

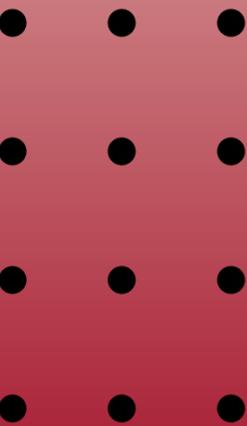
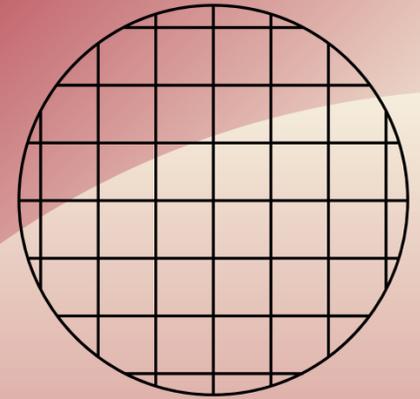
ML Augmented Prediction for Labor Exploitation Detection



x



August 14, 2024



Our Team



Enkhjin Munkhbayar
DSSG Fellow



Leon Reilly
DSSG Fellow



Kyler Shu
DSSG Fellow

Our Mentors



Dr. Benjamin Seiler
Technical Mentor



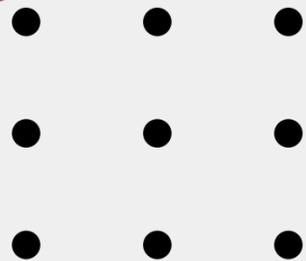
Dr. Mike Baiocchi
Faculty Mentor

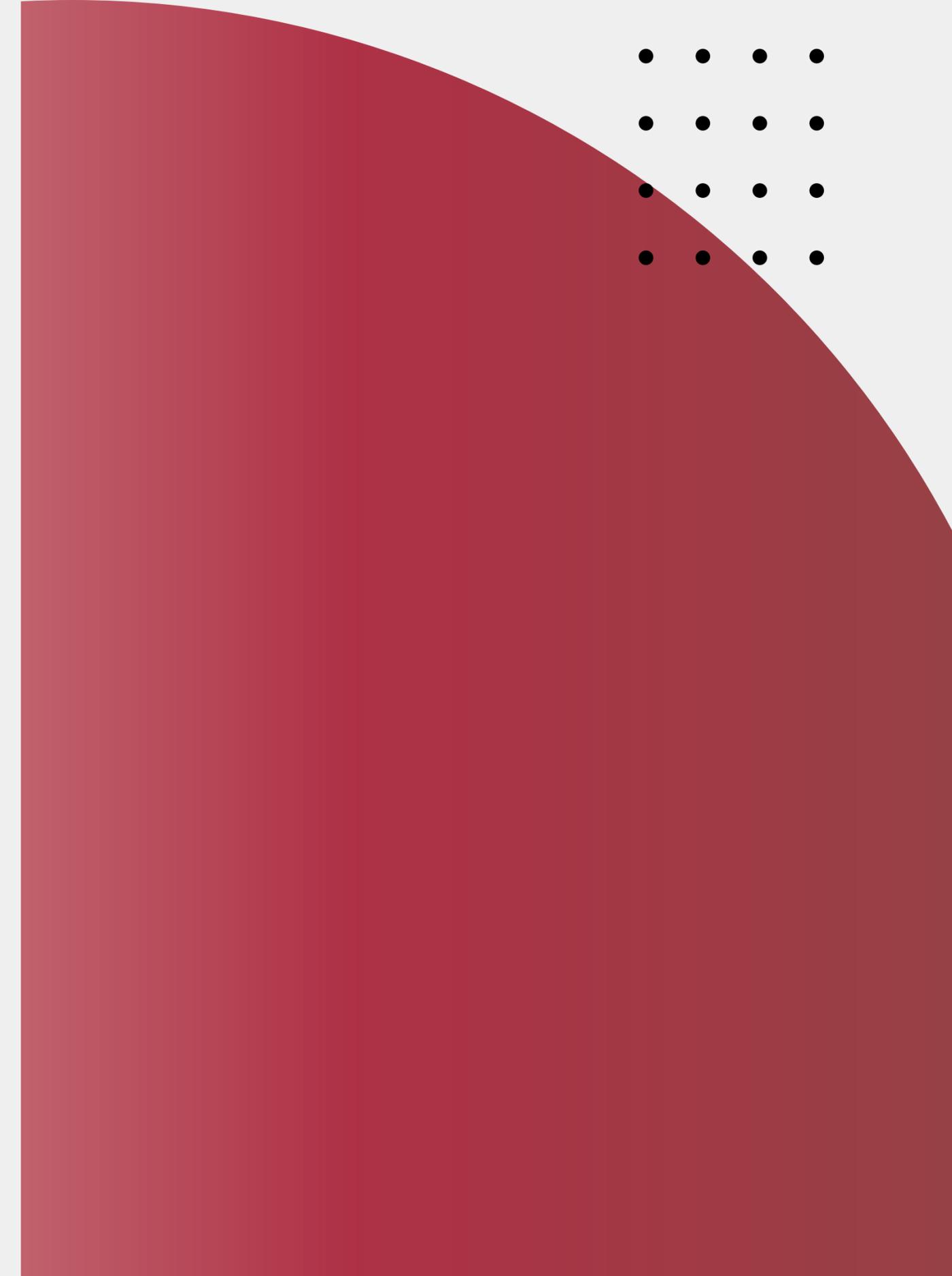
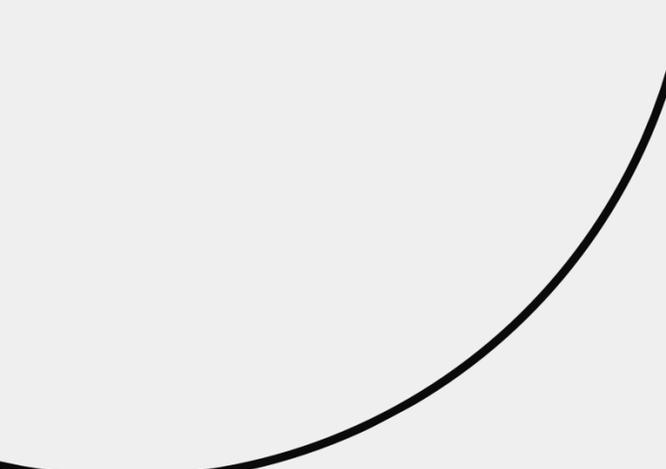


Dr. Kim Babiarz
Technical Mentor



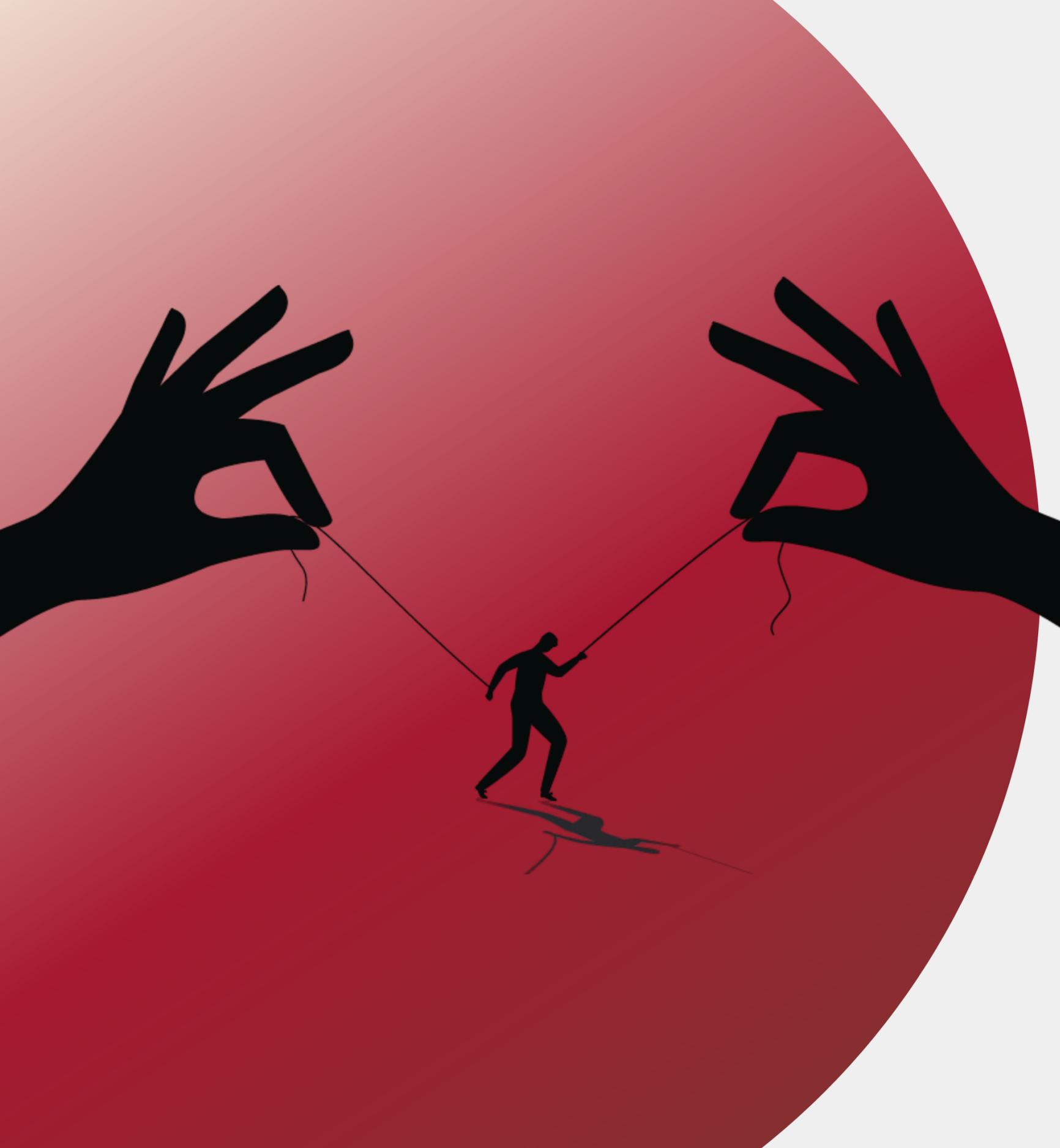
Jonas Junnior
Technical Mentor





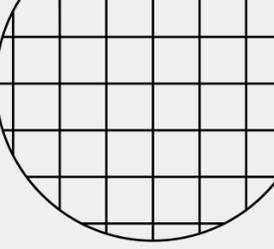
Background

Modern Slavery. Brazil. Charcoal. Our Goal.



