# Black Freedom Mapping: Application of Artificial Intelligence (AI) to Construct a National Database of Black Mobility During the Jim Crow Era

Sydney Moseley, BS[1], Pedro Pinheiro-Chagas[1,2,3], PhD, Muriel Taks Calle, MA[4], Paris B. Adkins-Jackson, PhD, MPH[4,5], & Tanisha G. Hill-Jarrett, PhD. [1,6]

[1]Memory and Aging Center, Department of Neurology, University of California San Francisco, [2]Bakar Computational Health Sciences Institute, University of California, San Francisco, [3]Center for Intelligent Imaging (ci2), University of California, San Francisco, [4]Department of Epidemiology, Mailman School of Public Health, Columbia University, [5]Department of Sociomedical Sciences, Mailman School of Public Health, Columbia University, [6]Global Brain Health Institute, University of California San Francisco

UCSF Weill Institute for Neurosciences
Memory and Aging Center

## Background

**Historical Context:**
- During the Jim Crow era, traveling while Black meant possible life-threatening danger due to widespread racial discrimination and the threat of violence in sundown towns.
- The Motorist's Green Book, published annually from 1937 to 1964 by Harlem postal worker Victor Green, served as a vital guide for Black travelers seeking safe spaces—i.e., businesses and residences—while on the road.

**Dataset Importance:**
- While research on how structural racism shapes the health of diverse populations has grown, there remains critical data gaps that prohibit the use of historical data to examine associations between structural racism and population health.
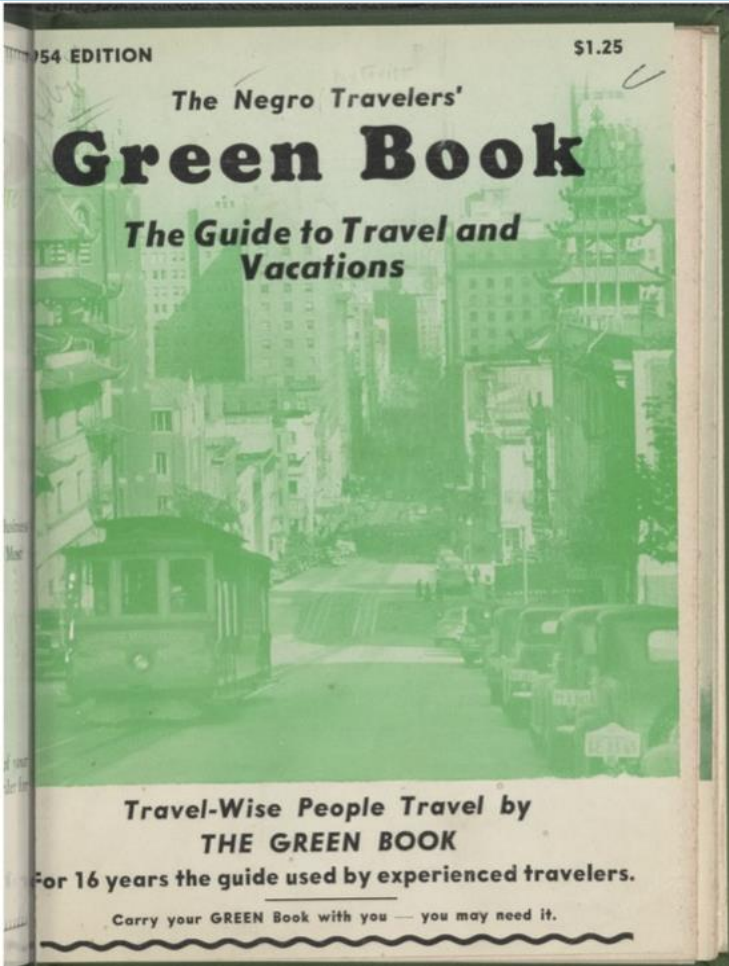
**Objective:**
- By extracting addresses from digitized editions of The Green Book, we aim to provide a novel database of Black mobility, while describing a reproducible method for context-aware data extraction.


