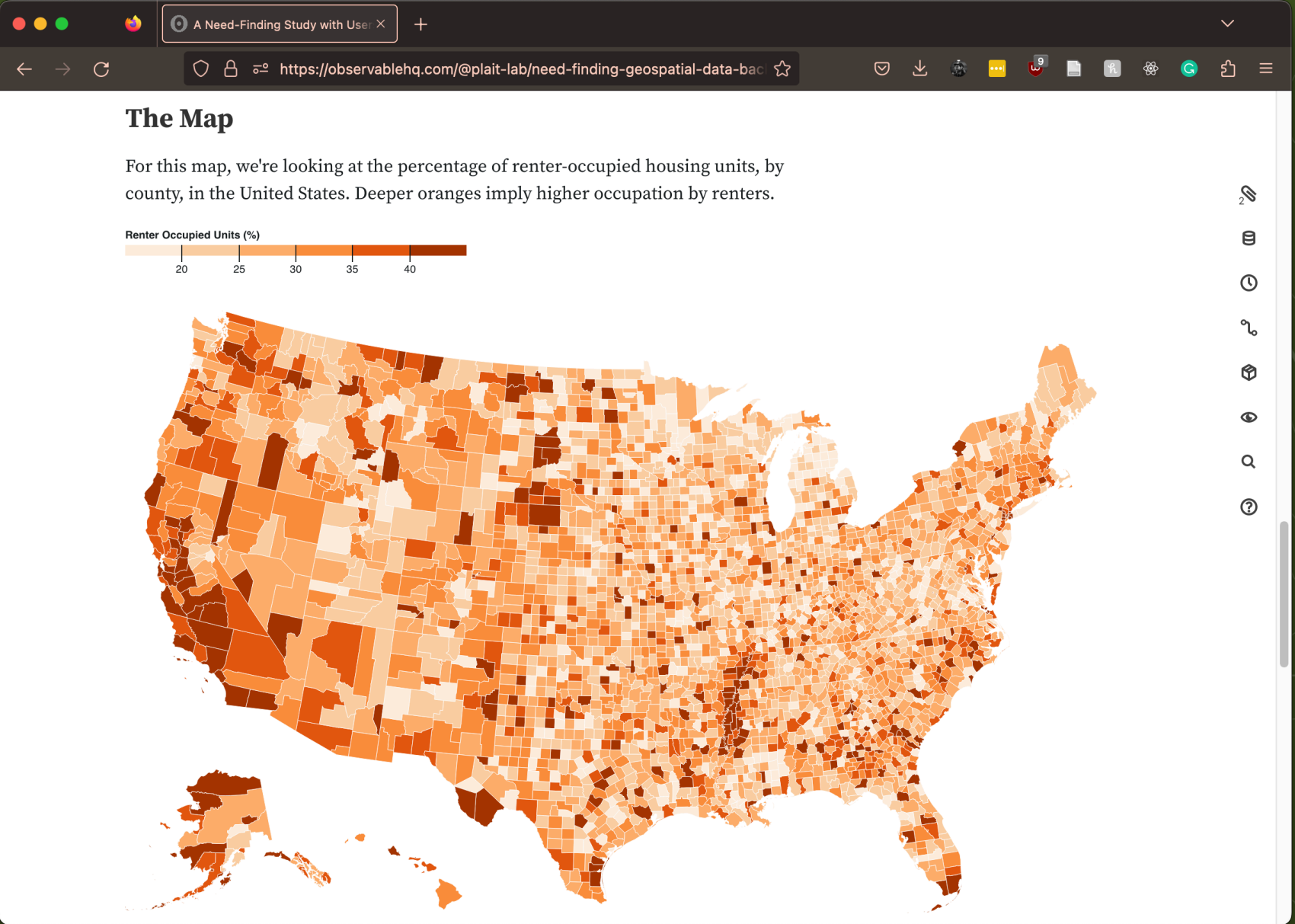
attributes

location



A screenshot of a web browser displaying a table of geospatial data. The browser's address bar shows the URL <https://observablehq.com/@plait-lab/need-finding-geospatial-data-bac>. The table has five columns: GEO_ID (string), NAME (string), OCCUPIED_U (integer), RENTER_OCCUPIE (integer), and PERCENT_RENTER_OCC (number). The data rows show information for various counties in Alabama, including Autauga, Baldwin, Barbour, Bibb, Blount, Bullock, Butler, Calhoun, Chambers, Cherokee, Chilton, Choctaw, Clarke, Clay, Cleburne, Coffee, Colbert, Conecuh, Coosa, Covington, Crenshaw, Cullman, Dale, Dallas, DeKalb, Elmore, and Escambia. The table indicates there are 3,221 rows in total.

	GEO_ID string	NAME string	OCCUPIED_U ... integer	RENTER_OCCUPIE ... integer	PERCENT_RENTER_OCC ... number
0	0500000US01001	Autauga County, Ala...	21,559	5,471	25.377
1	0500000US01003	Baldwin County, Ala...	84,047	19,331	23
2	0500000US01005	Barbour County, Ala...	9,322	3,547	38.05
3	0500000US01007	Bibb County, Alabama	7,259	1,831	25.224
4	0500000US01009	Blount County, Alab...	21,205	5,073	23.924
5	0500000US01011	Bullock County, Alab...	3,429	898	26.188
6	0500000US01013	Butler County, Alaba...	6,649	1,777	26.726
7	0500000US01015	Calhoun County, Ala...	44,572	13,202	29.619
8	0500000US01017	Chambers County, A...	13,582	4,449	32.757
9	0500000US01019	Cherokee County, Al...	10,836	2,359	21.77
10	0500000US01021	Chilton County, Alab...	17,140	4,262	24.866
11	0500000US01023	Choctaw County, Ala...	5,330	980	18.386
12	0500000US01025	Clarke County, Alaba...	9,323	2,562	27.48
13	0500000US01027	Clay County, Alabama	5,153	1,206	23.404
14	0500000US01029	Cleburne County, Al...	5,835	1,320	22.622
15	0500000US01031	Coffee County, Alab...	19,951	6,424	32.199
16	0500000US01033	Colbert County, Alab...	21,797	6,188	28.389
17	0500000US01035	Conecuh County, AL...	4,585	1,159	25.278
18	0500000US01037	Coosa County, Alaba...	4,016	760	18.924
19	0500000US01039	Covington County, A...	14,995	3,858	25.729
20	0500000US01041	Crenshaw County, Al...	5,011	1,169	23.329
21	0500000US01043	Cullman County, Ala...	31,733	7,999	25.207
22	0500000US01045	Dale County, Alabama	19,405	8,025	41.355
23	0500000US01047	Dallas County, Alaba...	15,409	6,678	43.338
24	0500000US01049	DeKalb County, Alab...	26,365	7,090	26.892
25	0500000US01051	Elmore County, Alab...	29,794	7,394	24.817
26	0500000US01053	Escambia County, Al...	12,931	4,338	33.547



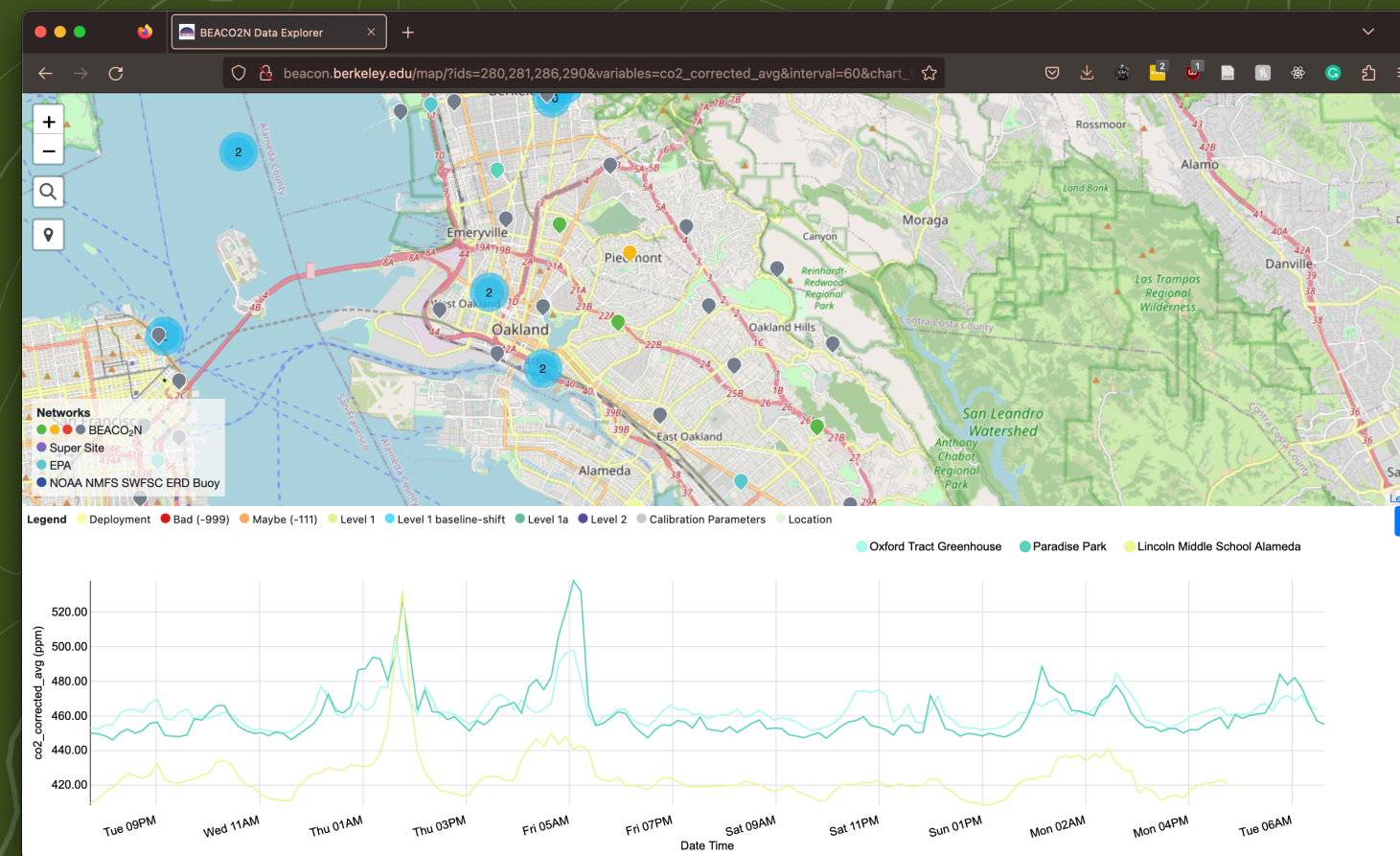
Background

Geospatial Data

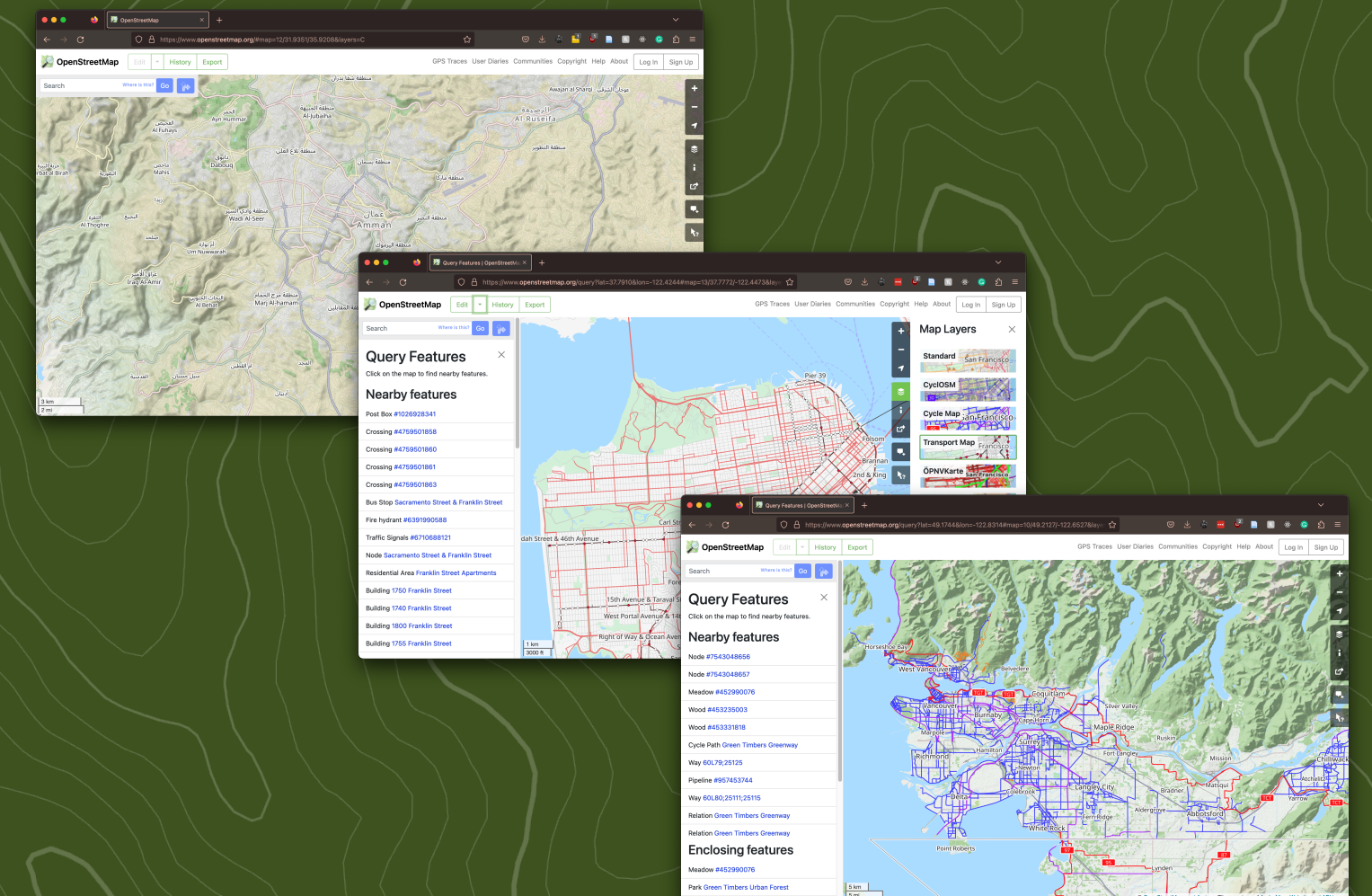
Geospatial data is everywhere today.