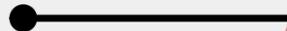
Modern Slavery

- **50 million** exploited annually.
- Victims **trapped** by threats and coercion.
- **Data gaps** hinder effective policy.
- Global efforts **lack data-driven impact**.



HTDLD's Brazil Focus

Over 1 million
trapped in
modern slavery.

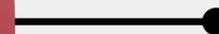


Robust data from
record-keeping and
transparency laws.

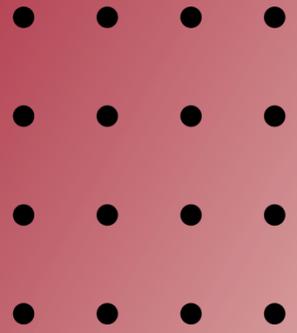


**Why
Brazil?**

Strong collaboration
with Brazil's Federal
Labor Prosecution Office.



HTDL's Charcoal Focus



Why Charcoal?



- Labor-Intensive Production Process
 - Exploitation Risk
 - Environmental and Economic Factors
 - Detection Challenges
 - Satellite Tracking
-



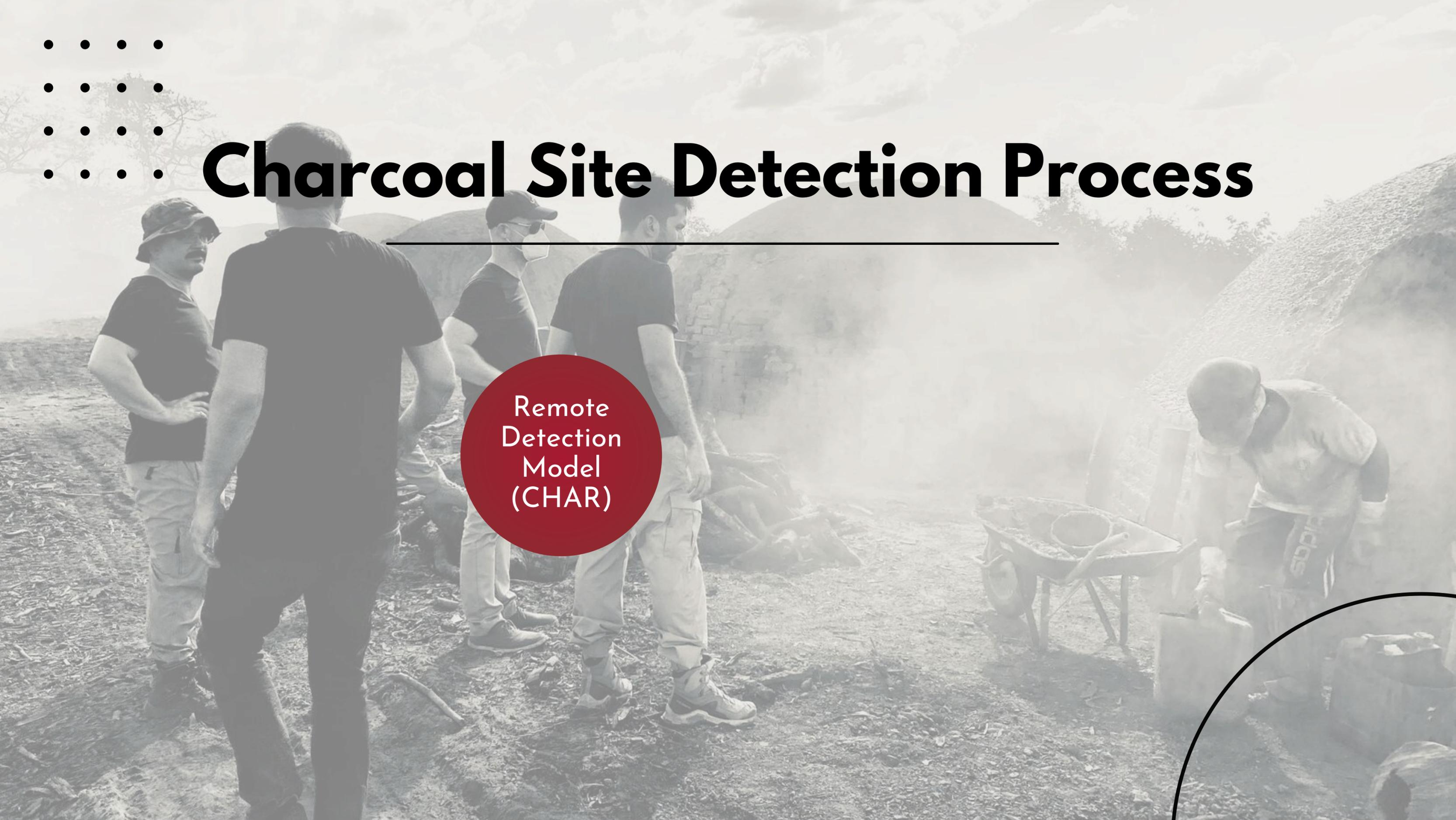
Charcoal Site Detection Process





Charcoal Site Detection Process

Remote
Detection
Model
(CHAR)