Figure 1. 1954 Edition Cover


Figure 2. Page from the 1940 Green Book


Figure 3. Page from the 1954 Green Book

## Method

**Introduction:**
- In order to create this database, one would typically manually convert the Green Books to text, an incredibly time-consuming task, or implement Optical Character Recognition (OCR).
- OCR uses pattern matching or feature extraction algorithms to convert images of text characters into a machine-readable format.
- While traditional OCR is imprecise in extracting unstructured and nonregular expressions from images of varying quality, this study introduces a novel application of Artificial Intelligence (AI) to extract historical addresses in inconsistent formats.
- We use a multimodal large language model (MLLM), which is an AI model that extends the capabilities of LLMs by processing multiple types of data, such as images, text, and audio.

## Method [cont'd]

**Utilizing AI:**
- To extract the addresses from the Green Books, we implemented structured prompting and knowledge retrieval prompting to extract and standardize historical address data from scanned PDFs.
- PDFs of the books were converted to images, then passed through GPT-4o, a multimodal large language model (MLLM), used to extract addresses from twenty-one Green Books and convert them to text.
- A second call to the model used structured outputs to constrain the responses to our specified JSON schema.
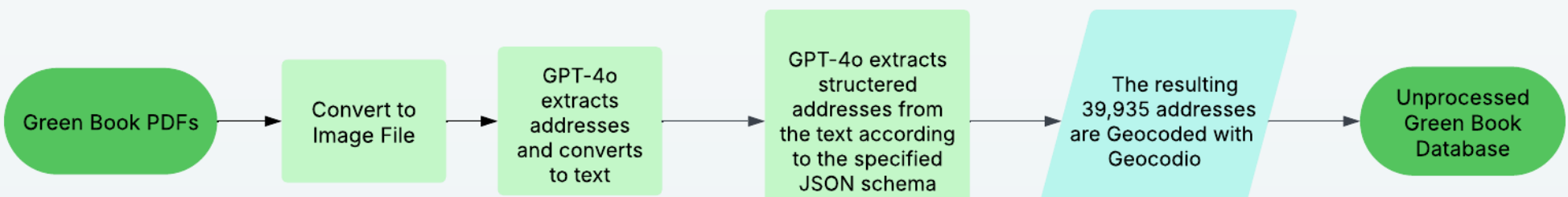