Satellite Imagery



Environmental Sensor Networks



OpenStreetMap

Background

Domain Experts and Geospatial Data



**Earth and
Climate
Science**



**Social
Sciences**



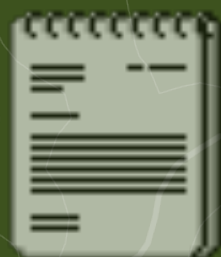
**Data
Journalism**

Background

Domain Experts and Geospatial Data



Earth and
Climate Science



Social Sciences



Data Journalism

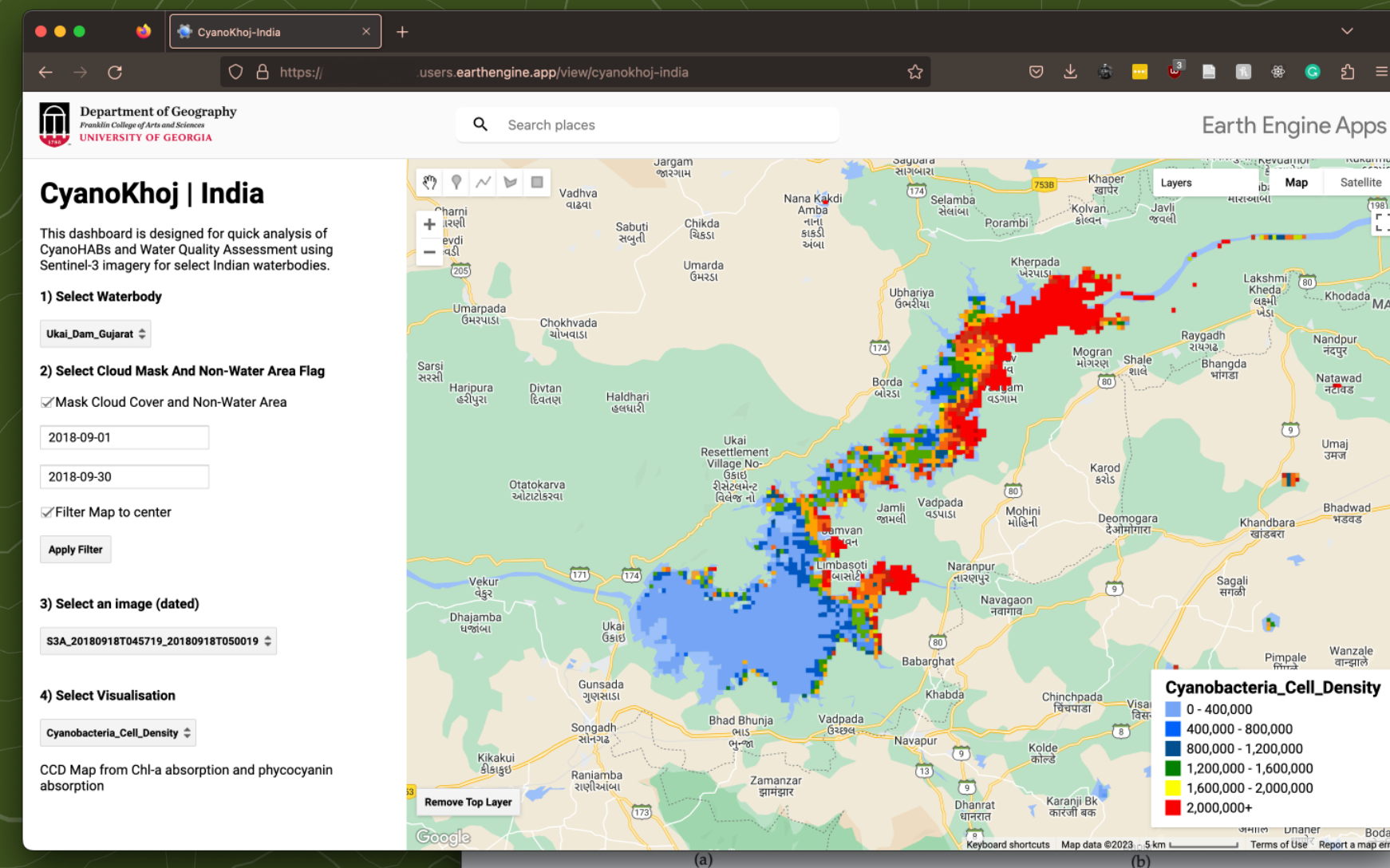
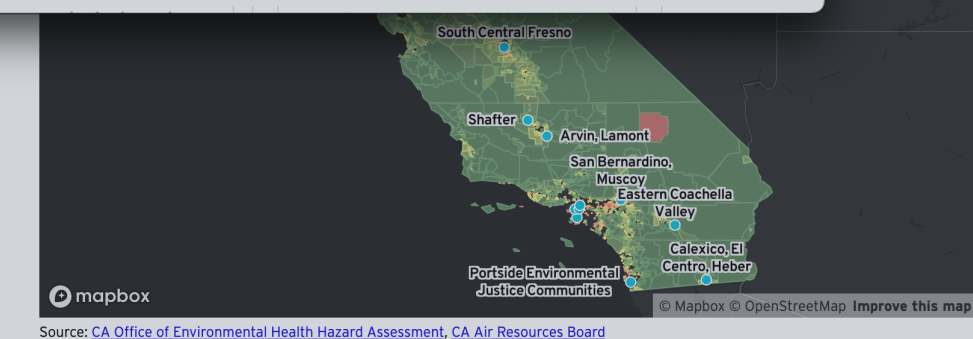
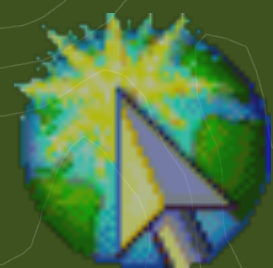


Figure 2: City Block Extraction. (a) We collect the set of streets for each administrative unit, formatted originally as individual LineStrings (shown here in varying colors for distinction), to form a regional street network formatted as a MultiLineString object. (b) This unioned object is buffered by a small amount to render the one-dimensional object as a two-dimensional object (shown here with an exaggerated buffer). This allows us to compute a set-theoretic difference between the GADM boundary and the buffered street network. (c) The difference between these two objects gives the geometric description of street blocks. Colors distinguish between adjacent street blocks. All subfigures show the street network and block geometries for Freetown, Sierra Leone.



Background

Domain Experts and Geospatial Data



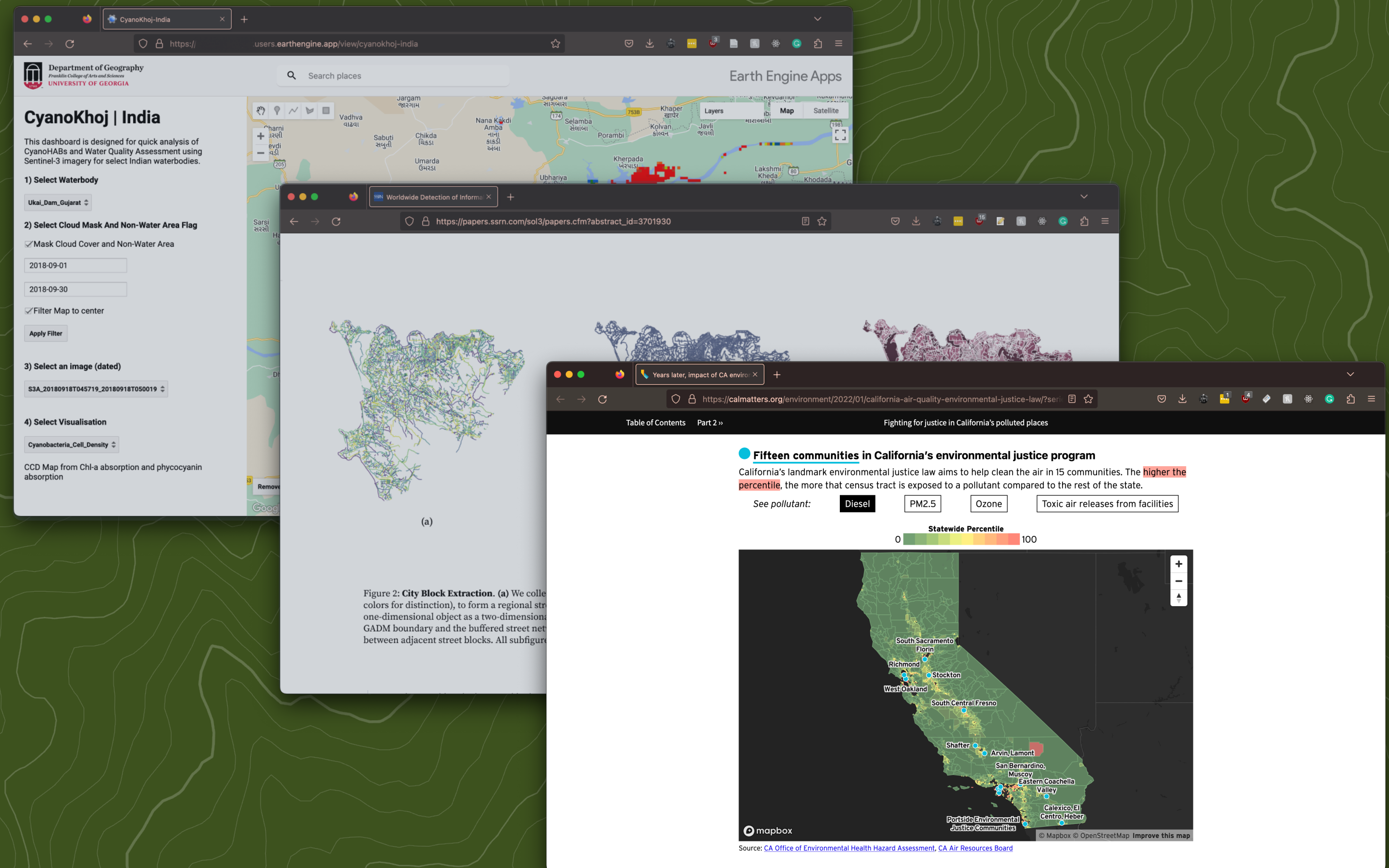
Earth and
Climate Science



Social Sciences



Data Journalism



Barriers to working with geospatial data are high.

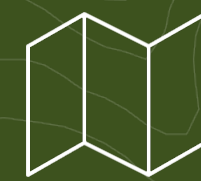
Barriers to working with geospatial data are high.

Geographic Information Systems

- Require significant background in **geospatial data theory**



Geography



Cartography

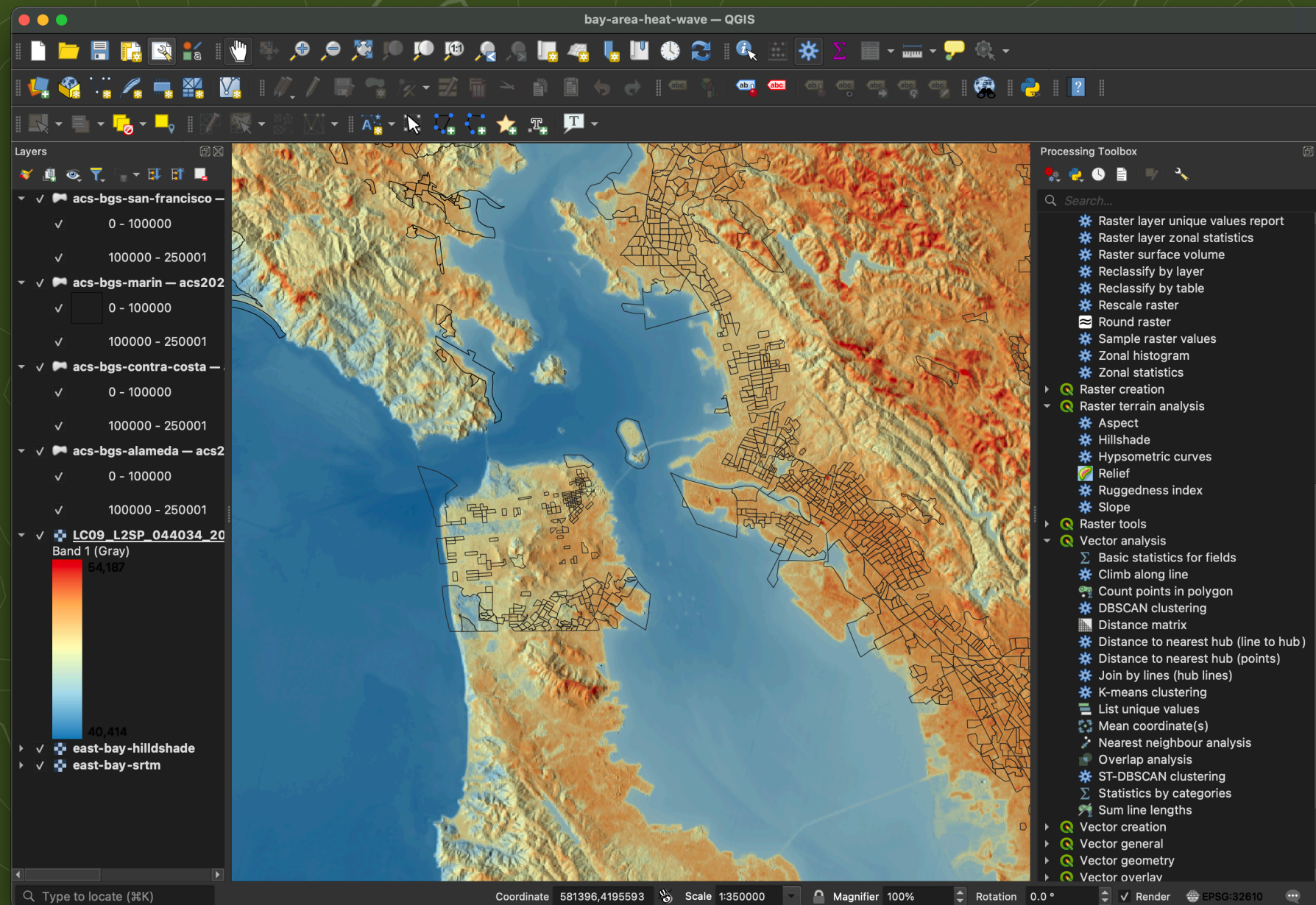


Databases



Statistics

- HCI research^{1, 2, 3} has shown that GISs are especially difficult for non-geographers to learn and use.



Example



QGIS

1. Traynor, C. and Williams, M.G. Why are geographic information systems hard to use? *Conference Companion on Human Factors in Computing Systems* (1995).
2. Traynor, C. & Williams, M. G. End users and GIS: a demonstration is worth a thousand words. in *Your wish is my command: programming by example* 115–134 (Morgan Kaufmann Publishers Inc., 2001).
3. Haklay, M. (Muki) & Skarlatidou, A. Human-Computer Interaction and Geospatial Technologies – Context. in *Interacting with Geospatial Technologies* 1–18 (John Wiley & Sons, Ltd, 2010). doi:10.1002/9780470689813.ch1.