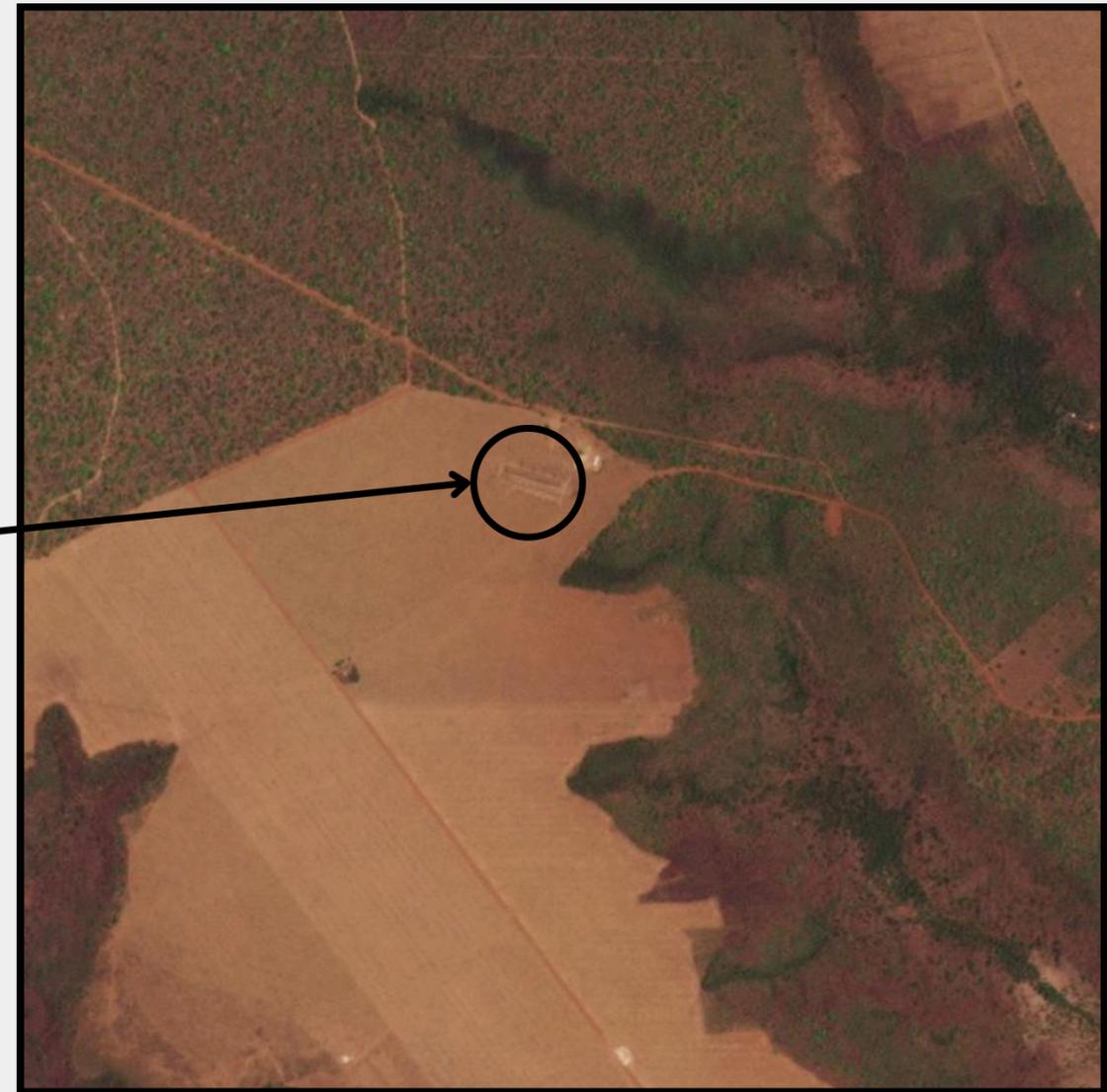
What We Seek



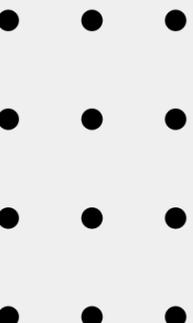
What We Seek



High Resolution Image



Training Image





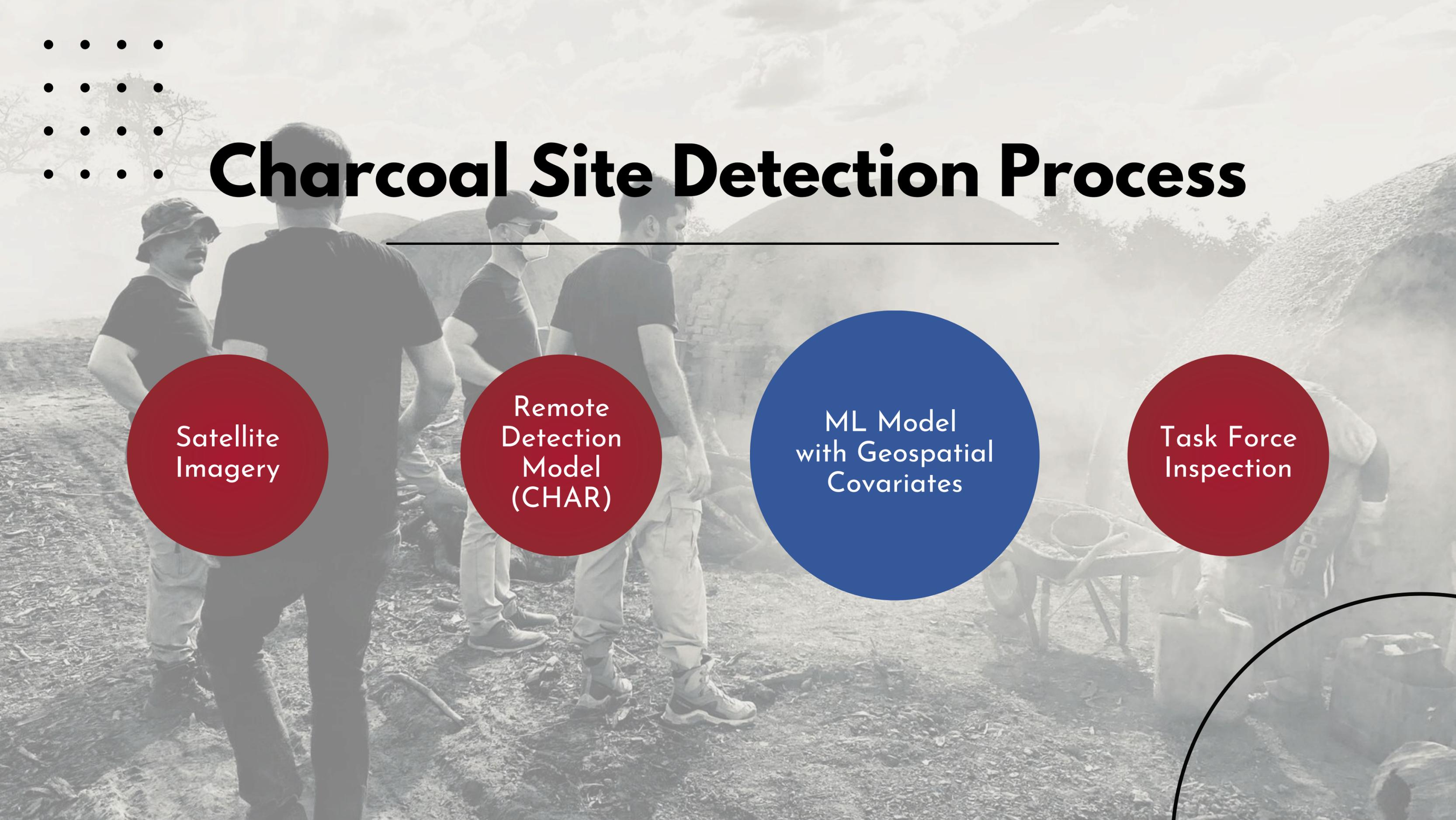
Charcoal Site Detection Process

Satellite
Imagery

Remote
Detection
Model
(CHAR)

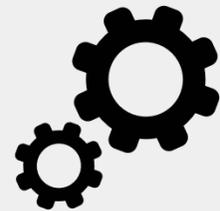
ML Model
with Geospatial
Covariates

Task Force
Inspection



Our Goal

Elevate the human post-processing.



Develop ML Models



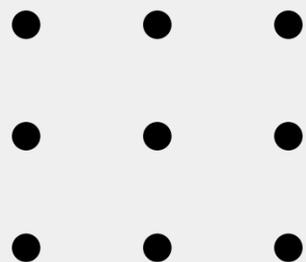
Explore Geospatial Data

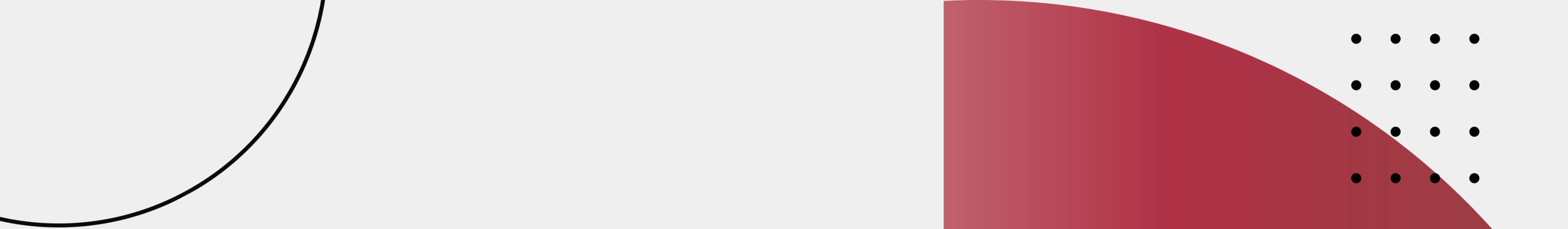


Analyze Feature Importance



Expand Training Data





Feature Engineering

GeoPandas Pipeline. Distance. Density.

Charcoal Site Data

5278 Sites

Flagged by the CHAR model from satellite images of Maranhão. Threshold of 0.9.

478 True Sites

Manually labelled and confirmed as charcoal sites.

Model Score

Model score from CHAR is included.

Month

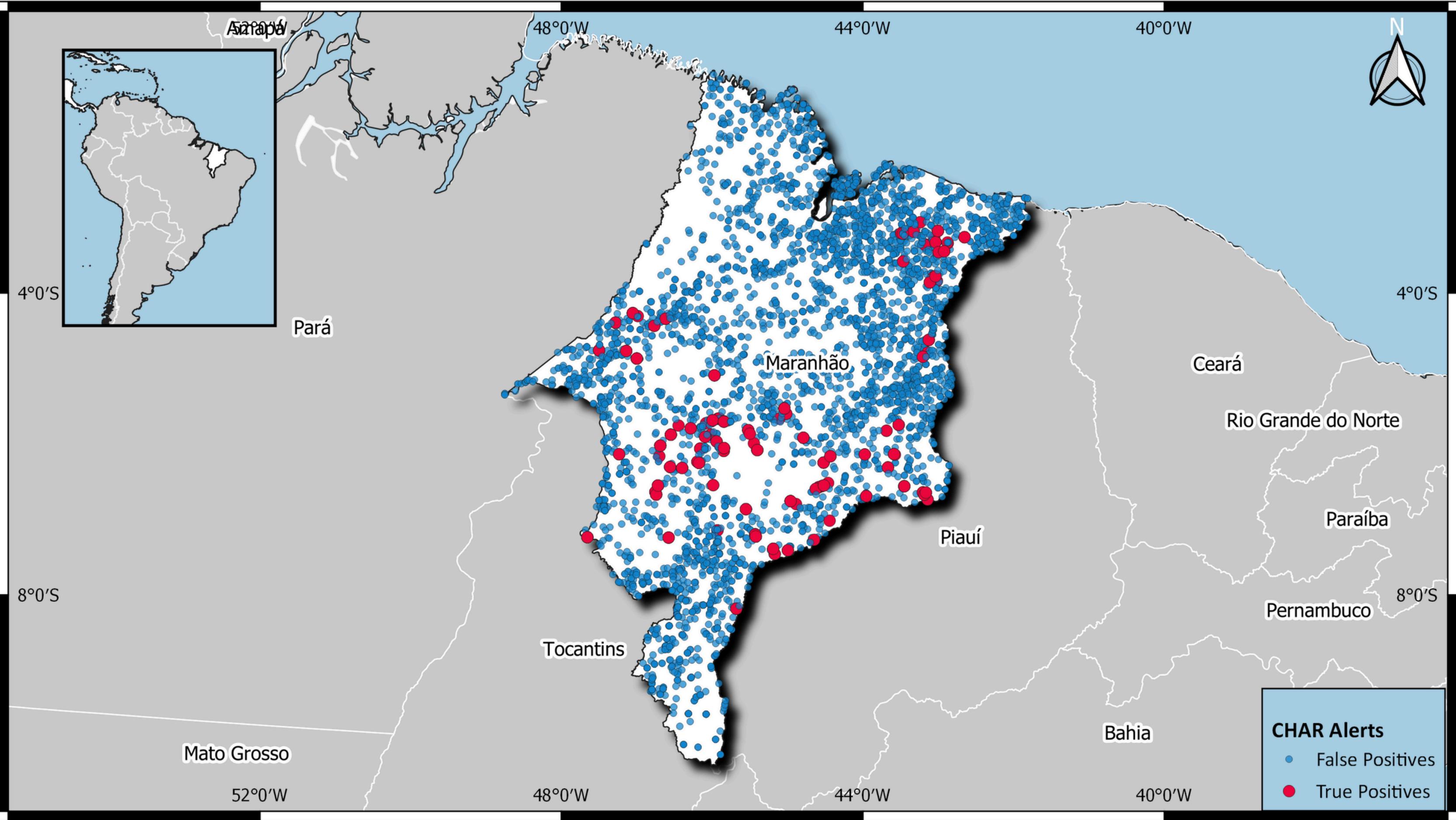
Images from 7/23 to 3/24.