Figure 4. Generalized workflow of using GPT-4.o to extract addresses

**Geocoding and Post-Processing:**
- Forward geocoding of the extracted Green Book addresses was completed in Geocodio, creating latitude and longitudinal coordinates for each address. These geospatial coordinates were linked to contemporary census data based on the 2000 U.S. Census.
- Accuracy Check: Fuzzy string matching with a threshold of 60% similarity was used as the accuracy cutoff for comparison between Geocodio's parcellation and the extracted address.
- With the remaining addresses not identified in Geocodio, Python pandas and geopy libraries were used to identify coordinates. These coordinates were matched to the U.S. Census Bureau's 2024 TIGER data using Nominatim.
- The remaining missing coordinates were manually examined for error types and corrected if applicable.
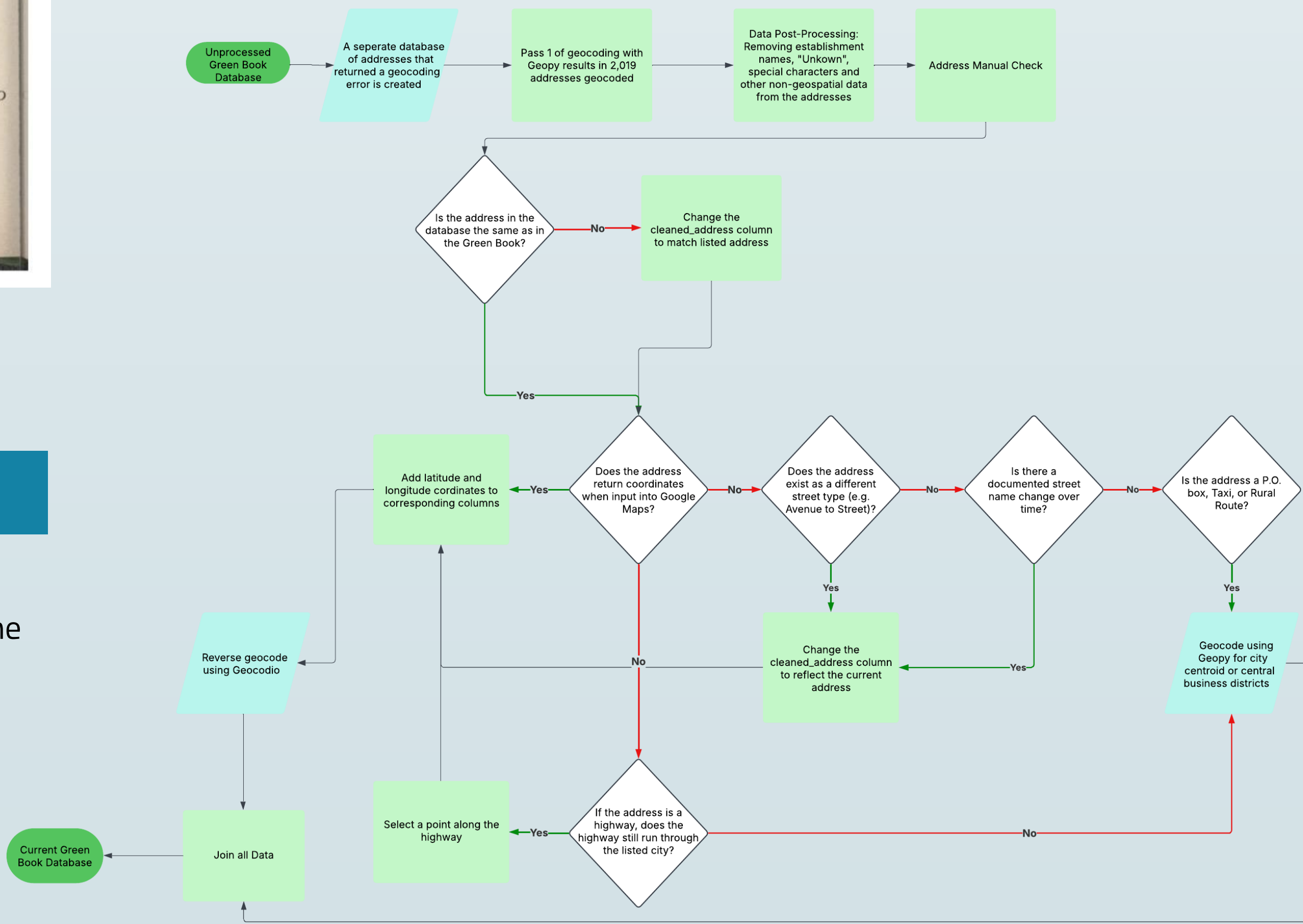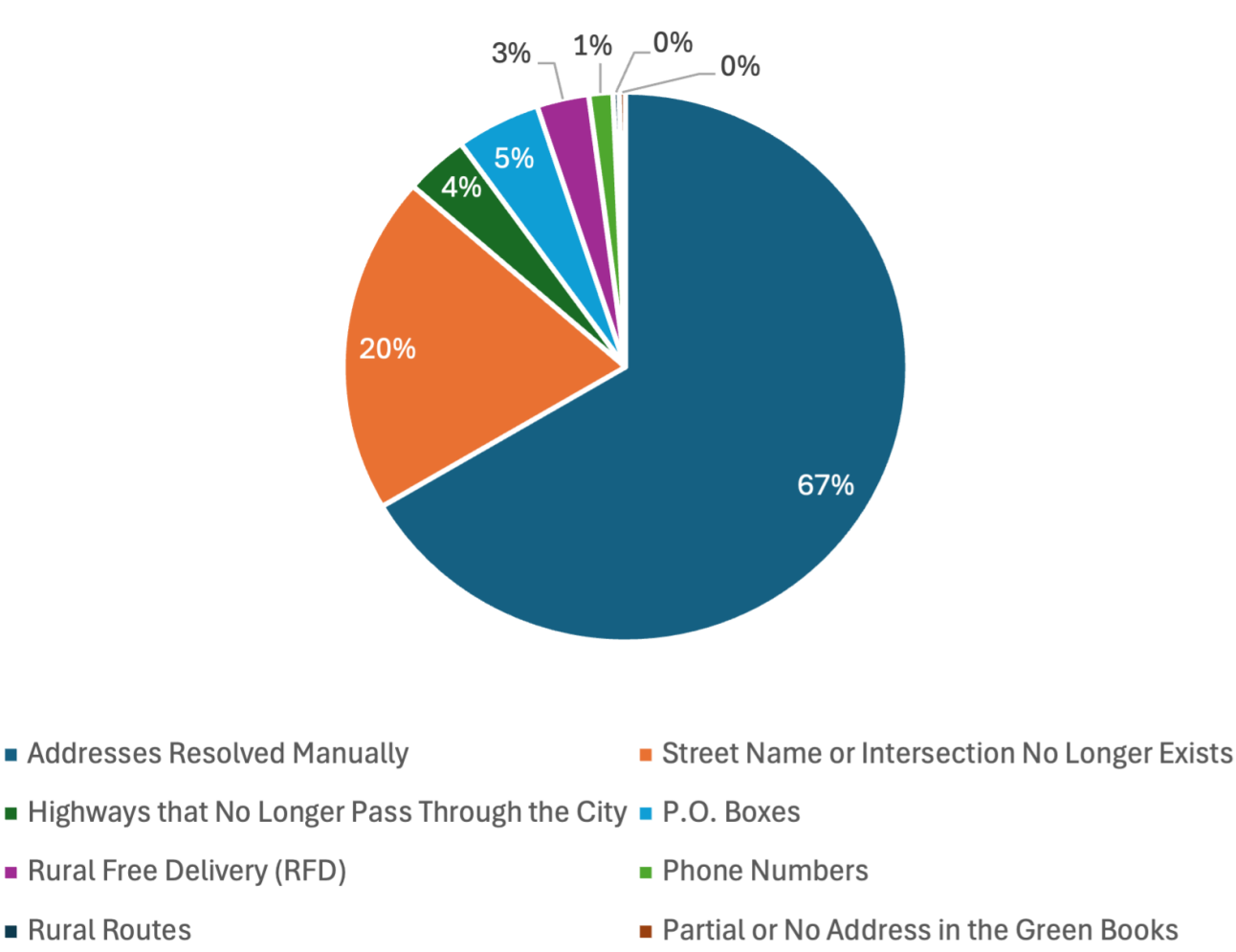