Barriers to working with geospatial data are high.

Programming Systems

- Geospatial programming abstractions are increasingly common in Python, R, and JavaScript



geopandas

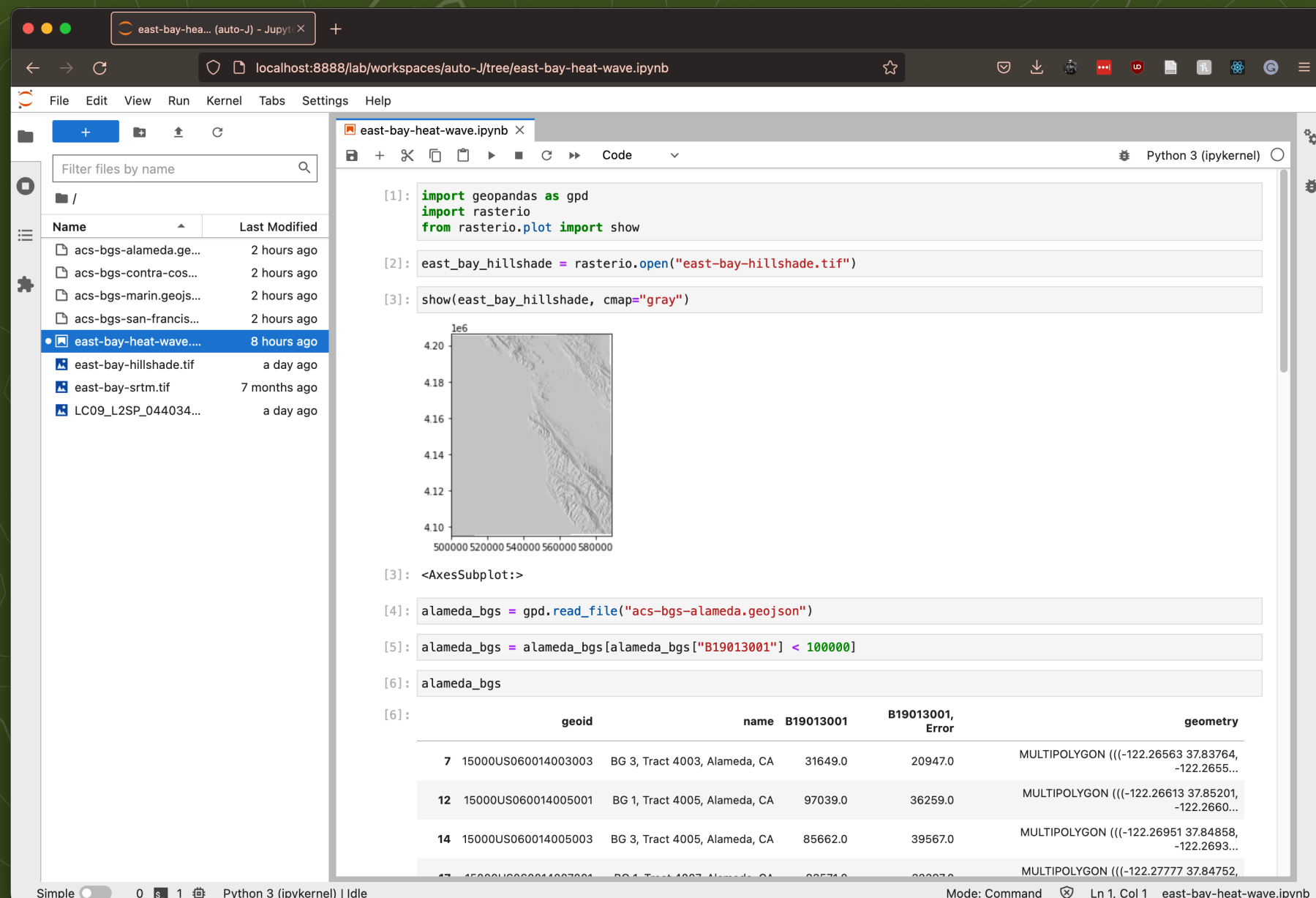


sf



mapbox

- Must develop proficiency with **programming languages** and **environments**



Example



Jupyter Notebooks

Background

Research has yet to explore the specific obstacles **domain experts** face in their work with geospatial data.



Computational Notebooks



Design Software



Geospatial Analysis and Visualization Libraries



Analysis
Visualization

Data Discovery
Data Transformation
Analysis Representation

Background

Contribution

The goal of this research is to **identify the computing needs of domain expert geospatial data users.**

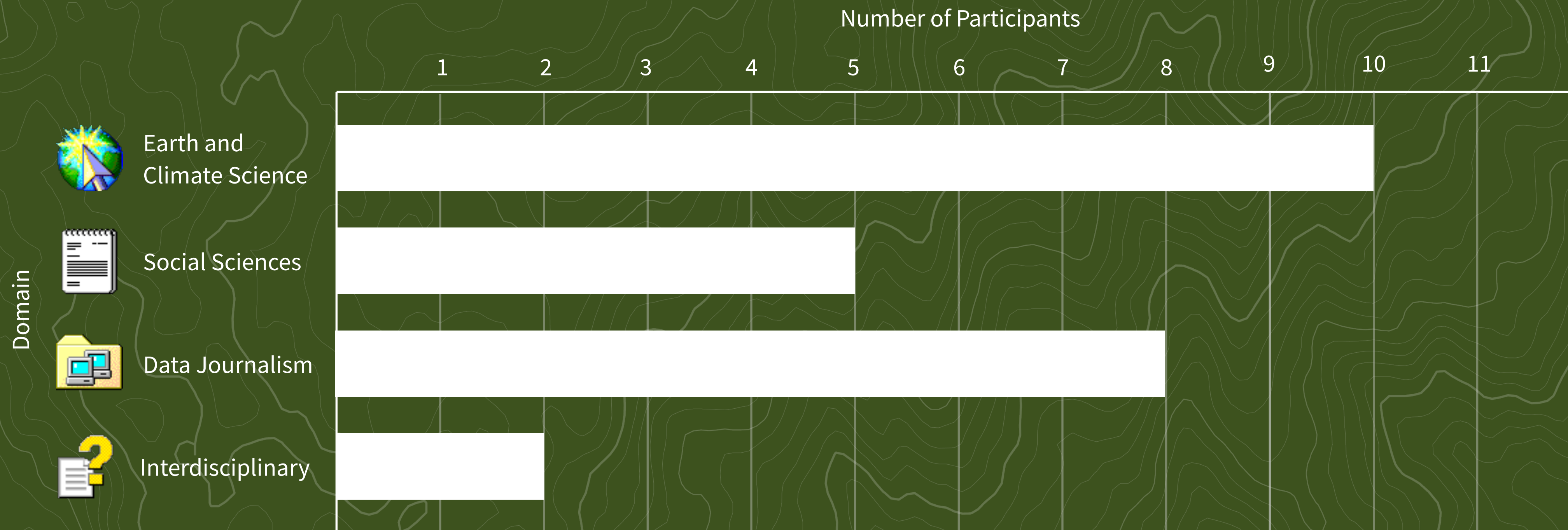
Roadmap



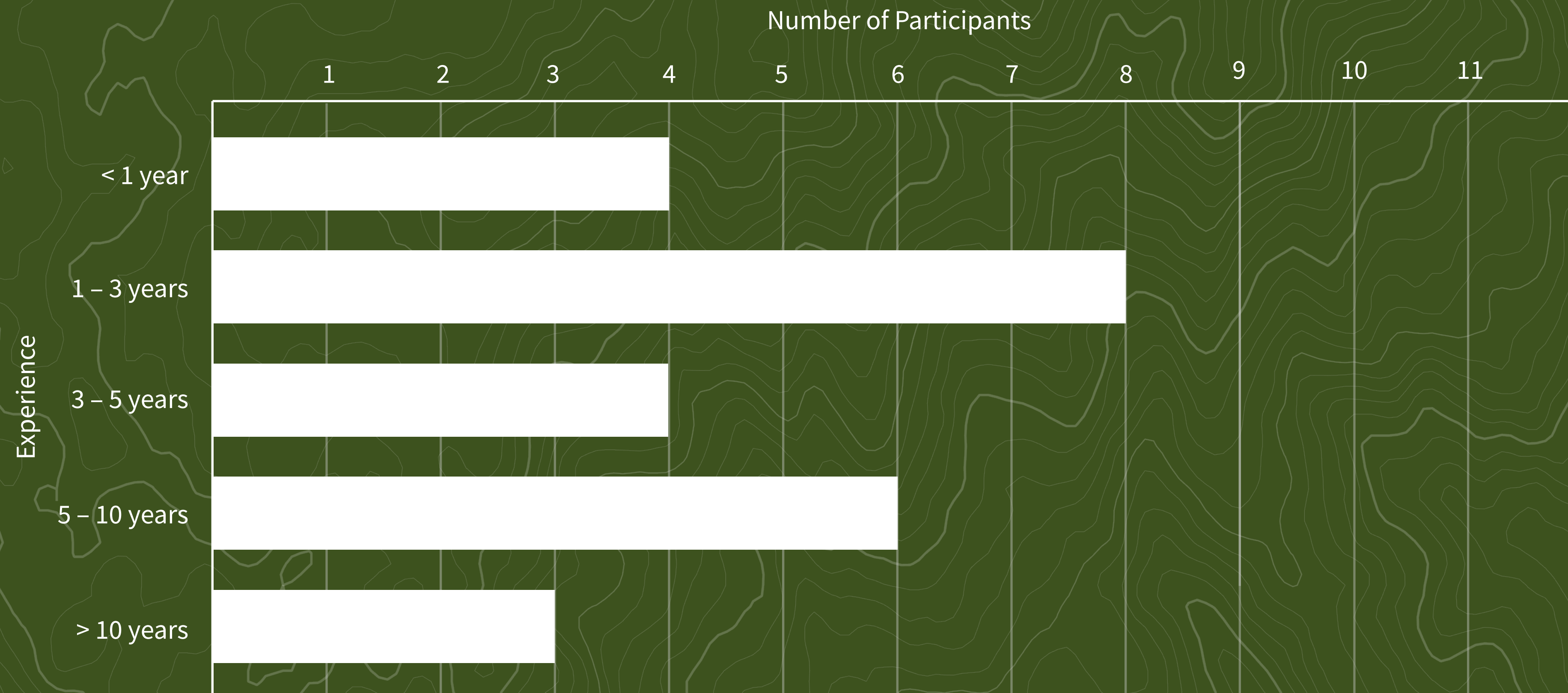
Roadmap



We conducted a contextual inquiry study with 25 participants.



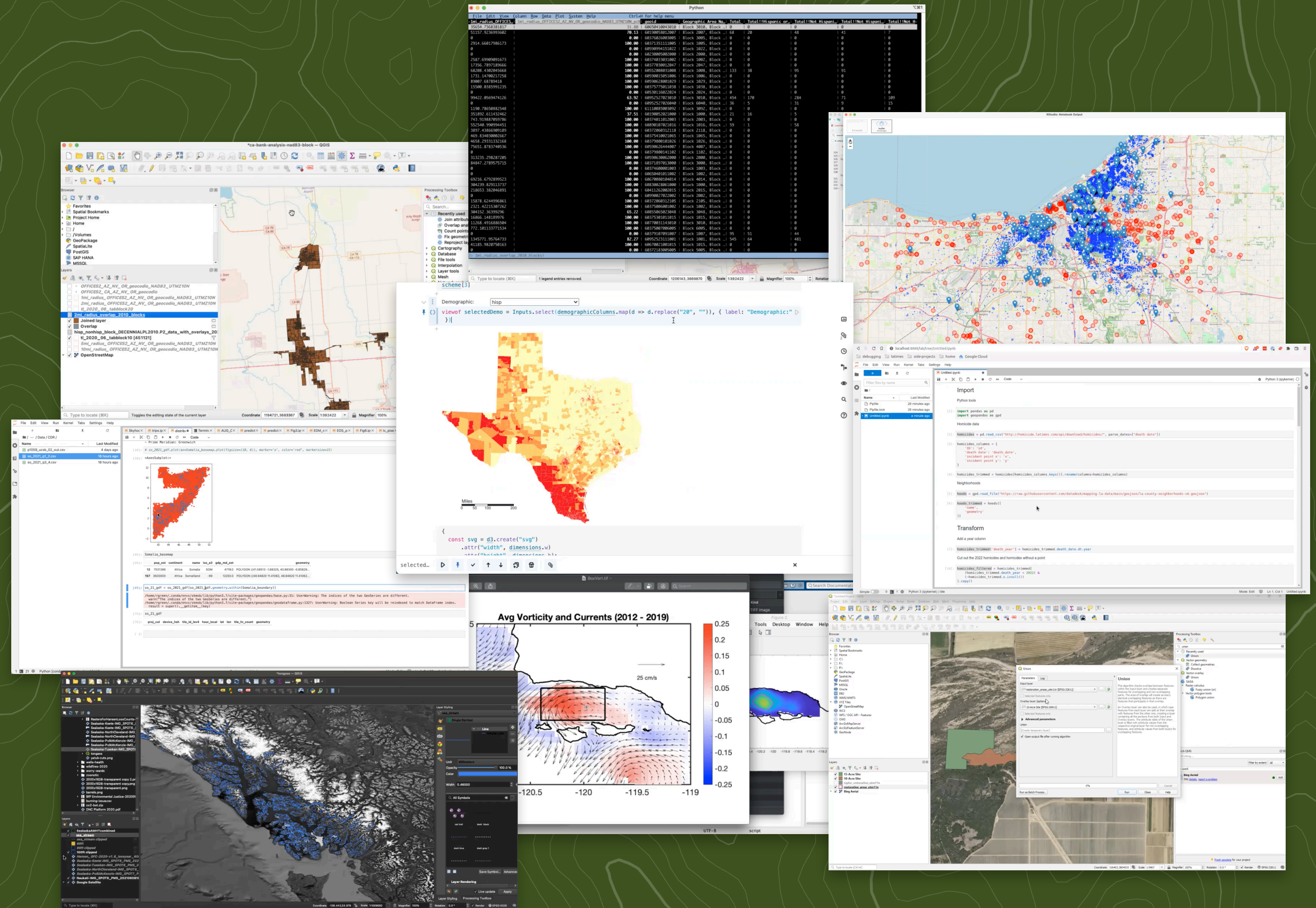
We conducted a contextual inquiry study with 25 participants.