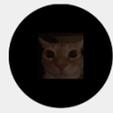
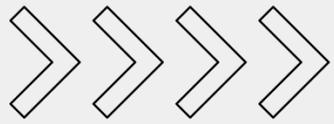
Geometry

Includes precise location of the flagged site.

Tiling

Satellite imagery grouped by unique tile ID.





Roads

Charcoal needs to be transported to steel mills.

Deforestation

Charcoal is made from cutting down trees.

**Model
Intuition**

Villages

Charcoal sites may want to be far away from villages to avoid detection.

Other Sites

Expect some clustering of charcoal sites.

Feature Construction

1

Data Source

EDA and thinking to determine relevant data which may have signal.

2

Appropriate Metric

Determine which metric to construct. Shortest distance to, feature count within a radius, within municipality, etc.

3

GeoPandas

Create pipeline to query database and construct features to be fed into the model.

Feature Construction

1

Data Source

EDA and thinking to determine relevant data which may have signal.

➤ **SmartLab:**

- Contains survey data of every municipality in Brazil. Includes data like literacy rate, poverty rate, number of workers rescued, and so on.

➤ **Geographic Features:**

- Geometries (locations) of roads, lakes, towns, indigenous lands, deforestation permits.

➤ **MapBiomas Alerts:**

- Geometries of deforestation alerts that are updated every two weeks by the Brazilian government.

Feature Construction

2

Appropriate Metric

Determine which metric to construct. Shortest distance to, feature count within a radius, within municipality, etc.

- It makes sense to ask how many charcoal sites might be within 10 km of a charcoal site.
- It makes sense to ask how many lakes are within 10 km of a charcoal site and how close a charcoal site is since the number and distance of lakes may have a bearing on whether to setup a charcoal site or not.

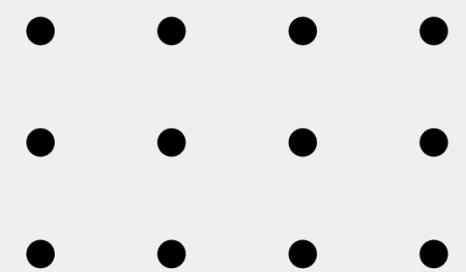
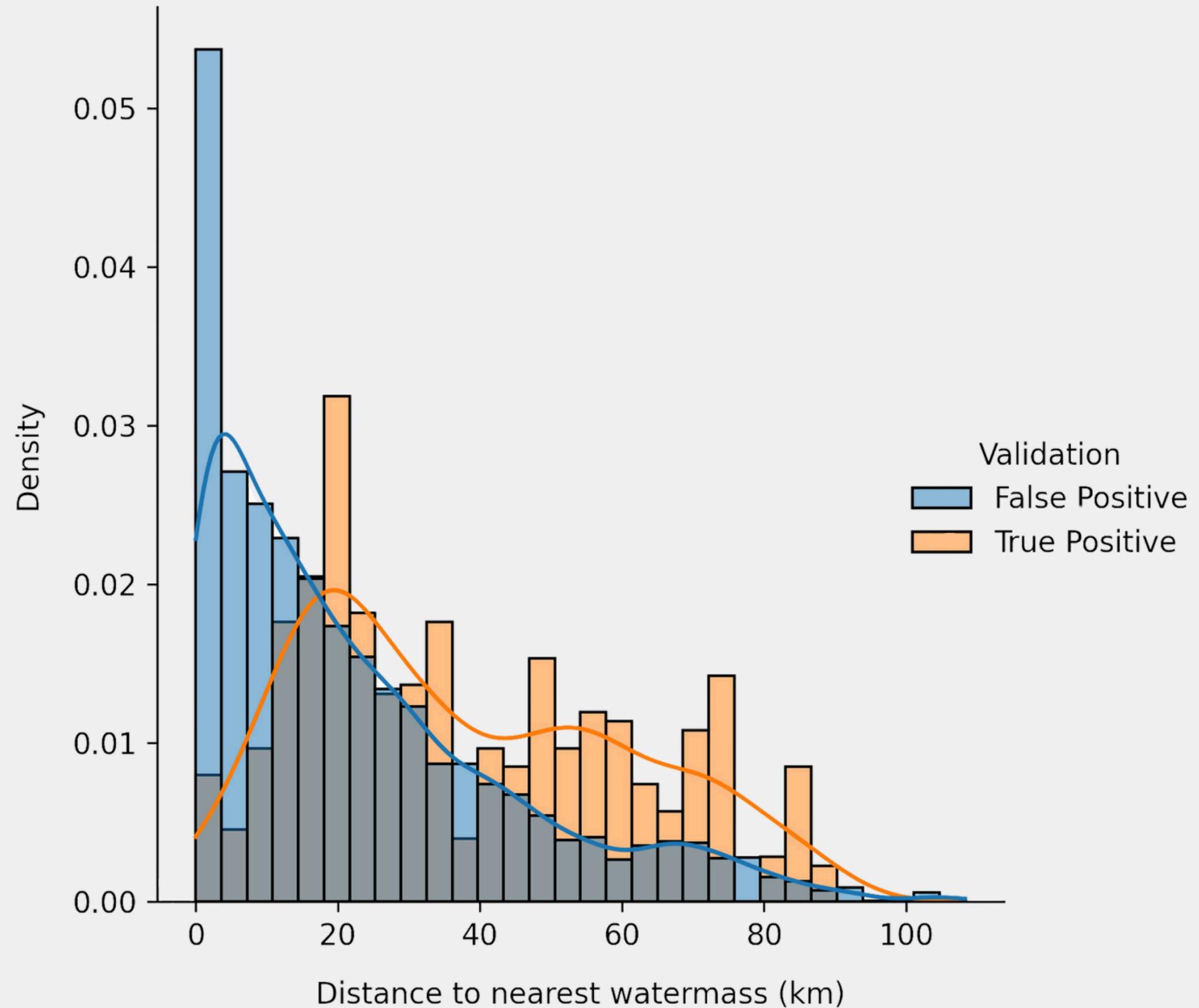
3

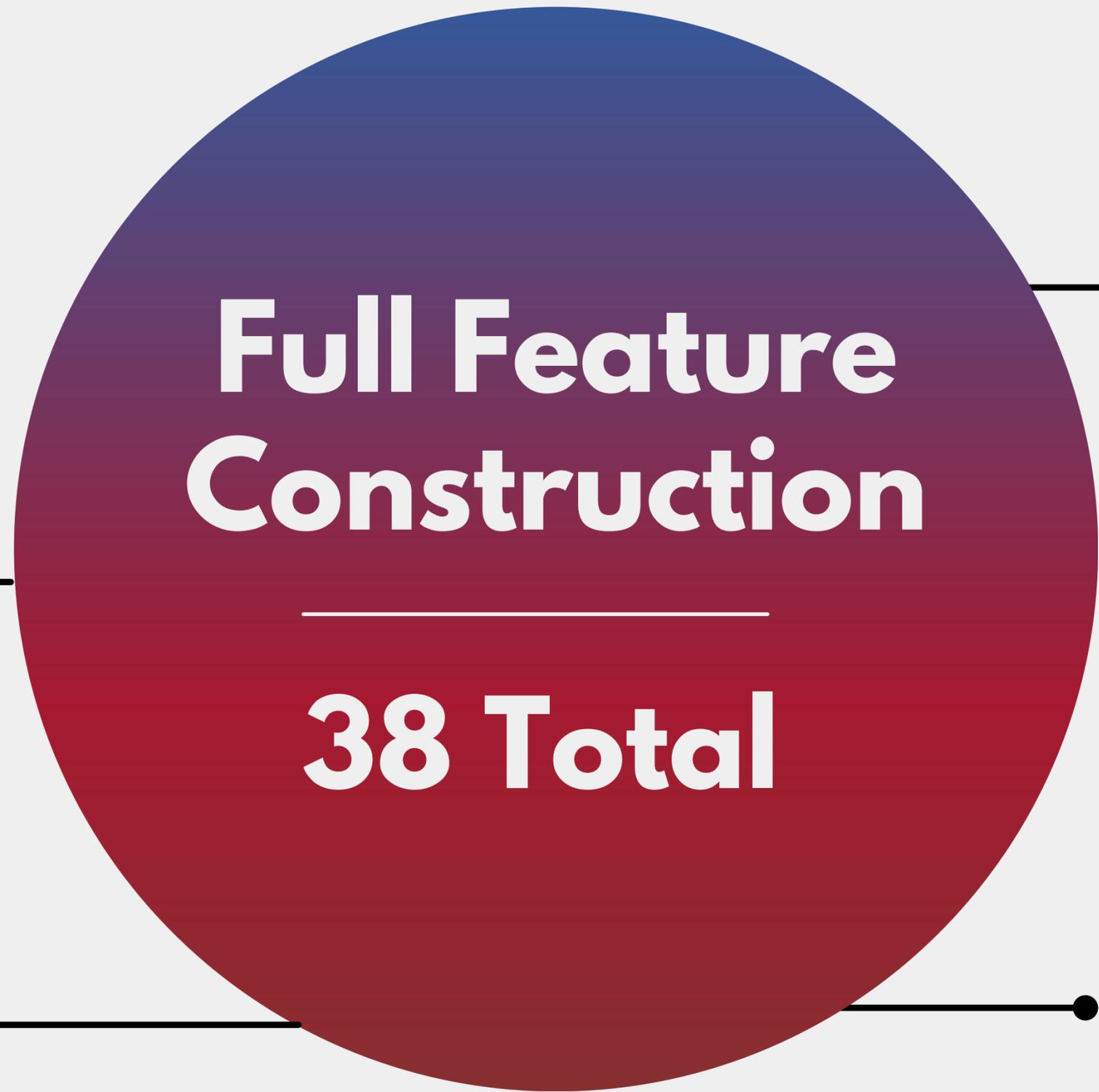
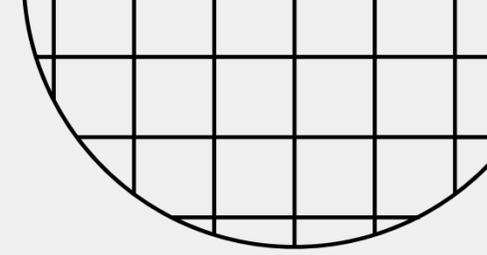
GeoPandas