Figure 5. Generalized workflow of the geocoding process

## Results

- A total of 39,935 addresses were extracted from Green Books across all years. Across all addresses extracted, 11.4% (N = 4,562) were unable to be forward geocoded due to preprocessing errors and Geocodio's lower tolerance level for incomplete or ambiguous addresses.
- Of the remaining addresses, 23.6% (N = 8,350) exceeded the threshold for address inaccuracy and 10.5% (N = 3,717) of addresses were extracted without a street name.
- A total of 44.3% (N = 2,019) of addresses with missing latitude and longitude coordinates were resolved via TIGER data. The remaining 55.7% (N = 2,543) of unresolved addresses revealed common errors including: incomplete addresses, changes in road type, digitization errors, and address changes over time.


Figure 6. Types of Addresses Unable to be Machine-Geocoded

Legend: Addresses Resolved Manually (67%); Highways that No Longer Pass Through the City (20%); Rural Free Delivery (RFD); Rural Routes; Street Name or Intersection No Longer Exists; P.O. Boxes; Phone Numbers; Partial or No Address in the Green Books

- OCR (N = 1,033) of these addresses were manually resolved by updating AI extraction errors, spelling, and updating differences in road type. Of all the extracted addresses, 2.1% (N = 851) were geocoded by city due to the address no longer existing or incomplete address information (e.g., a P.O. Box listed).