Study Design

Session Structure and Analysis

- 50–70 minute open-task **observations**
 - Followed by **semi-structured post-interviews**
- Inductive **thematic analysis** on the **29 hours** of video recordings



Roadmap

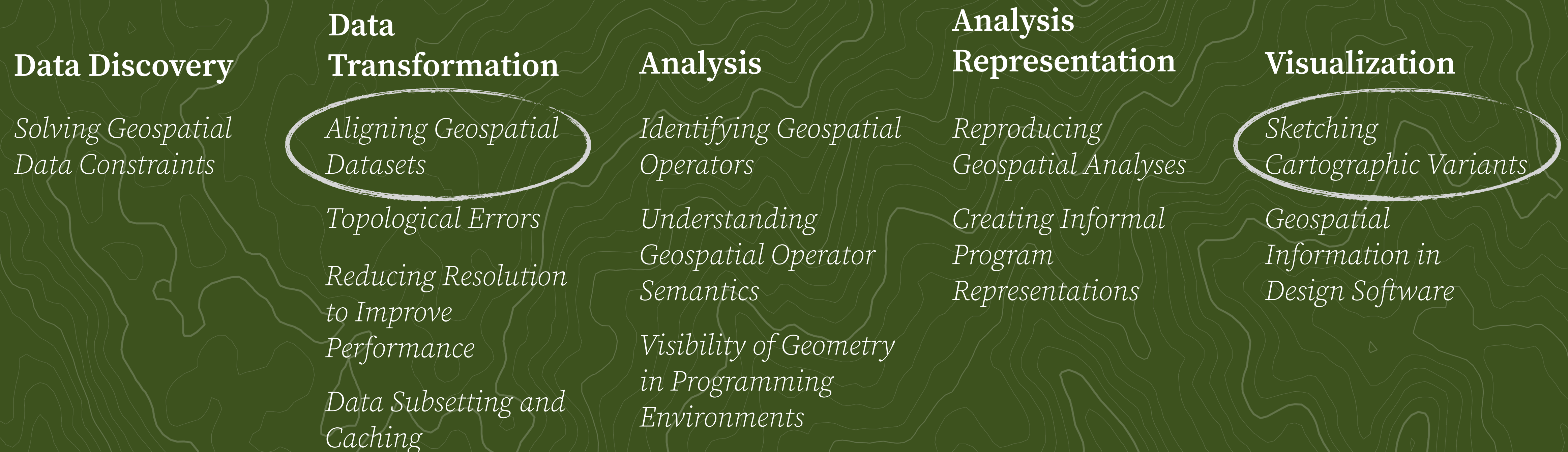


Roadmap



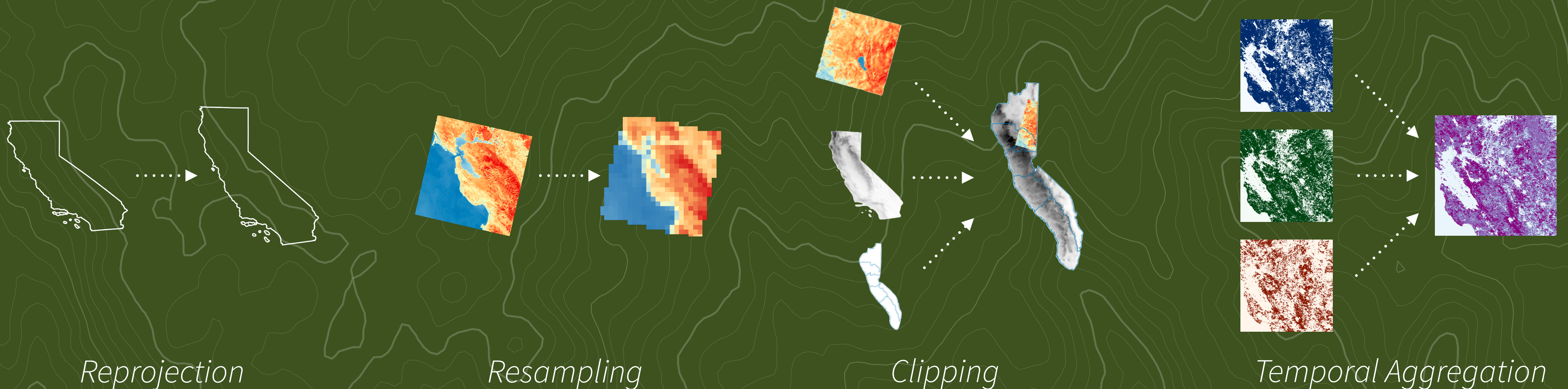
Findings

We identified **12 challenges** across **five phases** of participants' work with geospatial data.



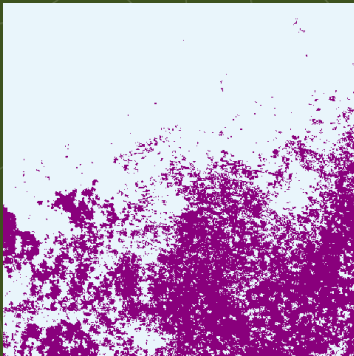


Aligning Geospatial Datasets

Participants needed to transform datasets to a **shared spatial and temporal reference** for analysis, but alignment required **complex preprocessing**.



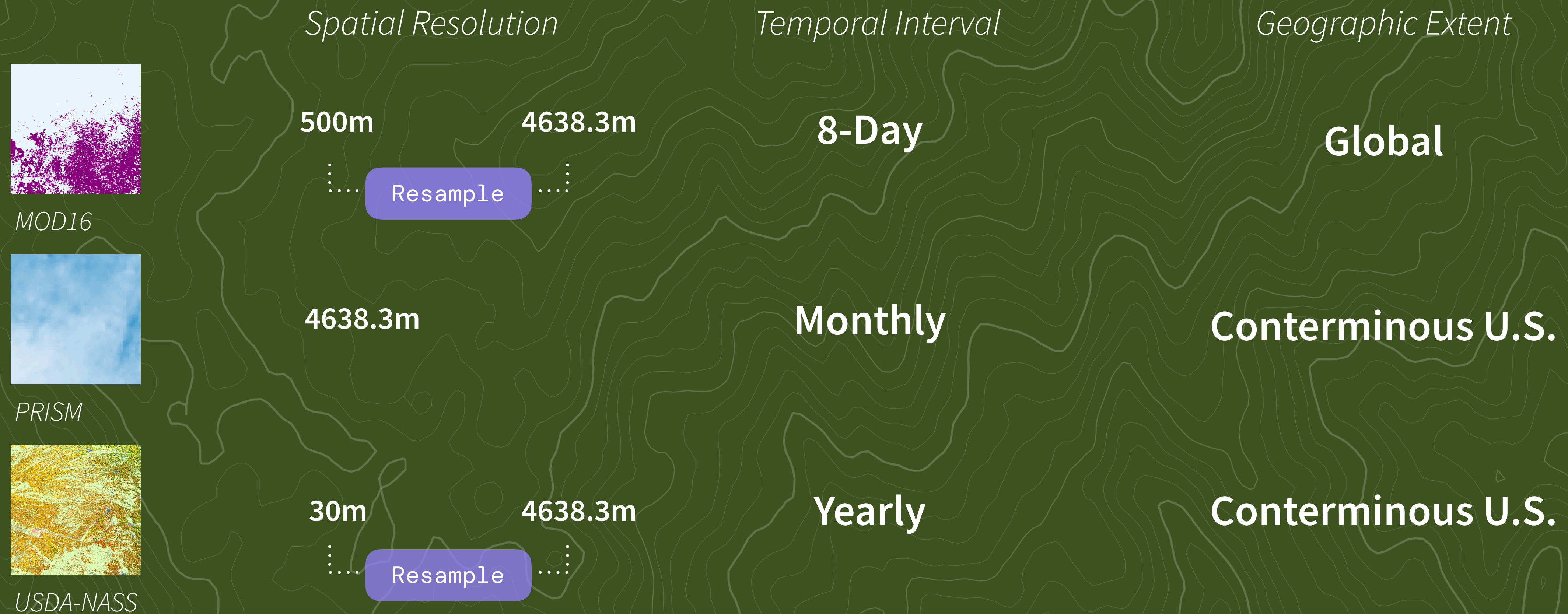
Aligning Geospatial Datasets

PE2's Task. Develop a model to predict groundwater withdrawal.

	<i>Spatial Resolution</i>	<i>Temporal Interval</i>	<i>Geographic Extent</i>
 <i>MOD16</i>	500m	8-Day	Global
 <i>PRISM</i>	4638.3m	Monthly	Conterminous U.S.
 <i>USDA-NASS</i>	30m	Yearly	Conterminous U.S.

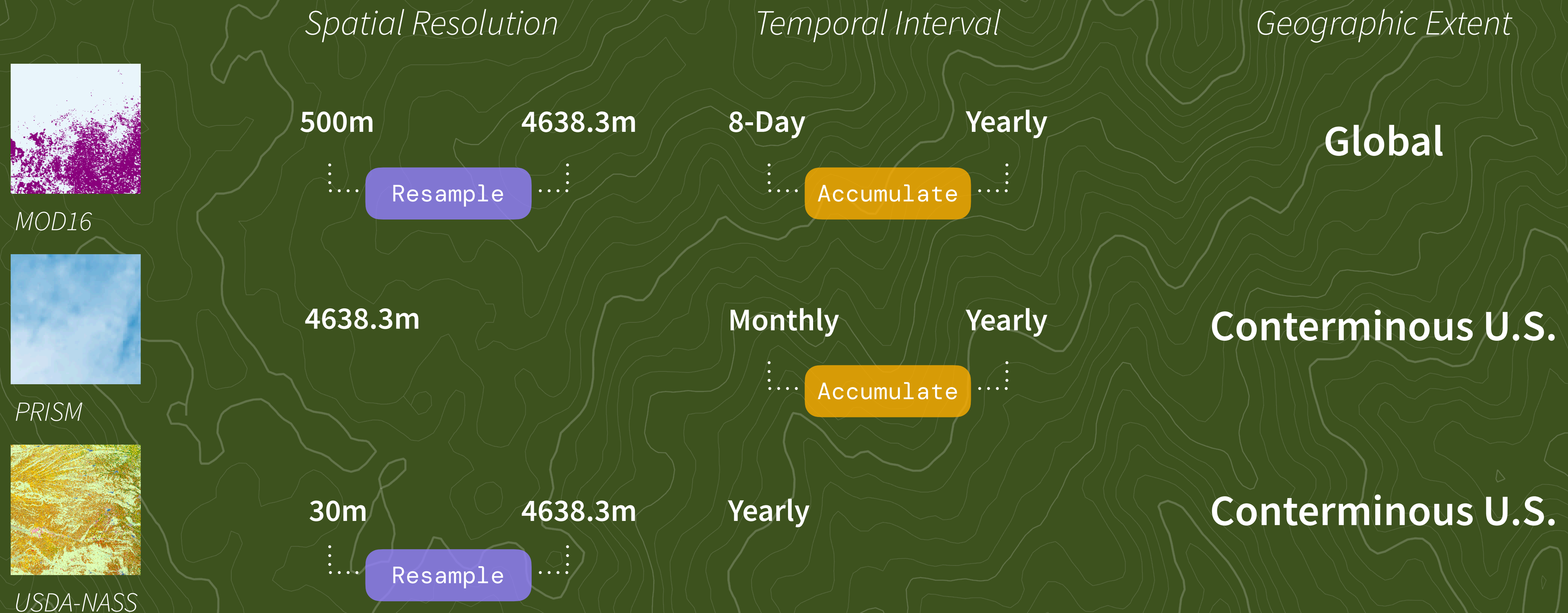
Aligning Geospatial Datasets

PE2's Task. Develop a model to predict groundwater withdrawal.



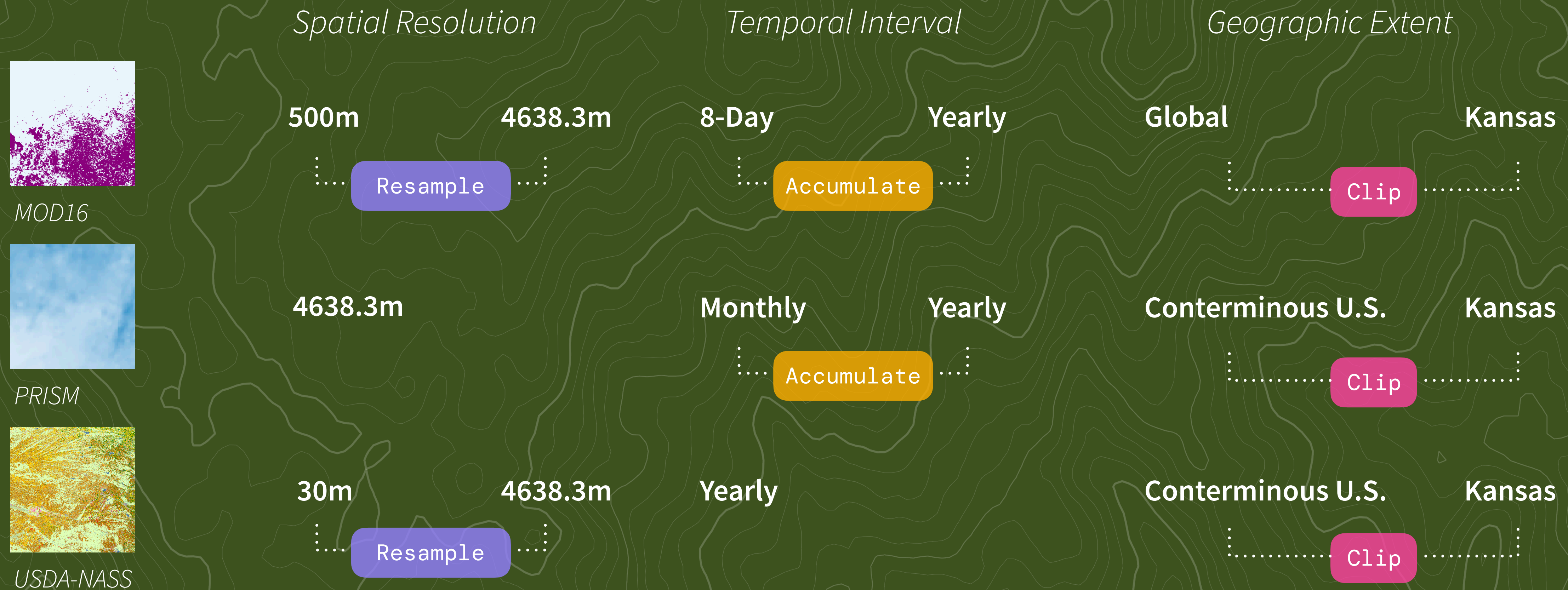
Aligning Geospatial Datasets

PE2's Task. Develop a model to predict groundwater withdrawal.