Create pipeline to query database and construct features to be fed into the model.

Initial Results

Discernment between FP and TP sites exists. Suggests there is signal here for the model to pick up on.





14 Distance Variables

Shortest straightline distance to feature

12 Landcover Categories

Forest plantation, savannah formation, etc.

3 Density Variables

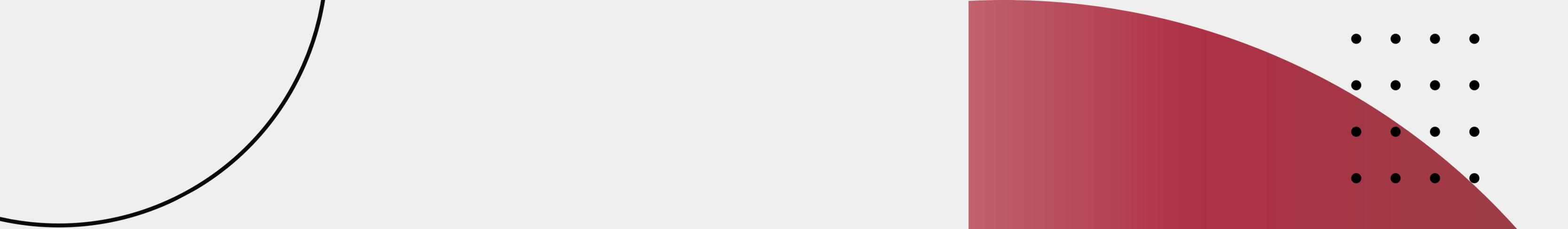
Feature count within 10 km radius

9 Survey Variables

From SmartLab data on poverty, literacy rate, rescued workers, etc.

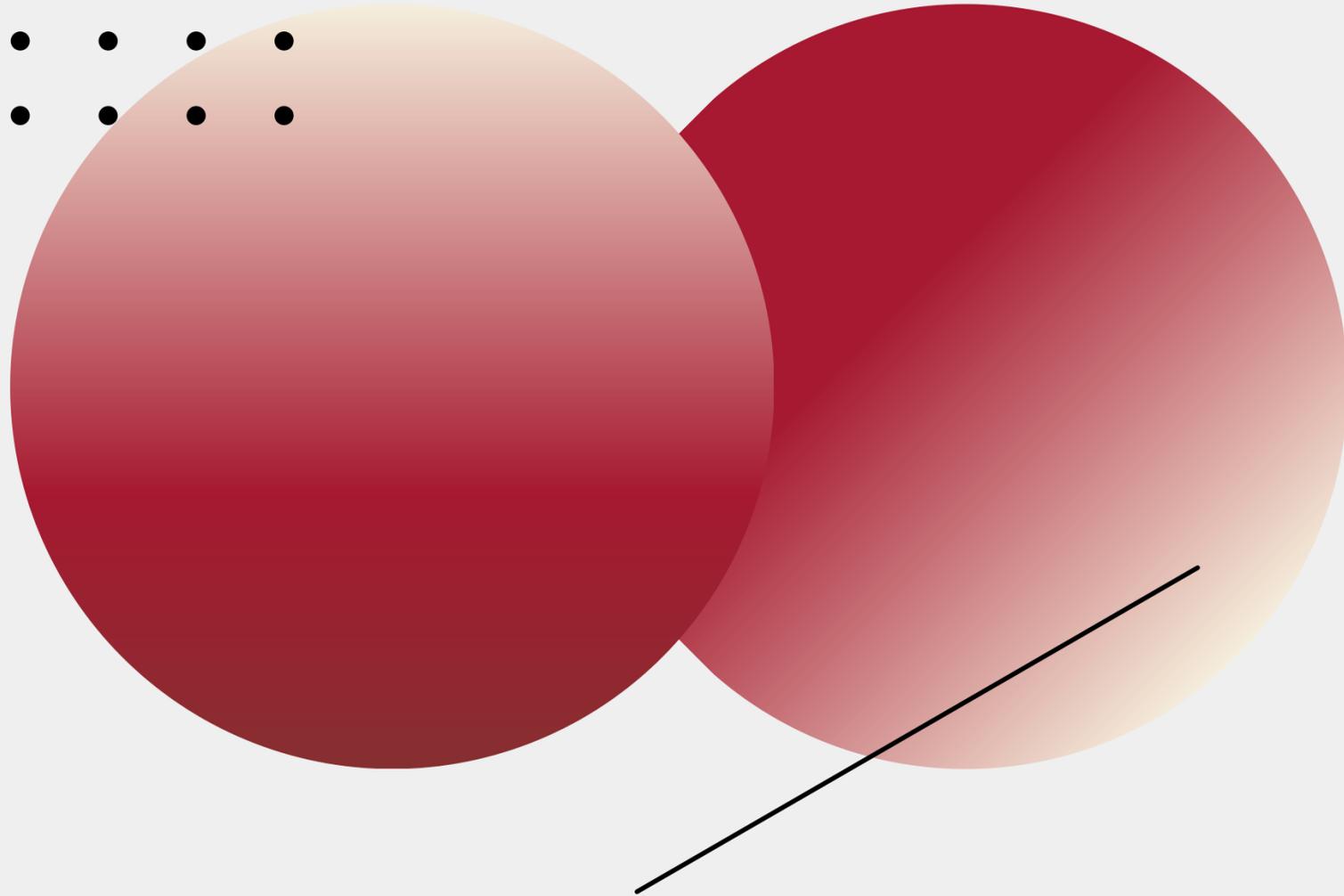
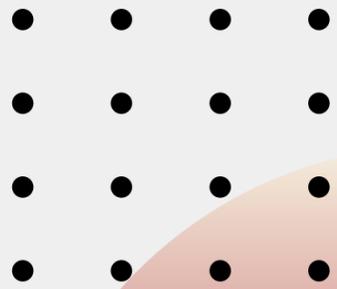
Full Feature Construction

38 Total



Machine Learning Model

Architectures. Analytics. Performance.



Implementation

- Grouping and Stratification
- Model Architectures
- Hyperparameter Tuning and Result Analysis

Data Handling

1

Grouping

Group datapoints by location to prevent train/test knowledge leakage.

2

Stratification

Balance by label to ensure sufficient training points and consistent evaluation.

3

Splitting

1/6 Holdout set, remaining 5/6 broken into 5-fold cross validation.