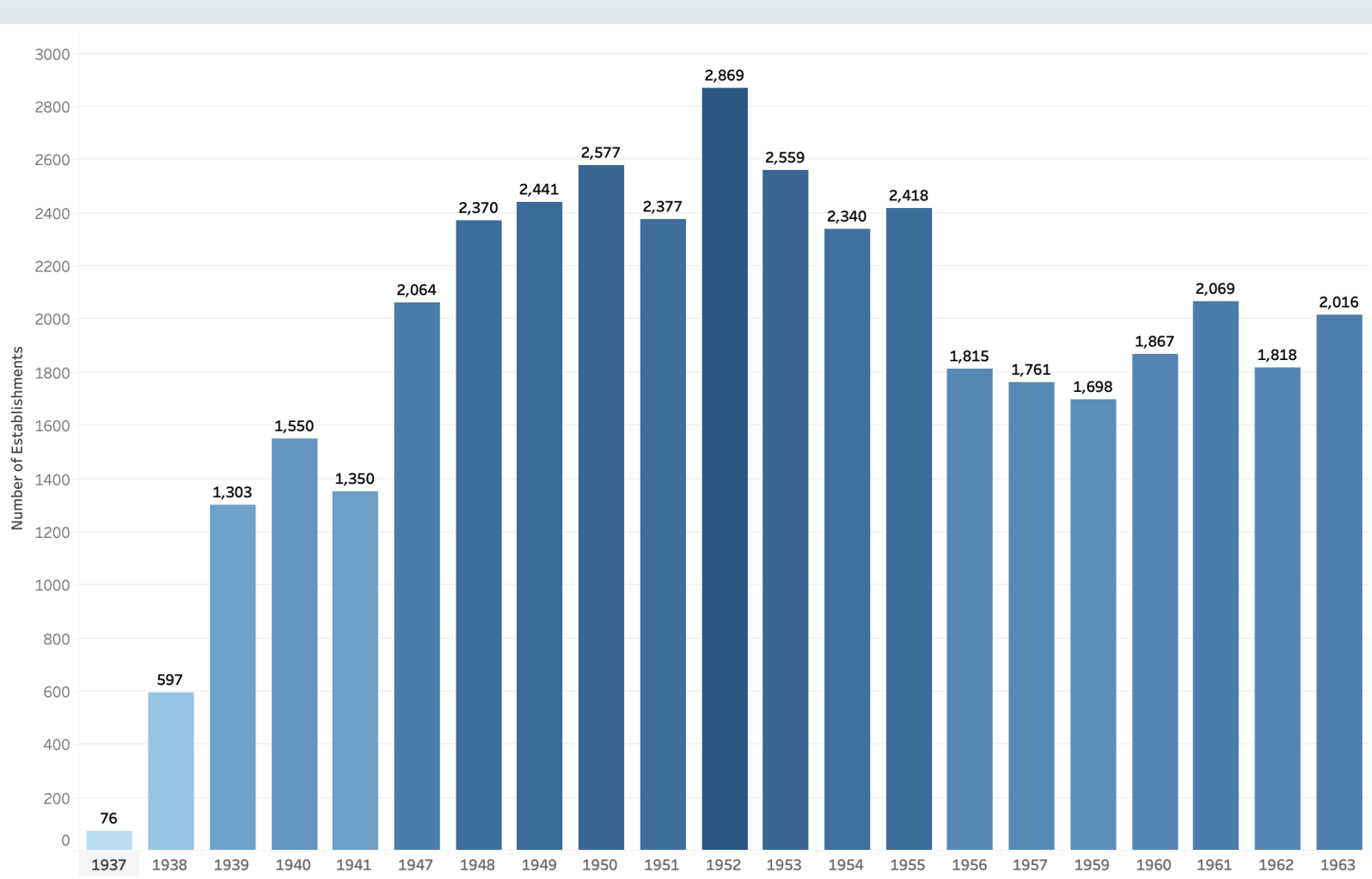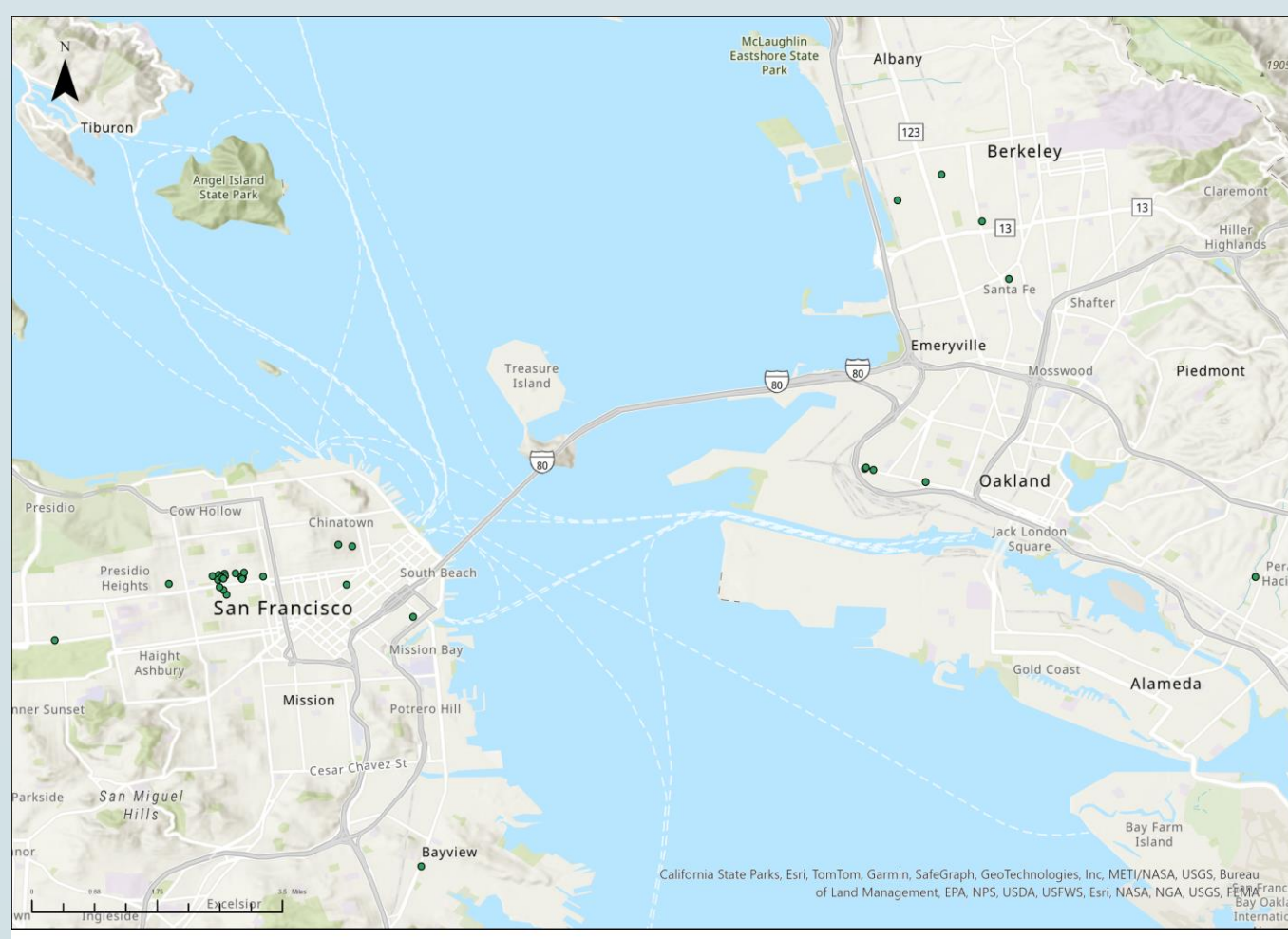Figure 7. Number of Green Book establishments listed per year