Aligning Geospatial Datasets

PE2's Task. Develop a model to predict groundwater withdrawal.



Aligning Geospatial Datasets

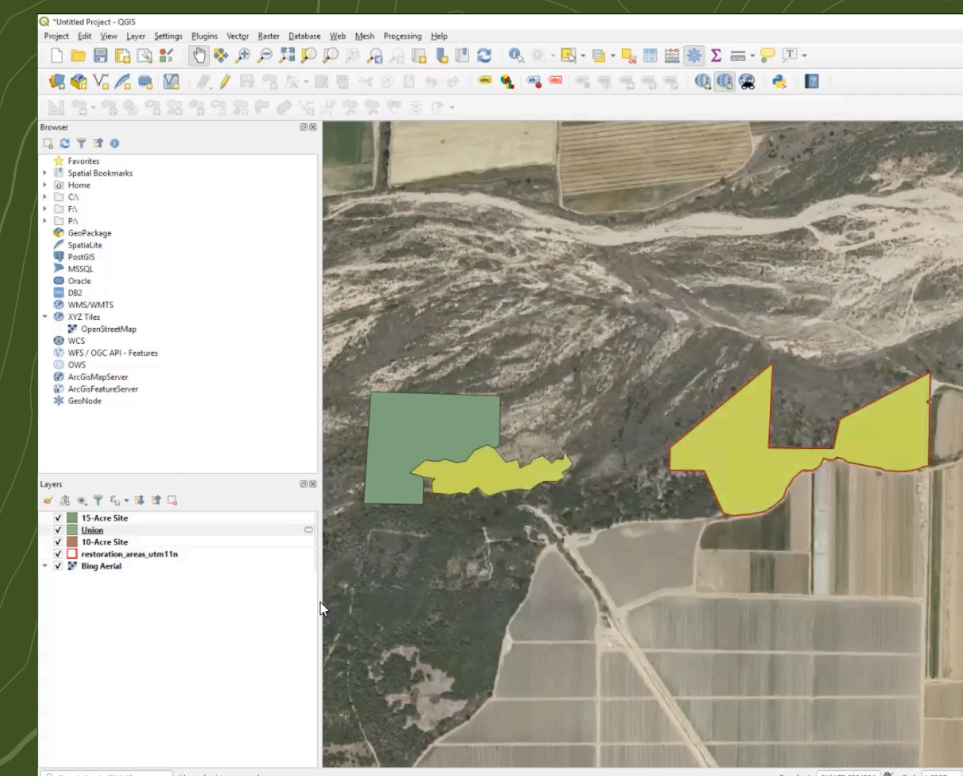
Aligning geospatial datasets required participants to have significant **fluency in geospatial data theory** in addition to **contextual information** about the datasets themselves.

Server Toolbox
Ready to Use Toolbox
Spatial Analyst Toolbox
Spatial Statistics Toolbox
... **+35 More**

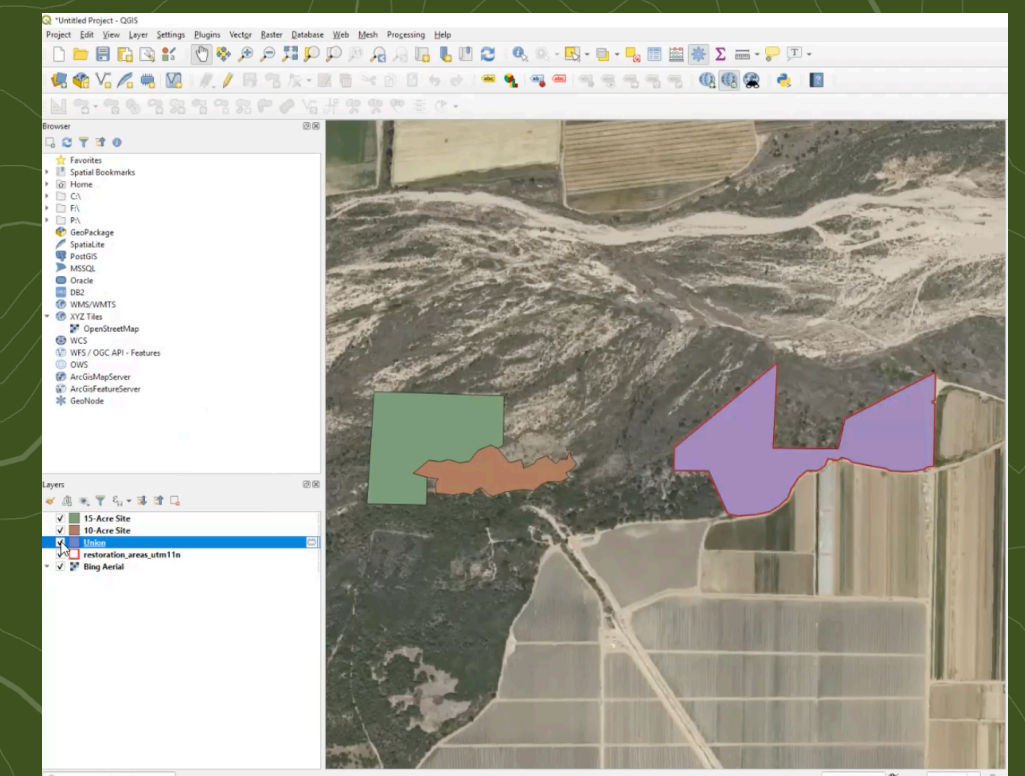
Bitwise Left Shift
Kriging
Raster Calculator
Iso Cluster Unsupervised
Fuzzy Overlay
Zonal Histogram
Darcy Flow
... **+200 More**

Identify the correct **sequence** of **transformations** among **hundreds of operators**

Expected



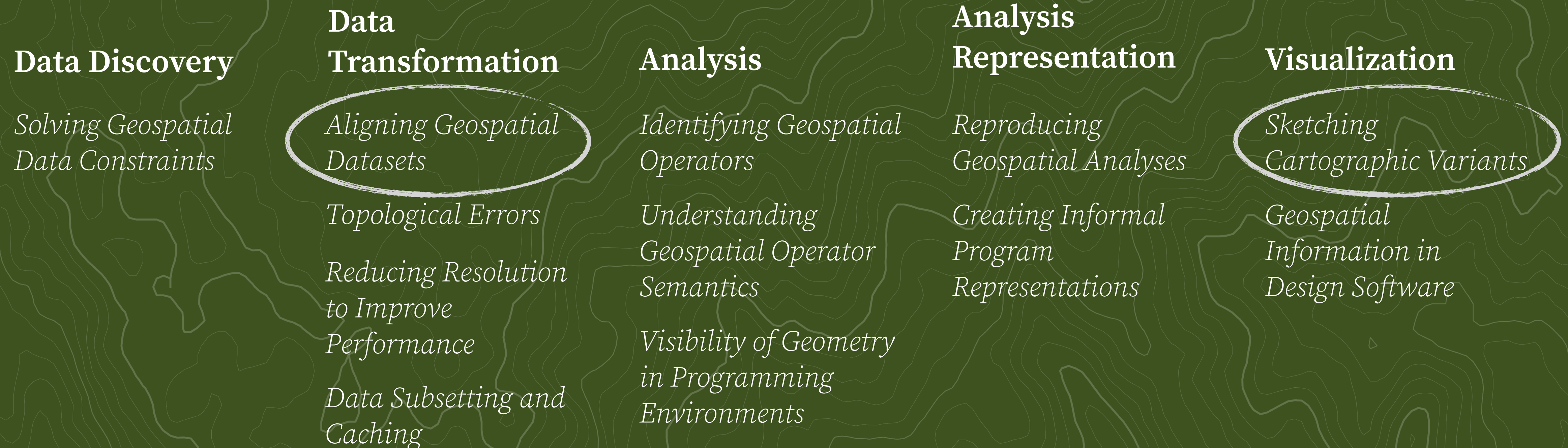
Actual



Determine when selected **transformations** produced **undesirable results**

Findings

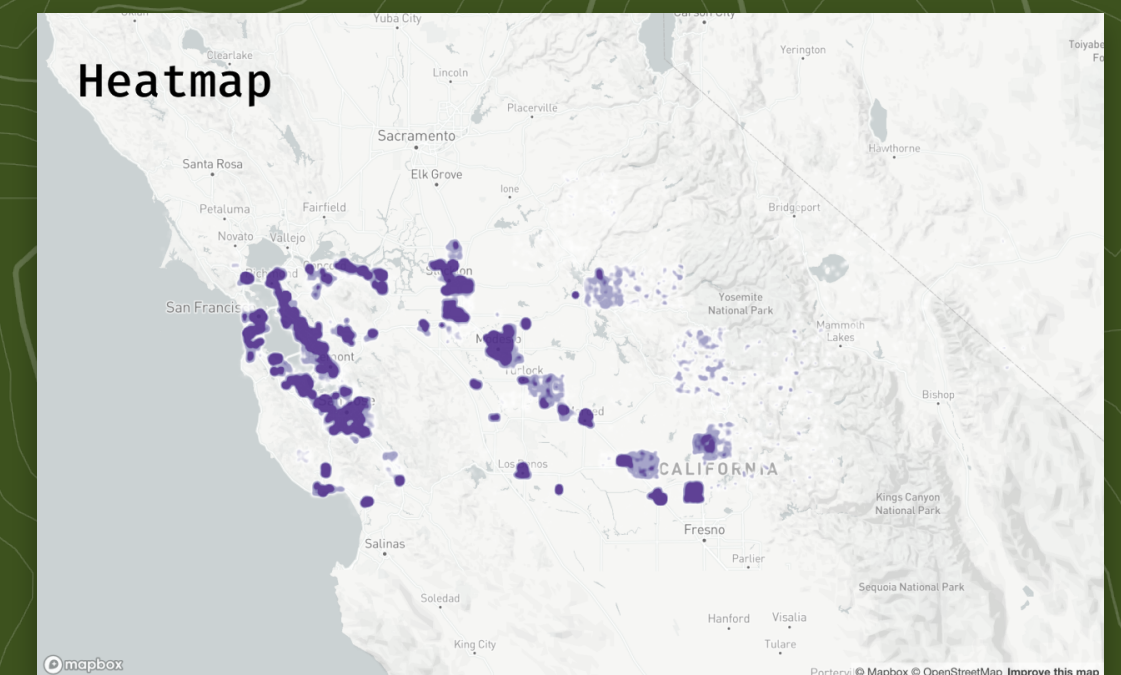
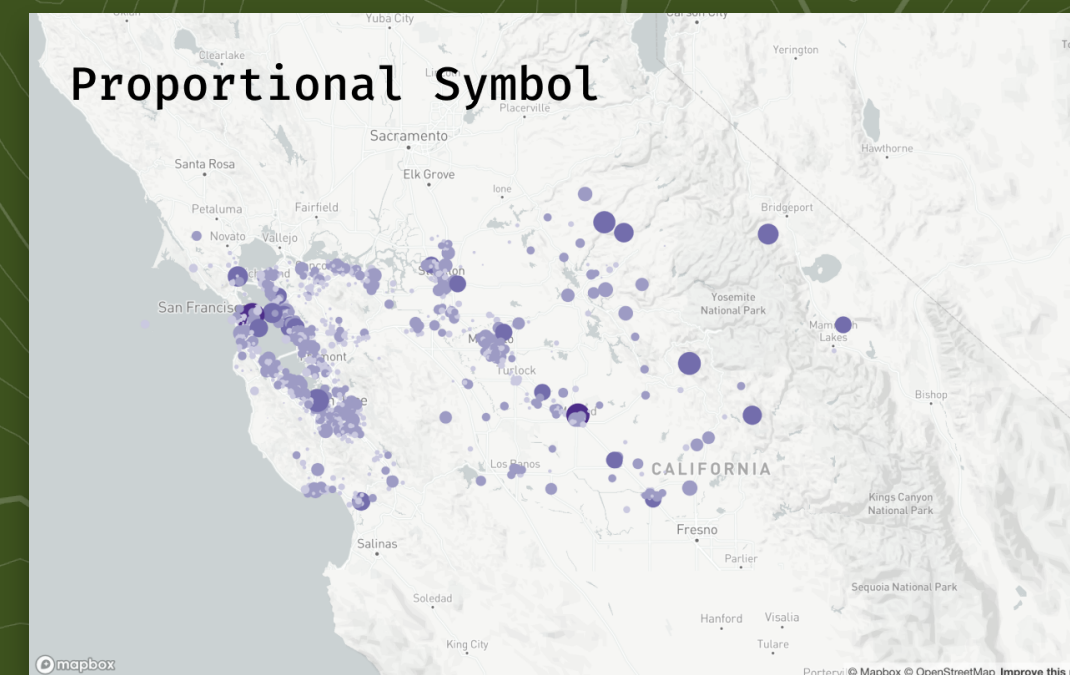
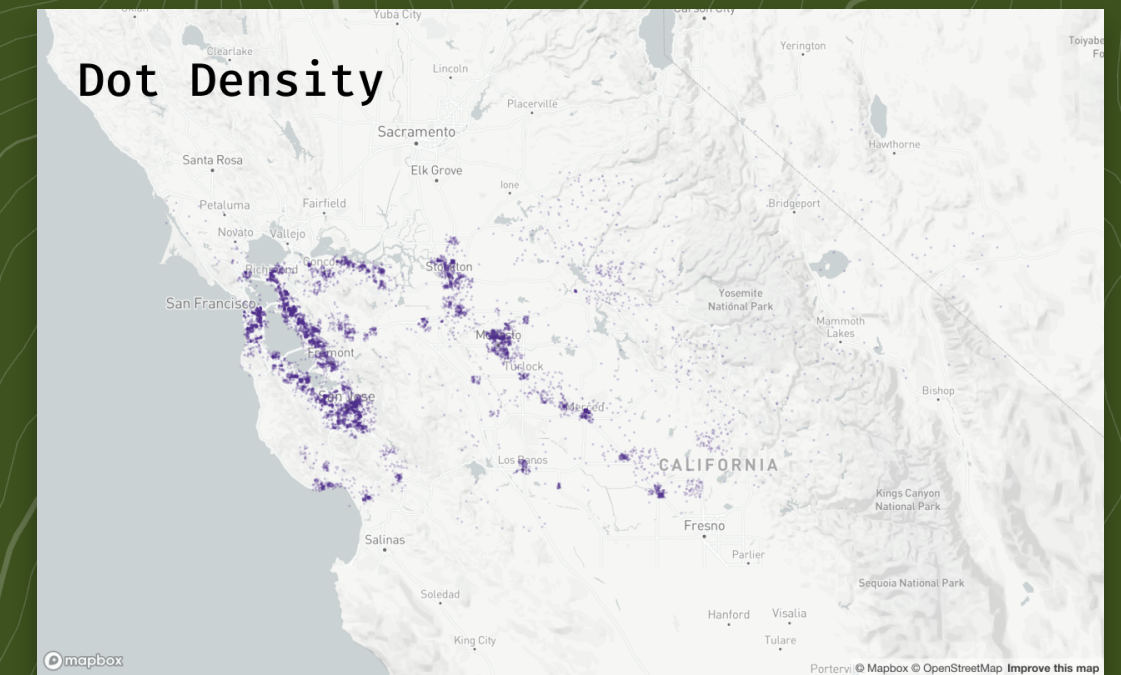
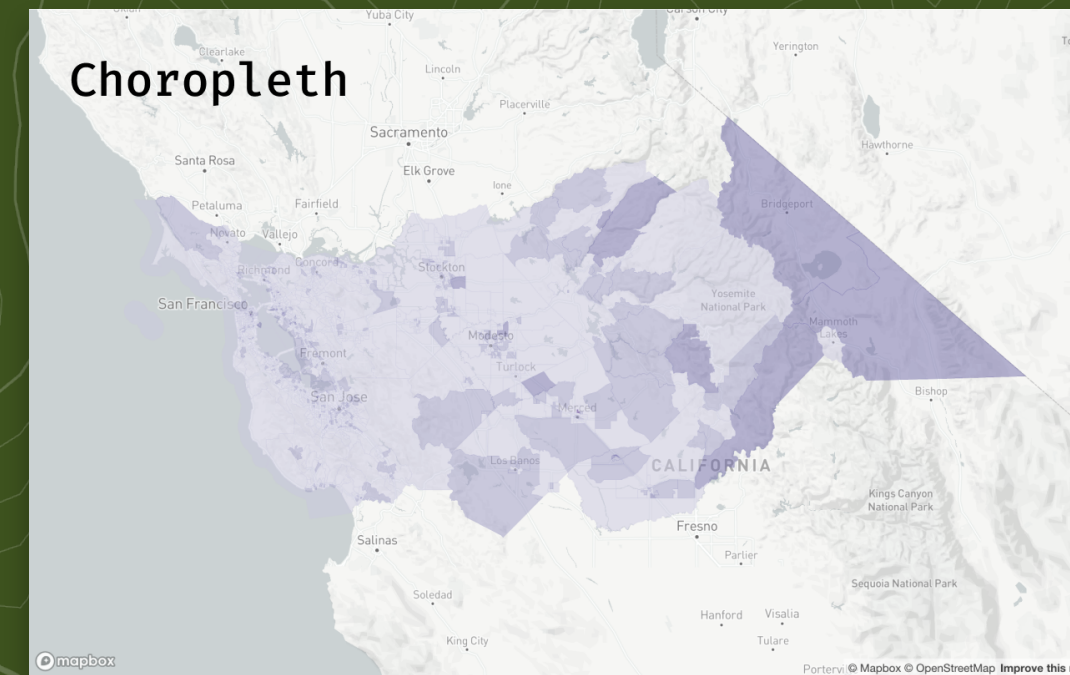
We identified **12 challenges** across **five phases** of participants' work with geospatial data.