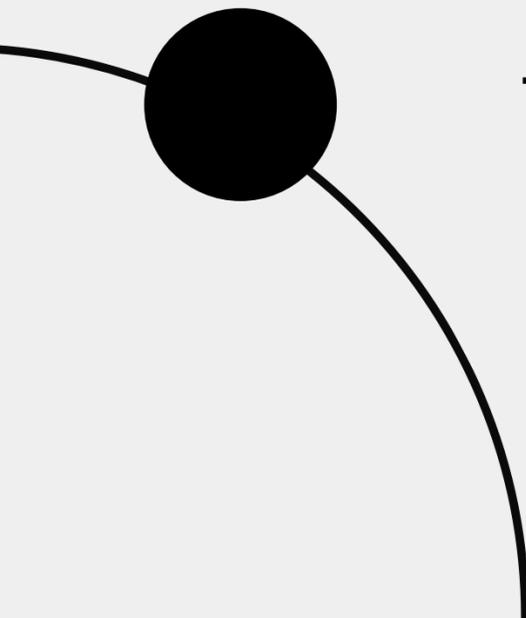
Model Architectures

Tree-based models

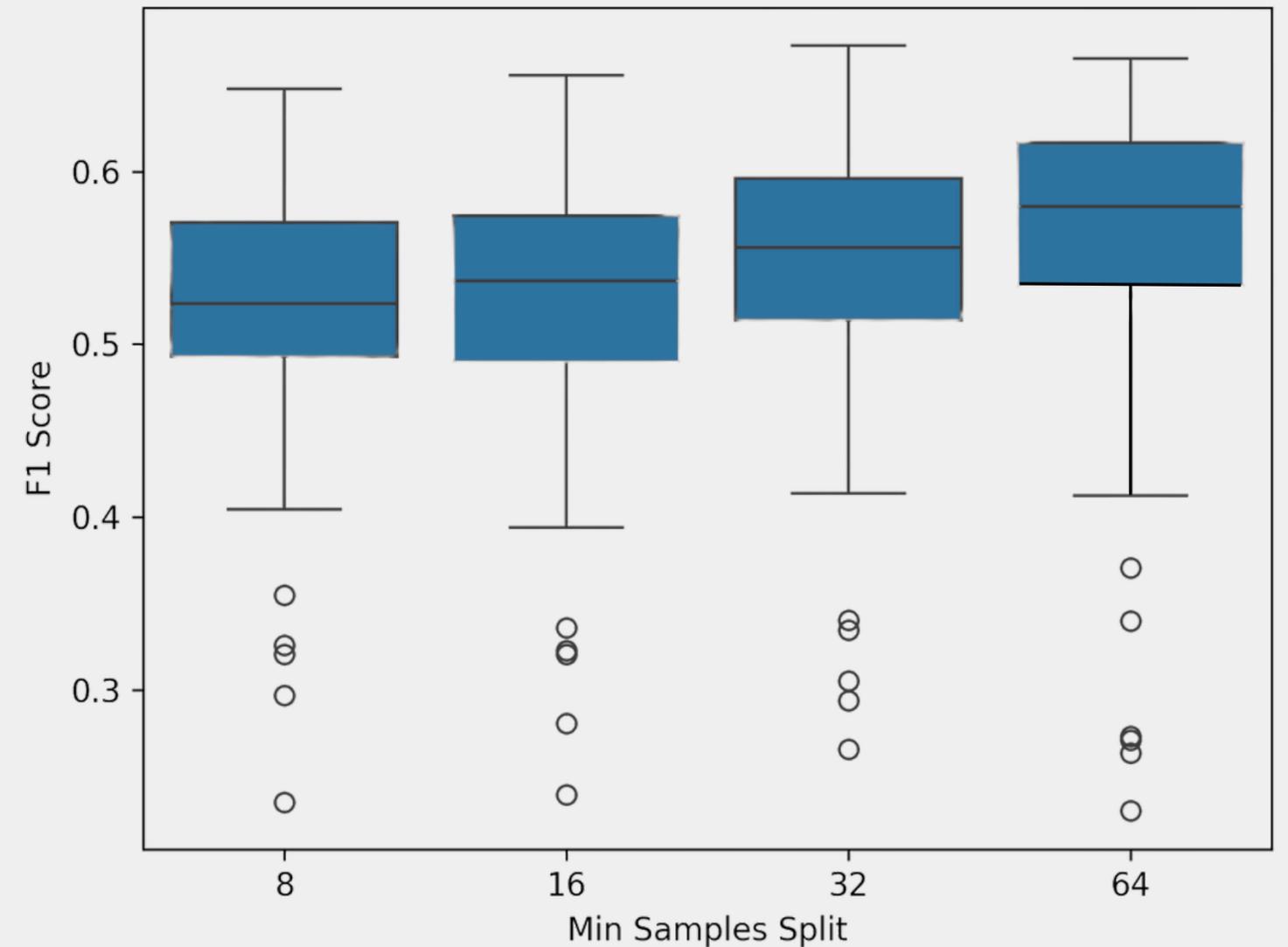
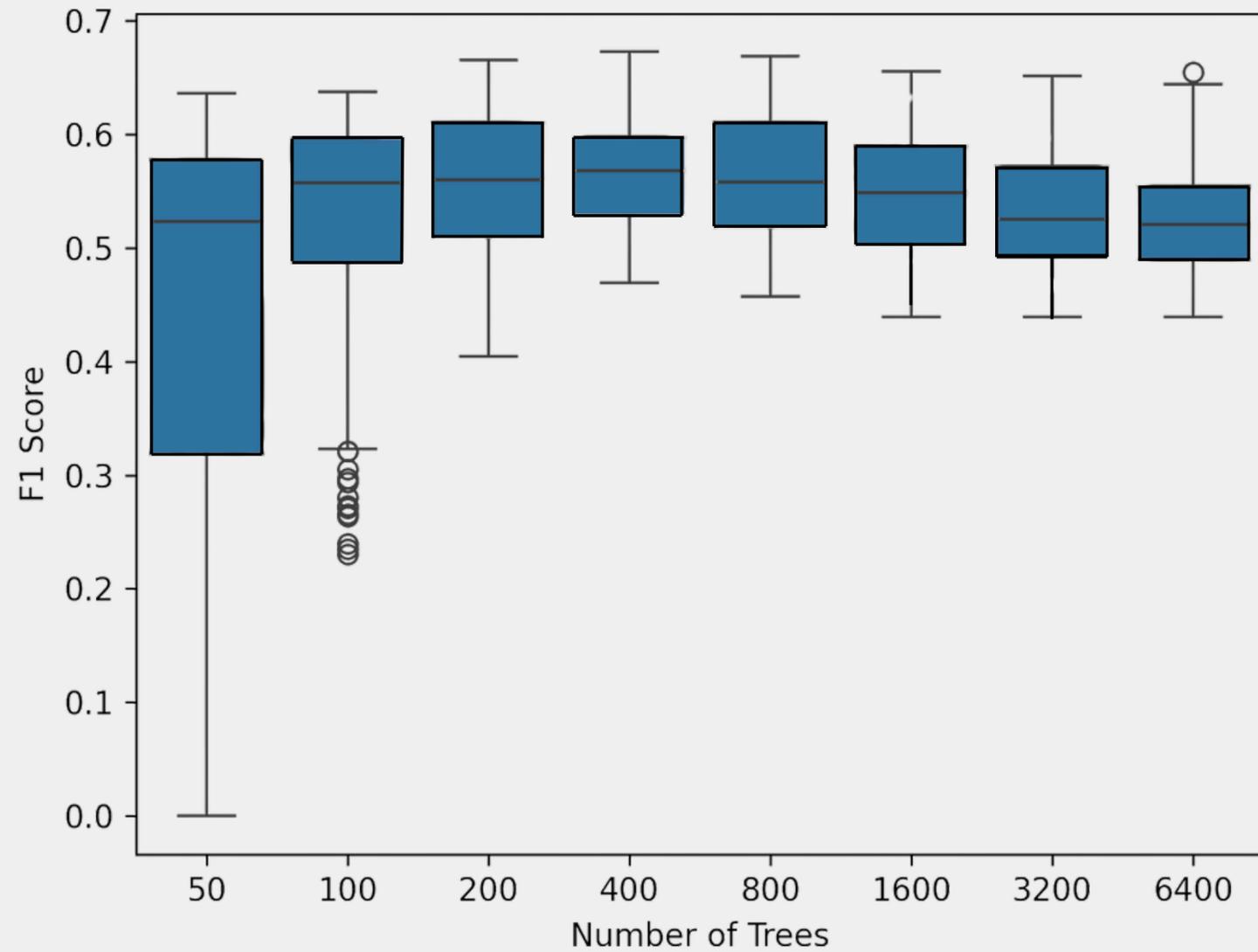
Gradient Boosting, Random Forest