Figure 8. Map of Bay Area Green Book establishments (1950)

## Discussion

**Challenges:**
- Errors in extracting city names due to page or column breaks, which resulted in "Unknown" or the establishment name being returned rather than the city.
- Special characters (¬Ω, Äô, Äù, Äù) were extracted in place of quotation marks and dashes in some addresses.
- Accounting for incomplete addresses, such as addresses containing P.O. boxes, taxicab phone numbers, and Rural Free Delivery (RFD).
- Working with a limited availability of block-group and block-level geospatial data for historical addresses to account for address changes over time.
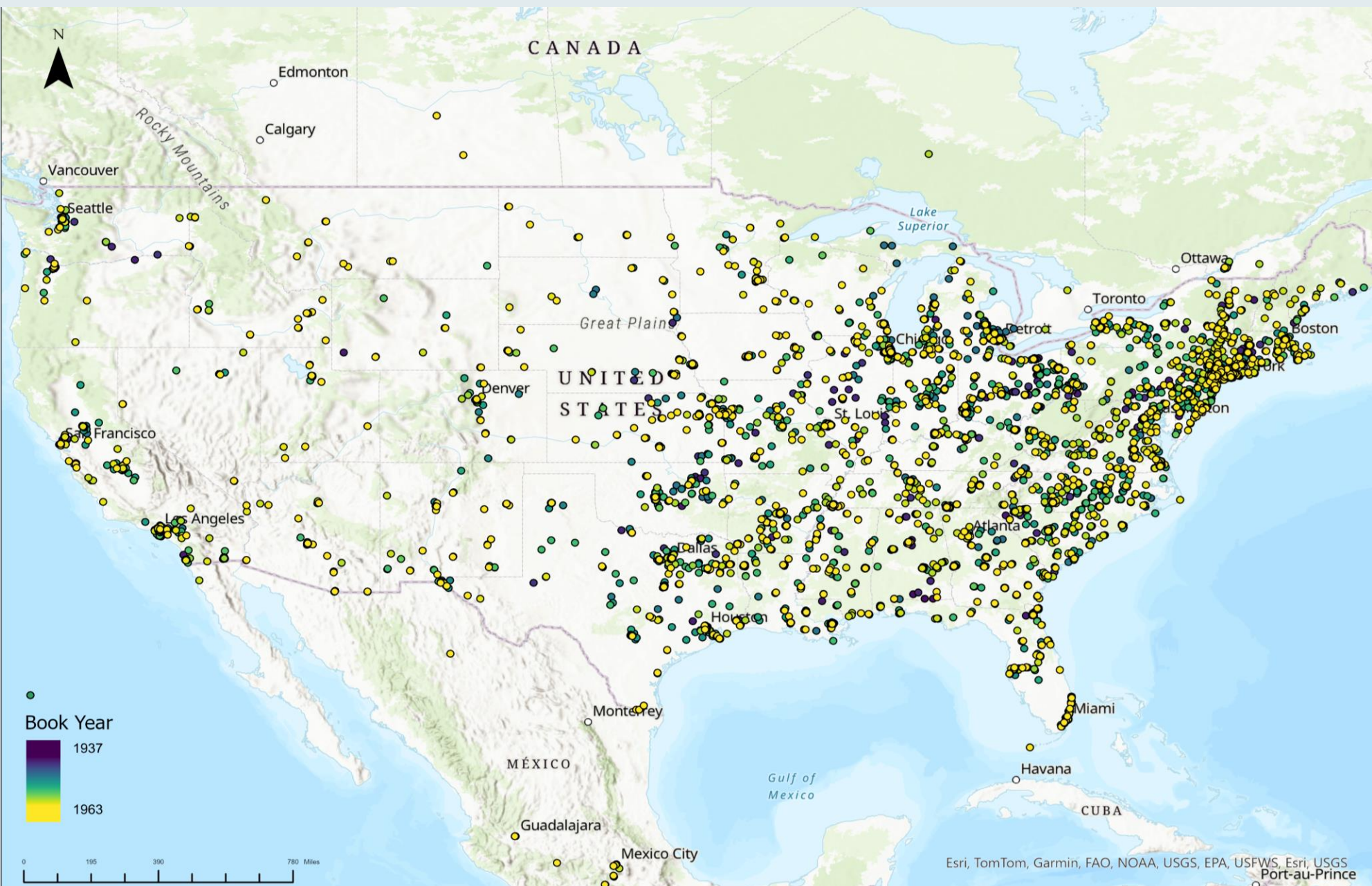
**Next Steps:**
- We plan on integrating business names and types to the database using few-shot learning to increase the consistency of model performance.
- Review addresses that failed fuzzy string matching to verify correctness and investigate discrepancies, which stem from the challenges listed above.

**Applications and Future Work:**
- Redlining and segregation associations with cognitive aging
- Combine with sundown town data to analyze stress exposure and terror on cognitive aging
- Food deserts and environmental hazards/pollution
- Economic mobility, urban renewal & gentrification/migration
- Loss of sense of place and community, time loss
- Access to transportation networks and commuting


Figure 9. Map of all Green Book establishments (1937-1963)

## Conclusions

- This is a seemingly feasible, efficient application of AI that allows us to begin to explore more nuanced associations between aspects of structural racism and health.
- Our next steps include geospatial linkages to the cognitive aging data in the Health and Retirement Study.
- We will investigate how Black countermapping practices and access to Black spaces during the Jim Crow era, designed in response to structural racism, influence cognitive aging trajectories and risk of diseases of aging over time.