Sketching Cartographic Variants

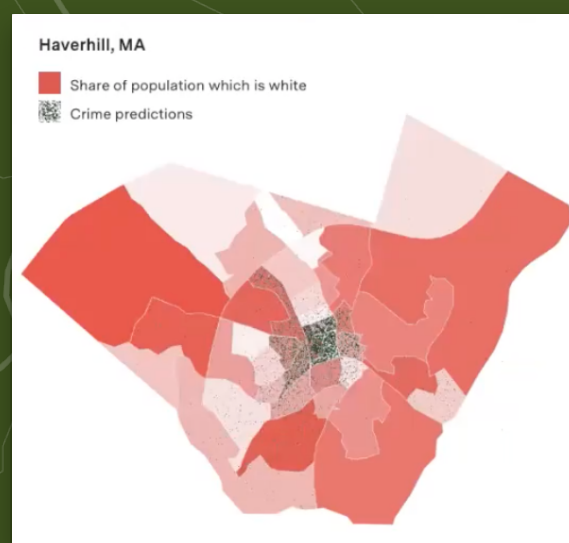
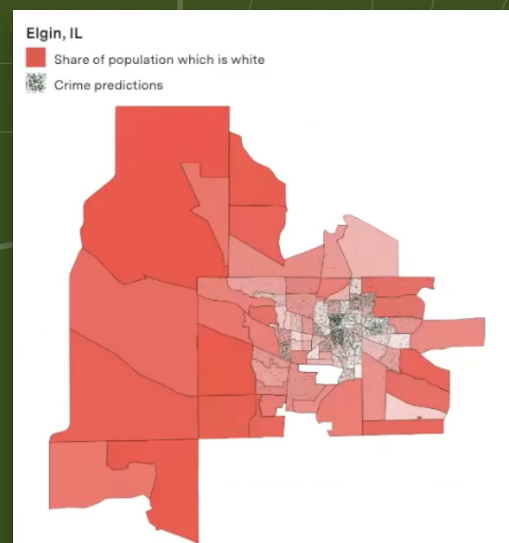
Participants wanted to visualize their data using many different **cartographic representations**.

- Identify the **map type** that represented their data most effectively
- Produce **tangible artifacts** for collaborators to evaluate

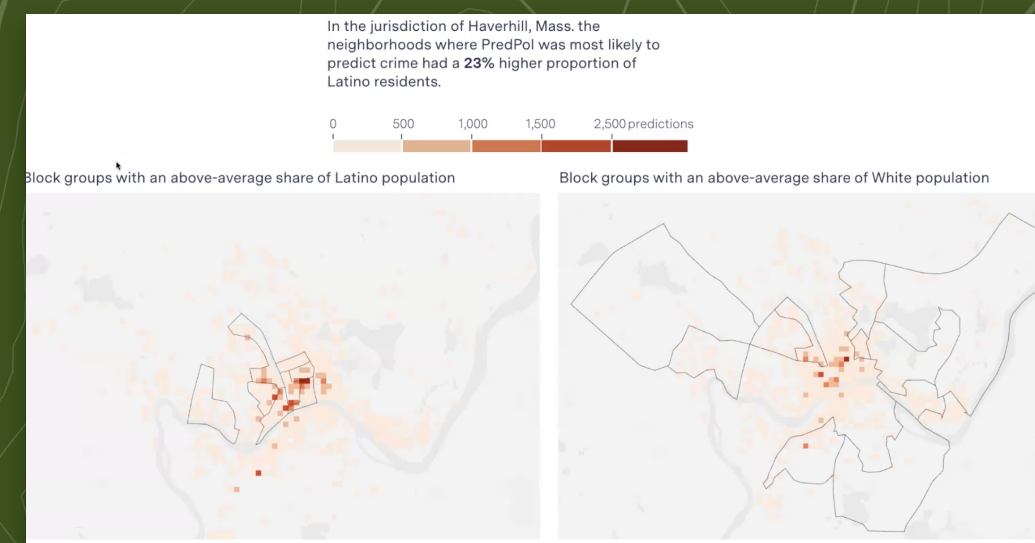
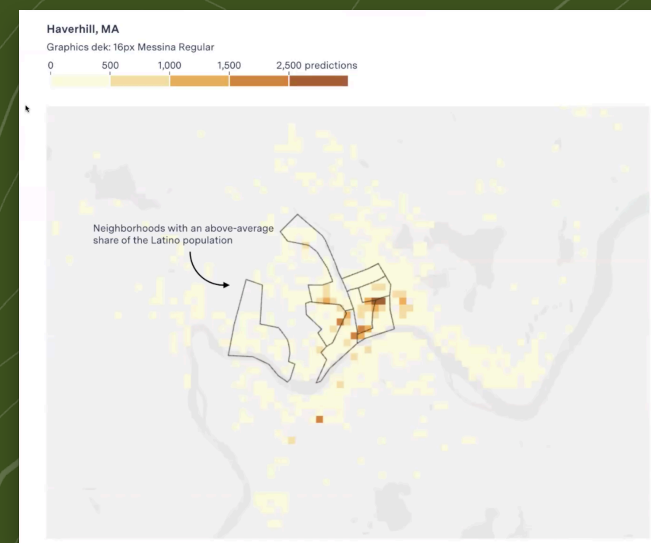


Sketching Cartographic Variants

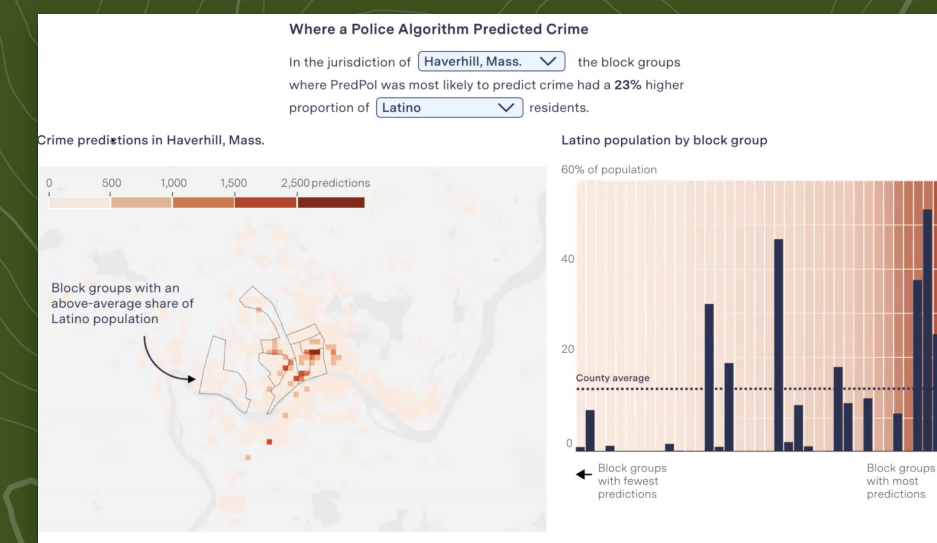
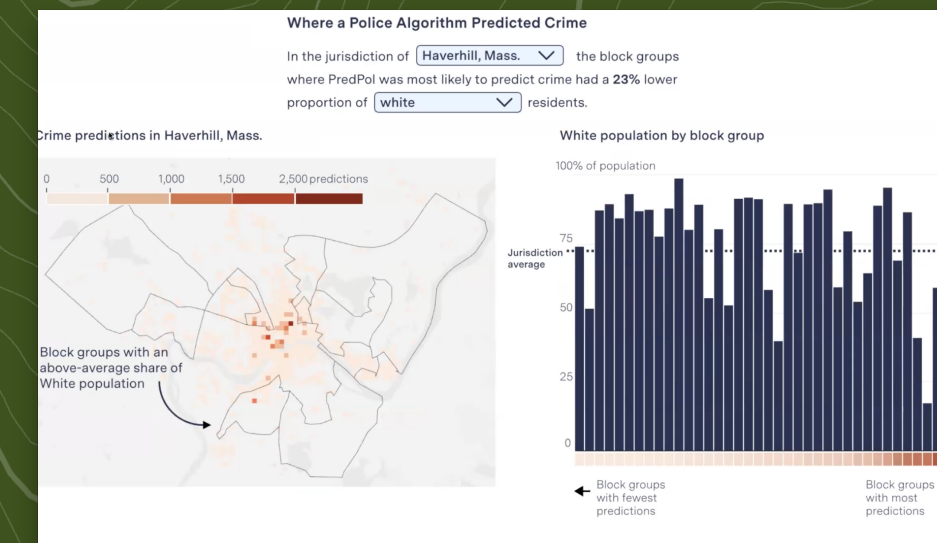
PJ5 created over **20 draft maps** for a story on biased predictive policing algorithms.



*Choropleth and
Dot Density*



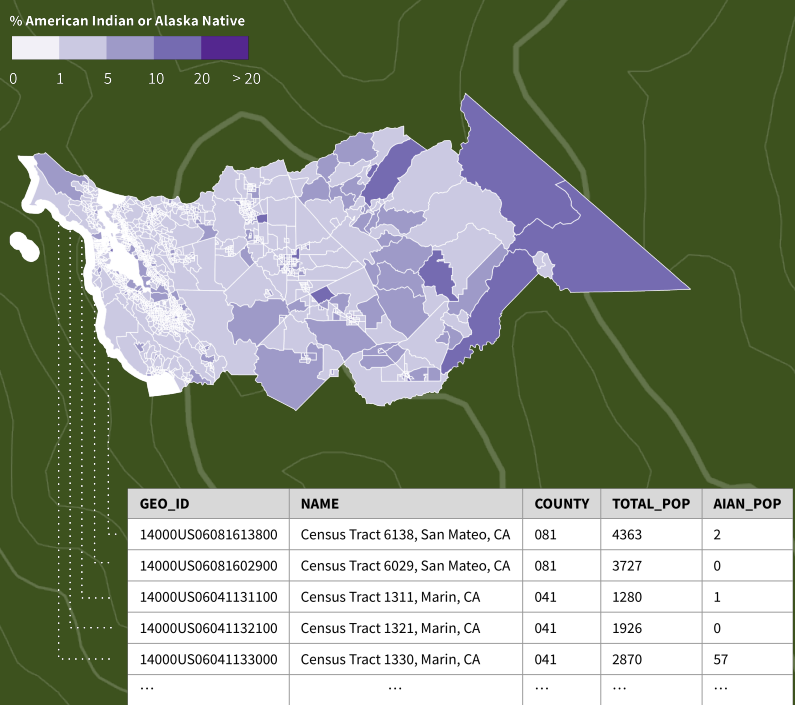
Gridded Heat Map



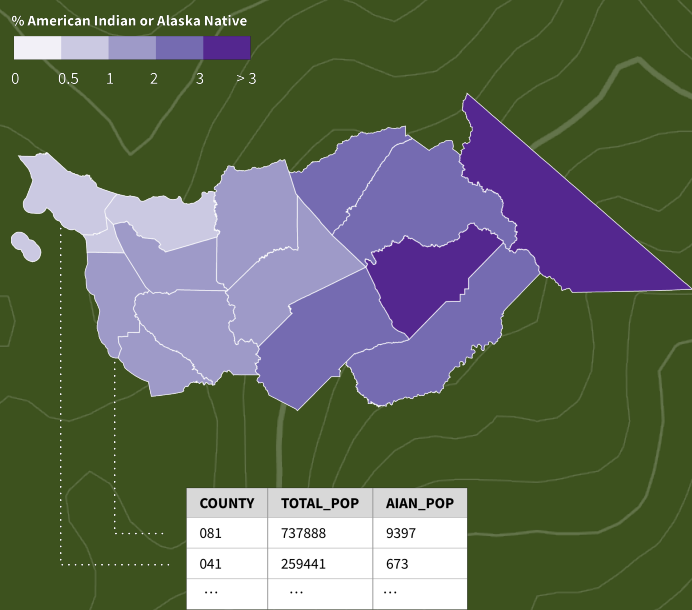
*Gridded Heat Map
with Bar Charts*

Sketching Cartographic Variants

Producing most map variants required going through **the entire analysis and visualization pipeline.**