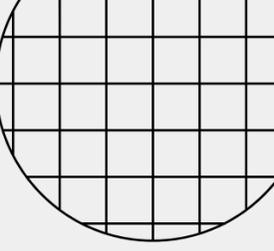
Transformer-based models

TabPFN



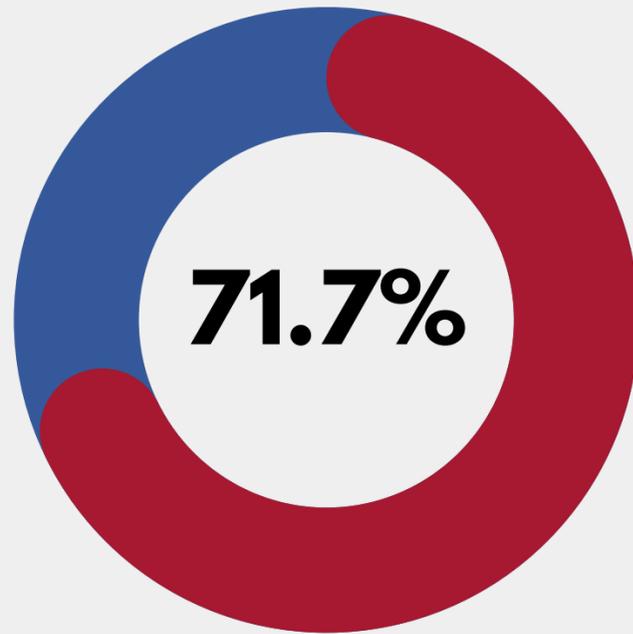
Hyperparameter Tuning



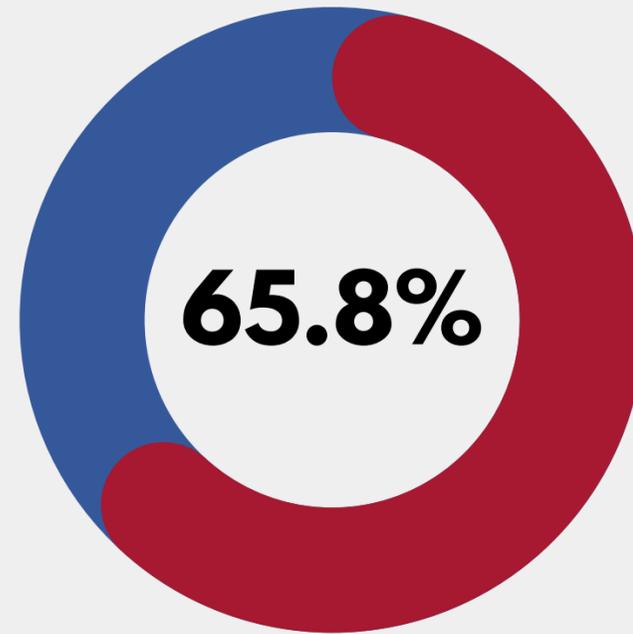


Model Performance

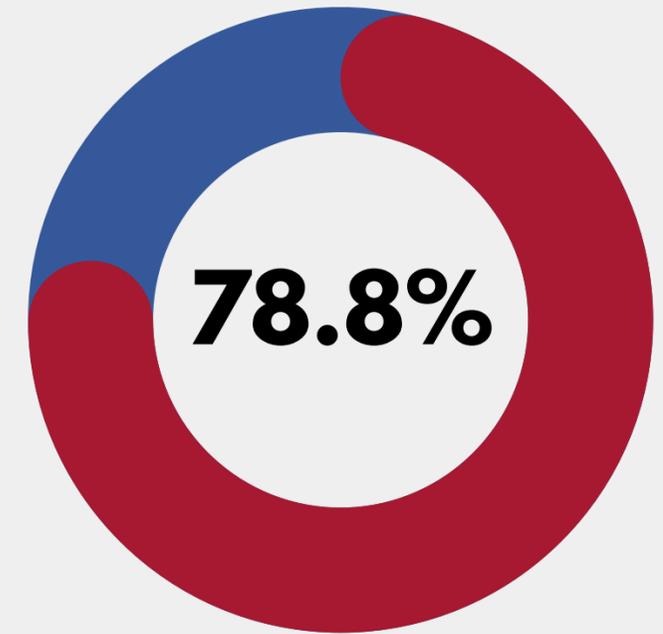
On validation set, at threshold 0.25.



F1 Score



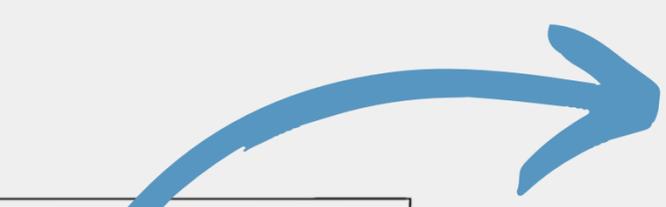
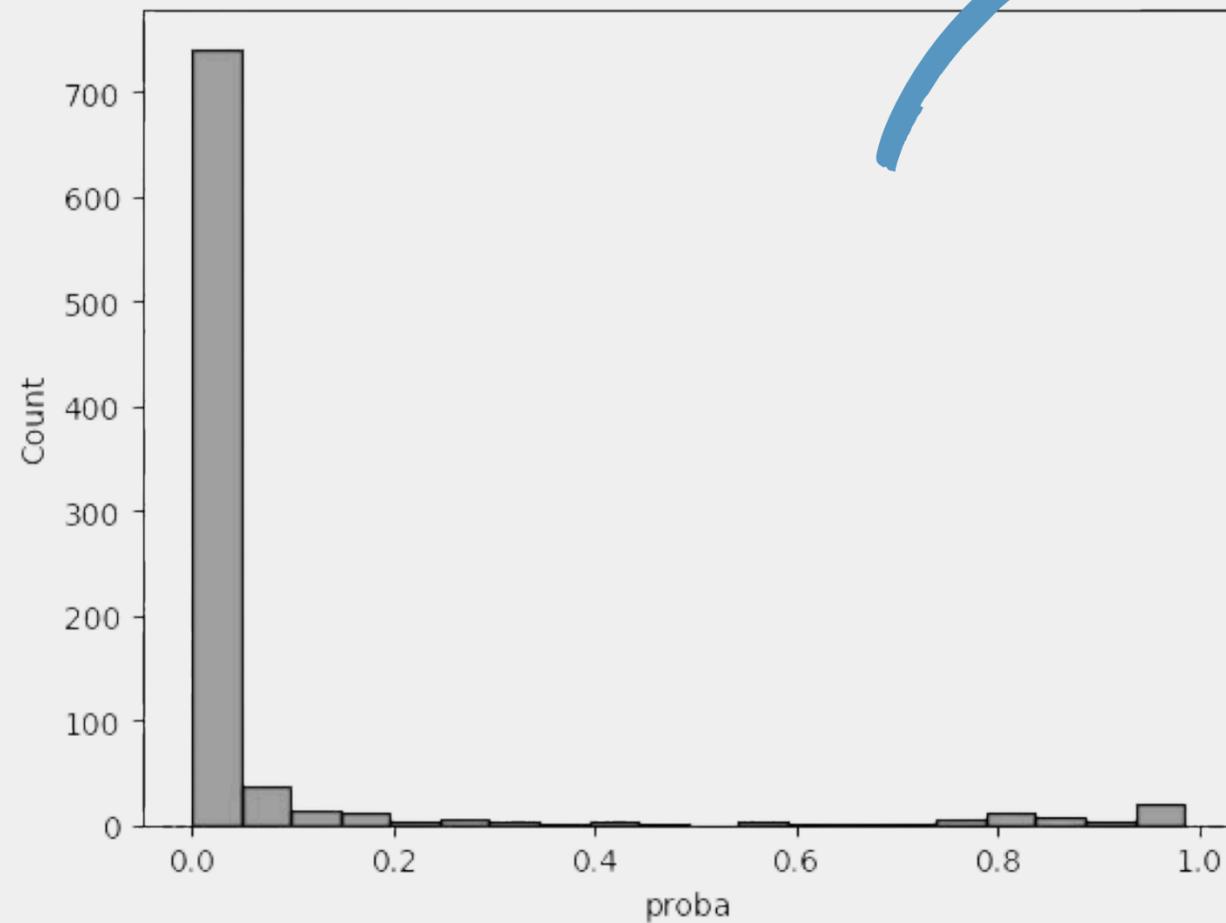
Precision



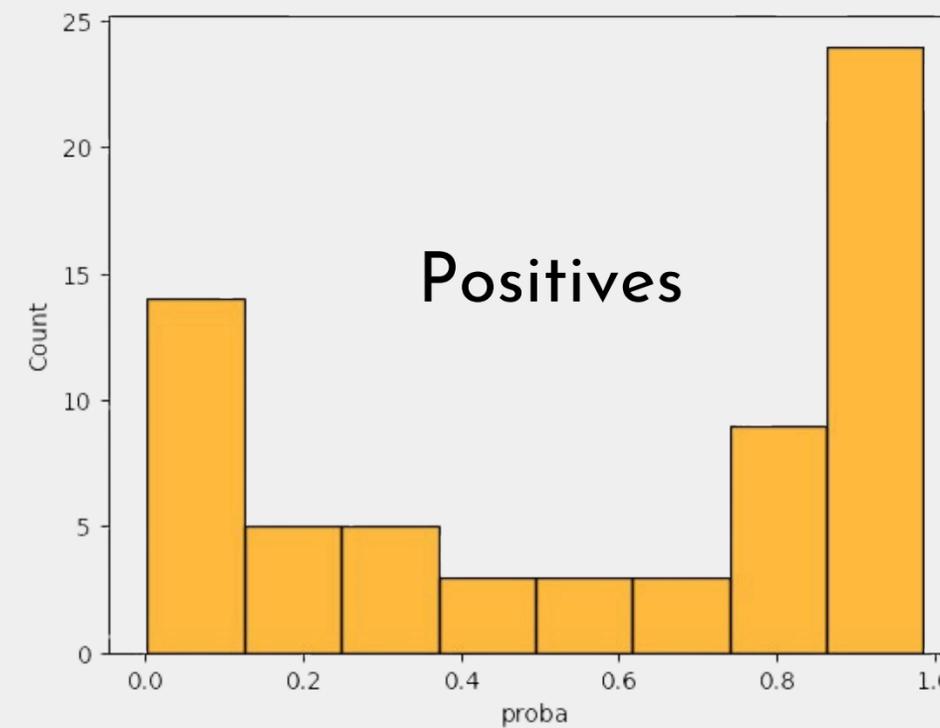
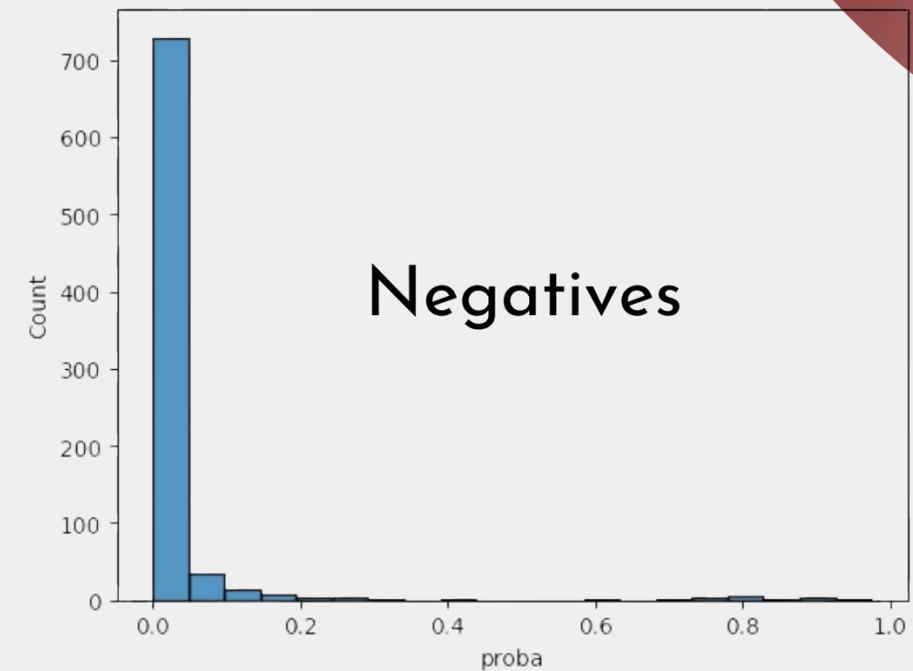
Recall