Census Tracts



Counties

Additional Data Transformation

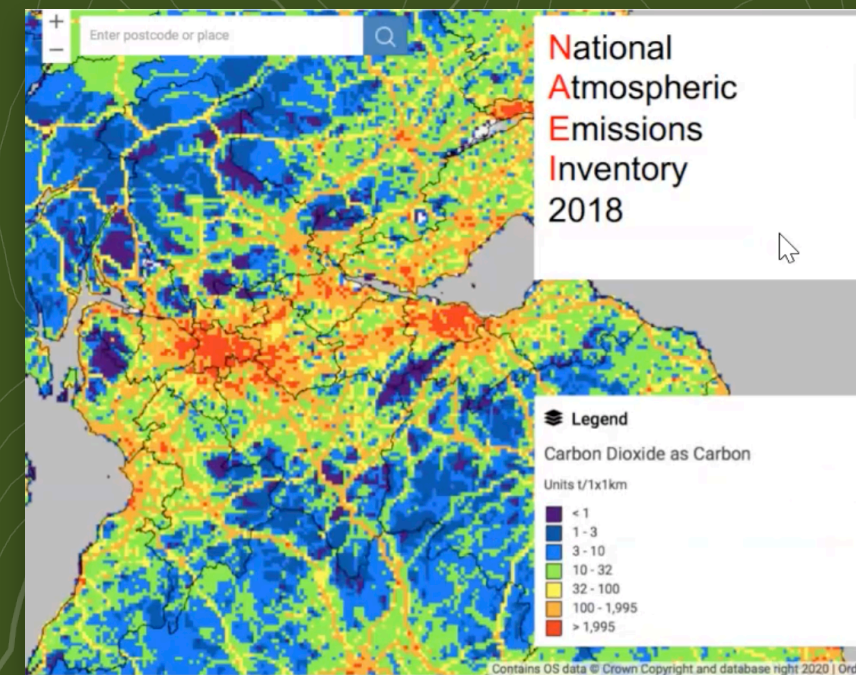


Across Multiple Tools

Sketching Cartographic Variants

Participants tried to **speed up** the drafting process in creative ways. One common technique involved **screenshotting in-progress maps**.

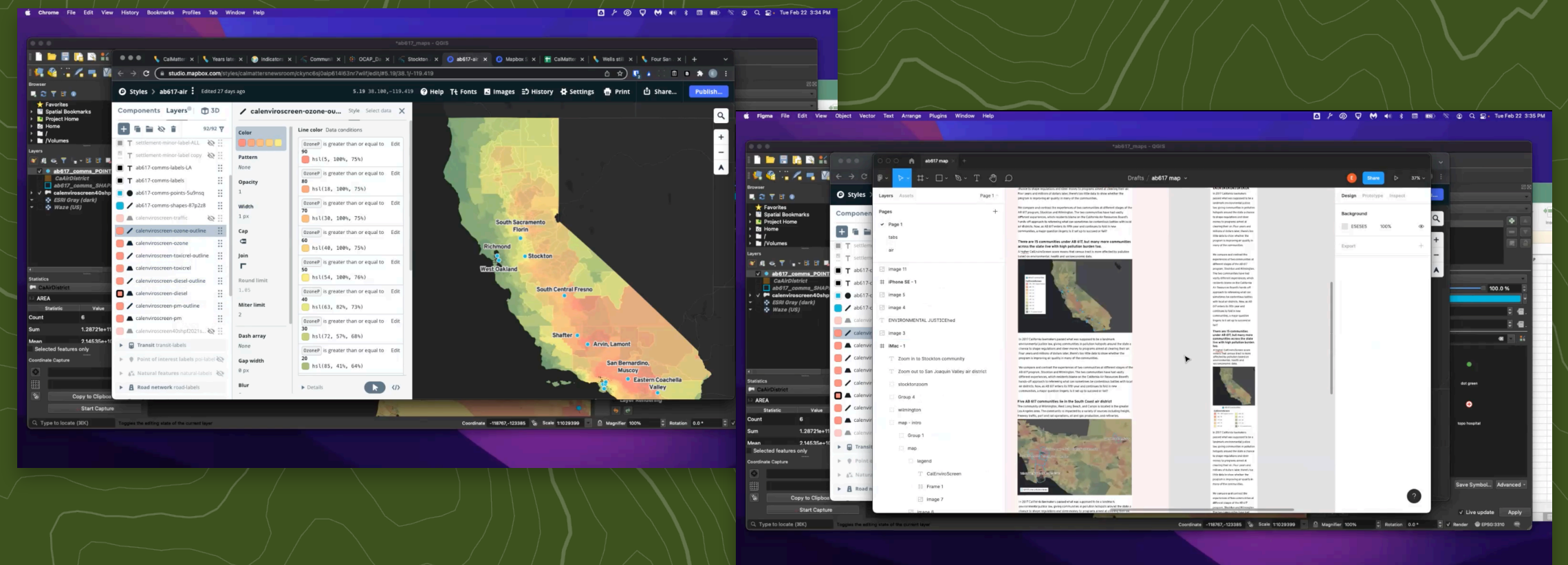
Participant E5



Participant S2

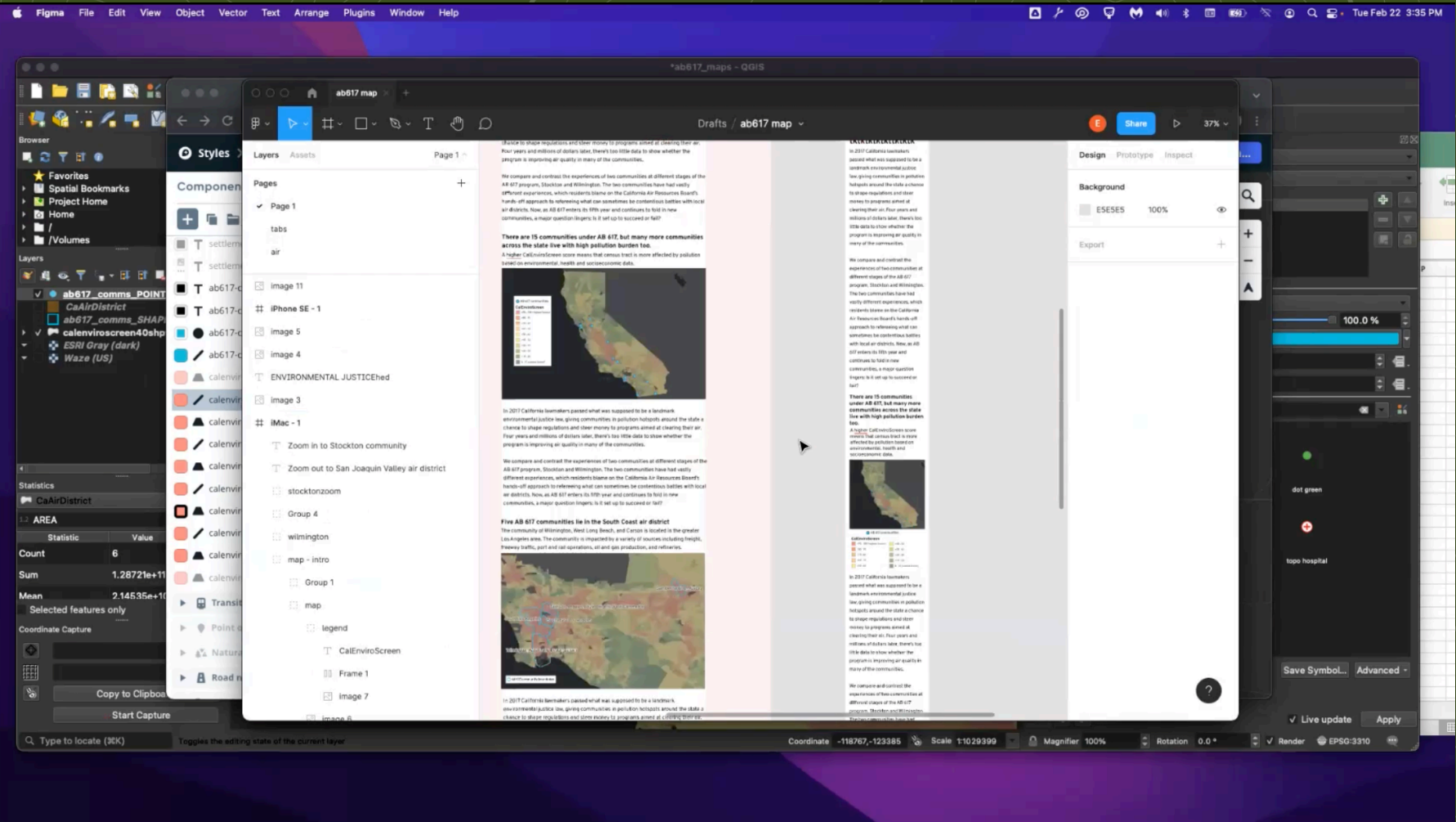
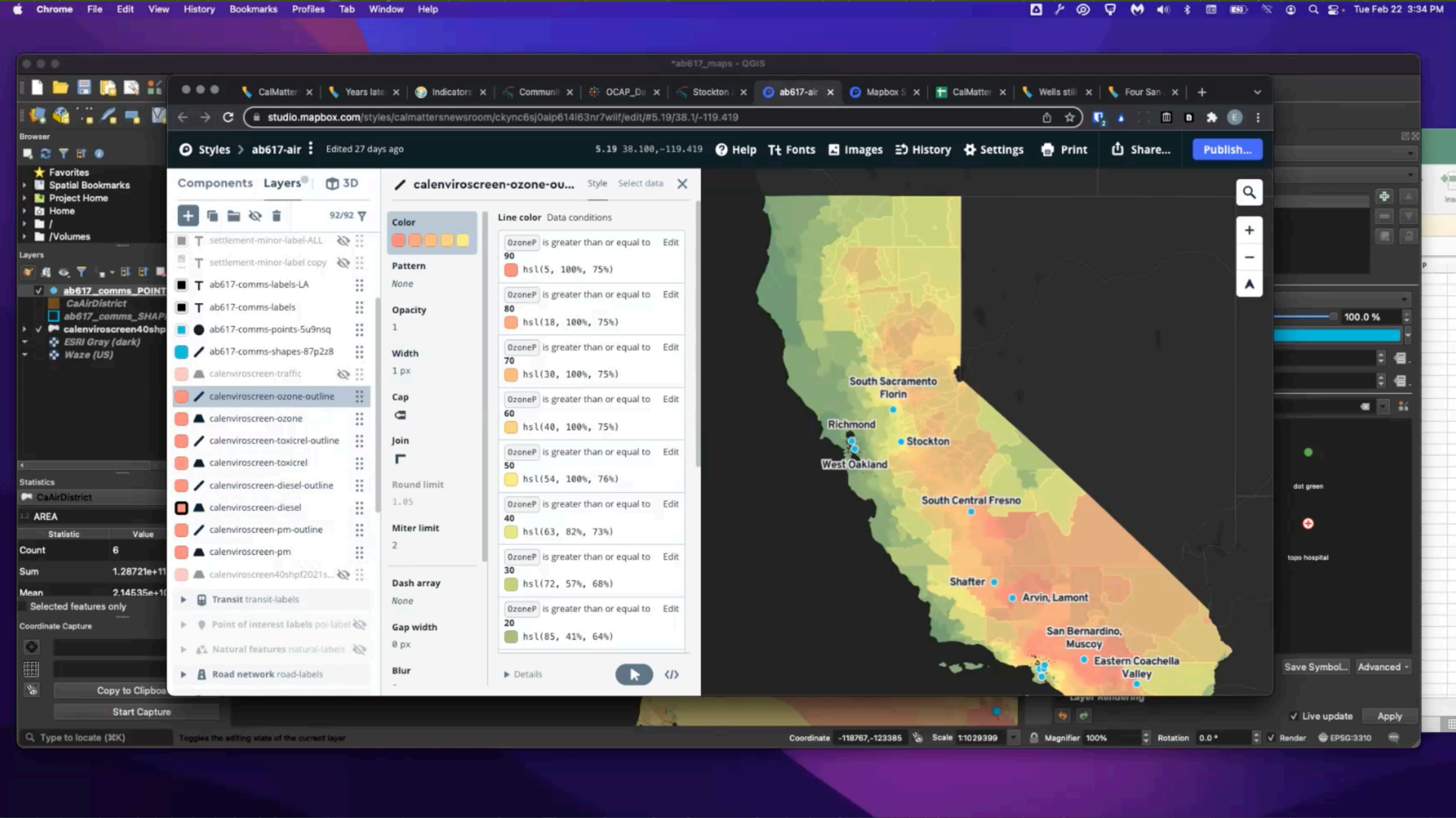


Participant J6



Sketching Cartographic Variants

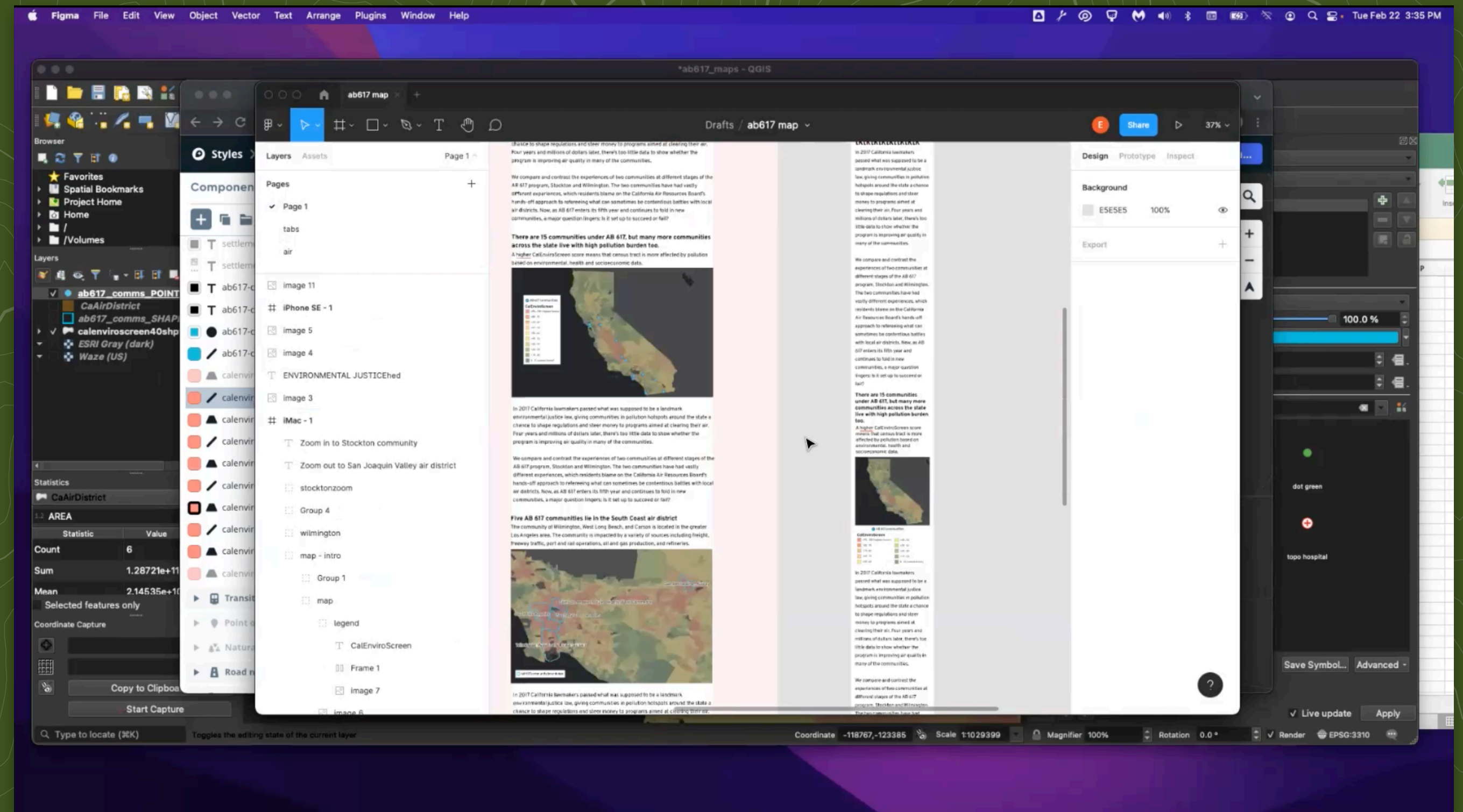
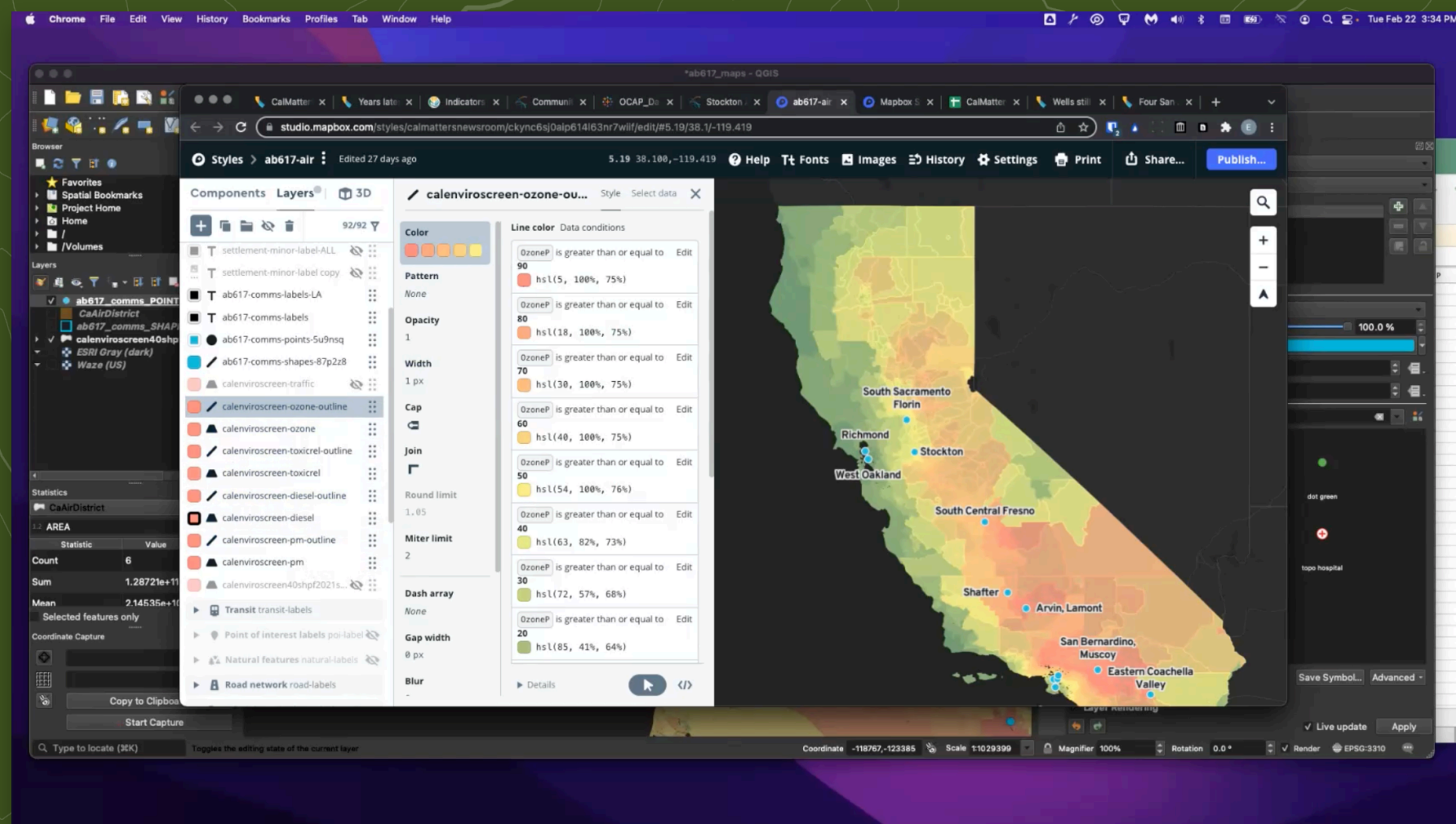
Screenshots



Sketching Cartographic Variants

Screenshots● Layouts

Allowed PJ6 to compare cartographic choices “before I code anything.”



Screenshotting came with limitations.

1. Only allowed users to capture cartographic changes **within** a map type rather than **across** map types

2. Once a final map design was chosen, participants had to **reproduce the selected draft in code**

Roadmap



Roadmap



Design Opportunities

We synthesized **six design opportunities** for designers and developers of geospatial analysis and visualization systems.

Solving Geospatial Data Constraints

Opportunity 1. Participants struggled to find geospatial data satisfying complex spatial and temporal constraints (Section 5.1). While many could describe their constraints succinctly, satisfying them involved constructing bespoke workflows to combine, align, and simplify their raw datasets (Section 5.2). These challenges suggest an opportunity for tools that (1) offer alternative programming abstractions to express data constraints and (2) infer geospatial data queries and transformations from constraints.

Assistive Tools for Constructing Geospatial Analysis Pipelines

Opportunity 2. Participants could describe the target outputs of their geospatial analyses but struggled to construct pipelines to produce them (Section 5.3). This suggests an opportunity for tools that (1) accept non-code specifications of analysis intent, (2) synthesize analysis programs that satisfy specifications, and (3) support users in editing programs.

Opportunity 3. Participants relied on running operators and manually inspecting outputs to understand operator semantics (Section 5.3.2). This was computationally expensive and time-consuming, suggesting an opportunity for tools that surface information on operator semantics without requiring execution across entire inputs.

Reproducible, Shareable Geospatial Workflows

Opportunity 4. Participants using GISs struggled to create reproducible, shareable geospatial workflows (Section 5.4.2). Limitations in existing history interfaces made it difficult to recover information on the current analysis state or revisit past analysis decisions (Section 5.4.1). These struggles suggest opportunities for tools that (1) support efficient search through system history and (2) distill history into a portable and executable representation.

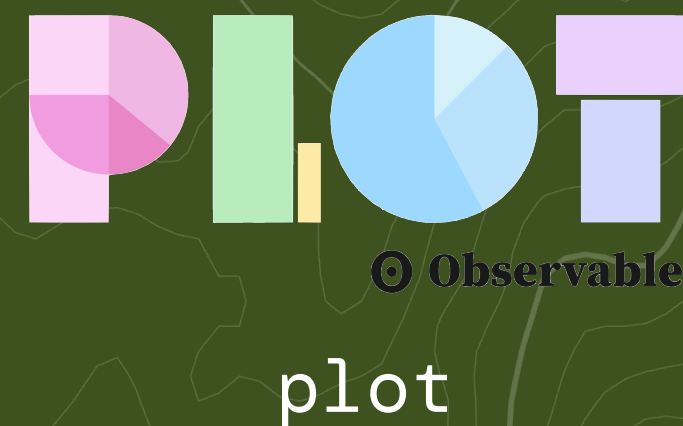
Exploring the Cartographic Design Space

Opportunity 5. Participants wanted to visualize their geospatial data using multiple cartographic representations, but transitioning between representations required engineering each one from scratch (Section 5.5.1). This suggests an opportunity for cartographic design tools that reduce the viscosity [8] of switching between map types.

Opportunity 6. Many participants used direct manipulation design software to visualize geospatial data. These tools discard all geographic information, making it difficult to refactor an analysis once visualization work has begun (Section 5.5.2). This suggests an opportunity for tools that (1) bridge geospatial analysis and cartographic design and (2) maintain the underlying geospatial data representation of graphical elements while supporting direct manipulation.

Design Opportunities

Opportunity. Cartographic design tools could focus on **reducing the “viscosity” of map type transitions.**

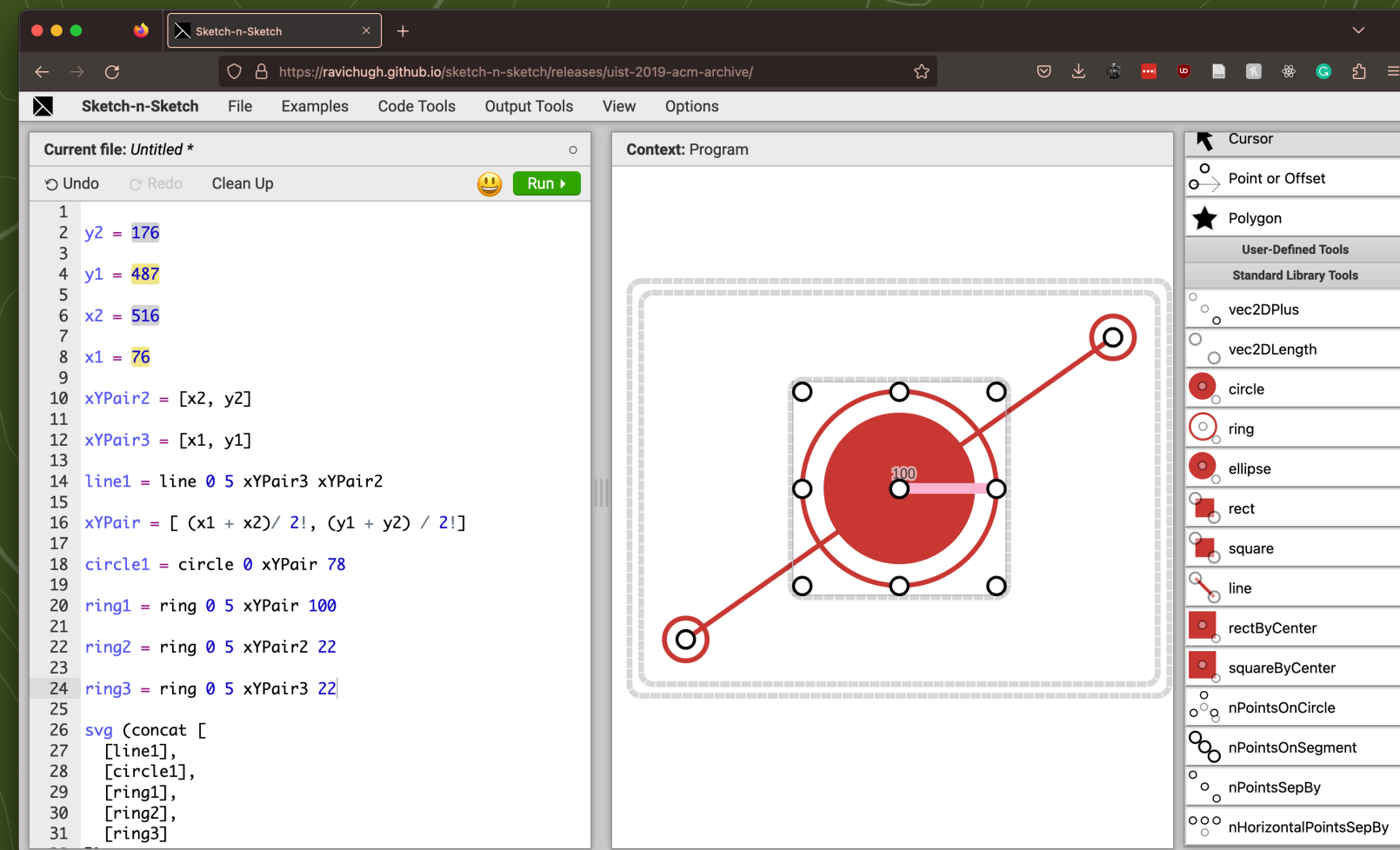


- Restrict **geospatial file formats, data models, and map types**
⇒ Could not express many of the maps participants made

Possible Solution. *Grammar of Graphics*

Design Opportunities

Opportunity. Cartographic design tools could **pair programmatic and direct manipulation paradigms** for map construction.



Sketch-n-Sketch

- Edit **source** or **output** and propagate edits **bidirectionally**
⇒ Design maps using **direct manipulation** while giving access to **program representations**

A Need-Finding Study with Users of Geospatial Data



Parker Ziegler

peziegler@cs.berkeley.edu
<https://parkie-doo.sh/>



Sarah E. Chasins

schasins@cs.berkeley.edu



Learn about all **12 challenges**, all **six design opportunities**, and hear from our participants in the paper.

CHI' 23 • Working with Data • April 25, 2023

Berkeley
UNIVERSITY OF CALIFORNIA