Model Analytics

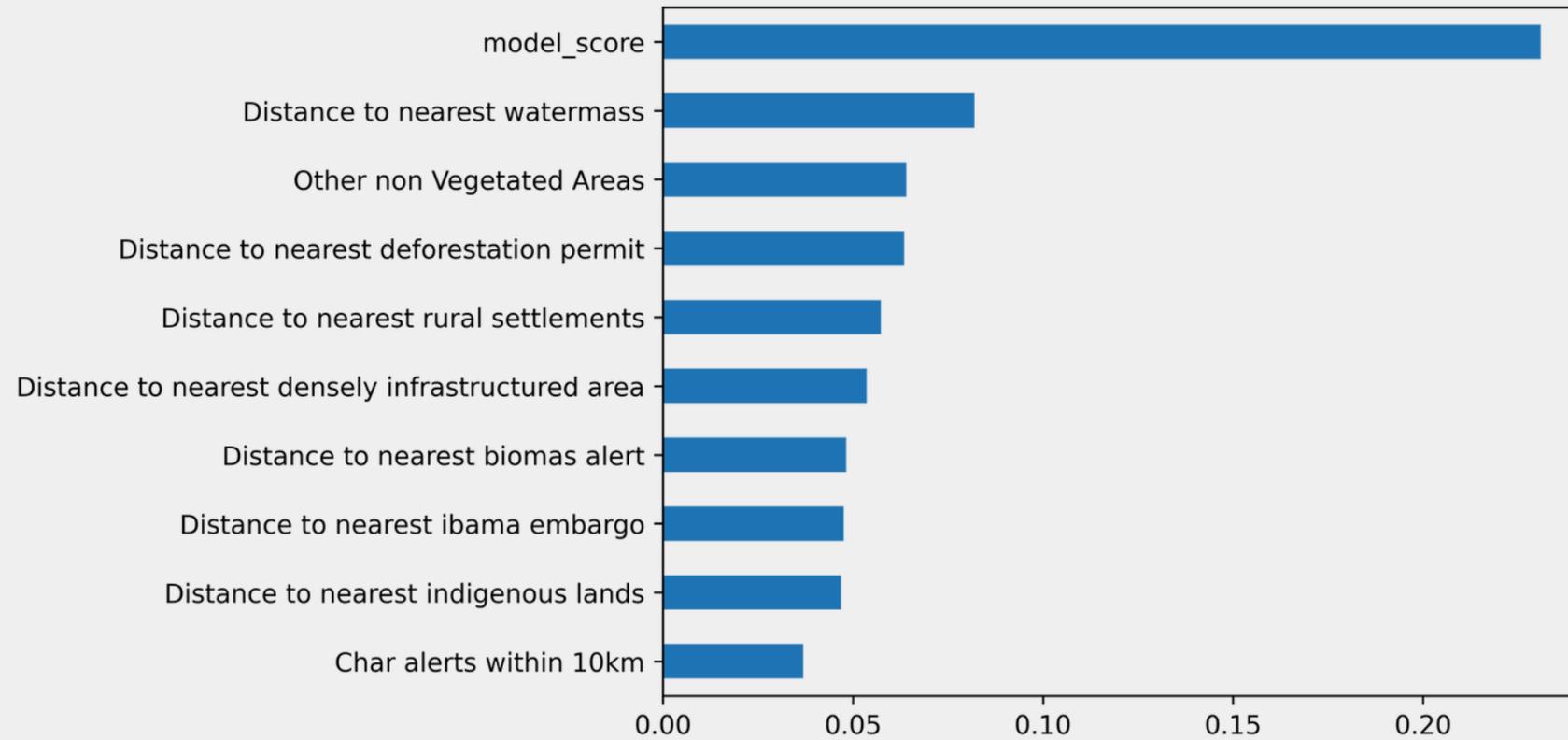


Ground Truth

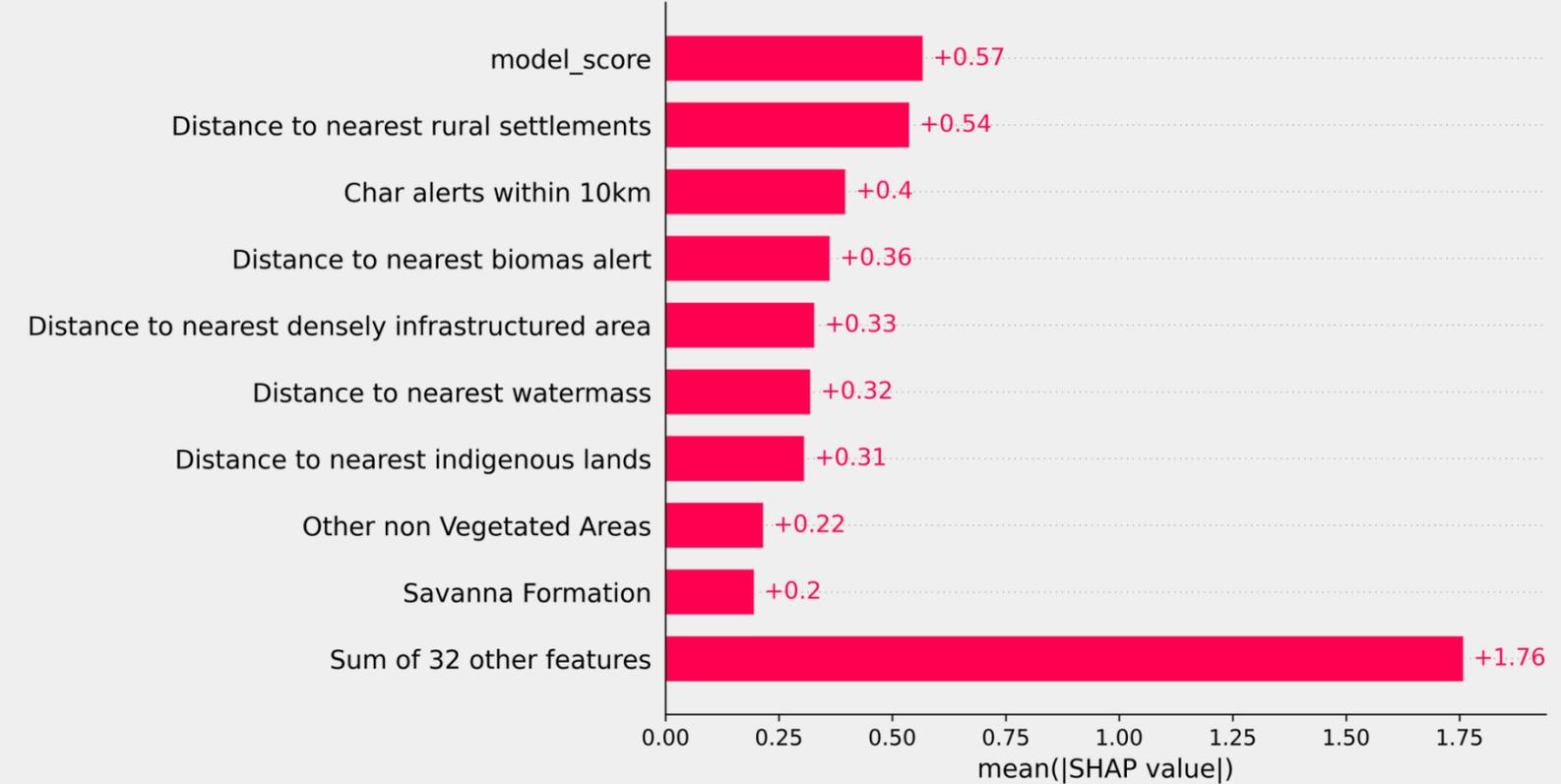


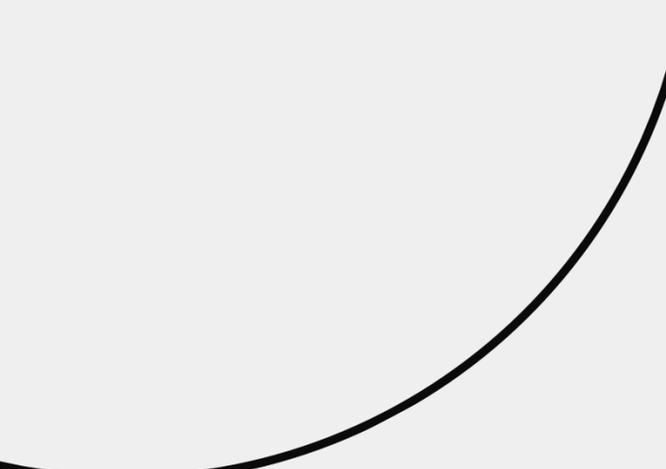
Feature Importance

Built-in Gradient Boost (top 10)



TreeSHAP



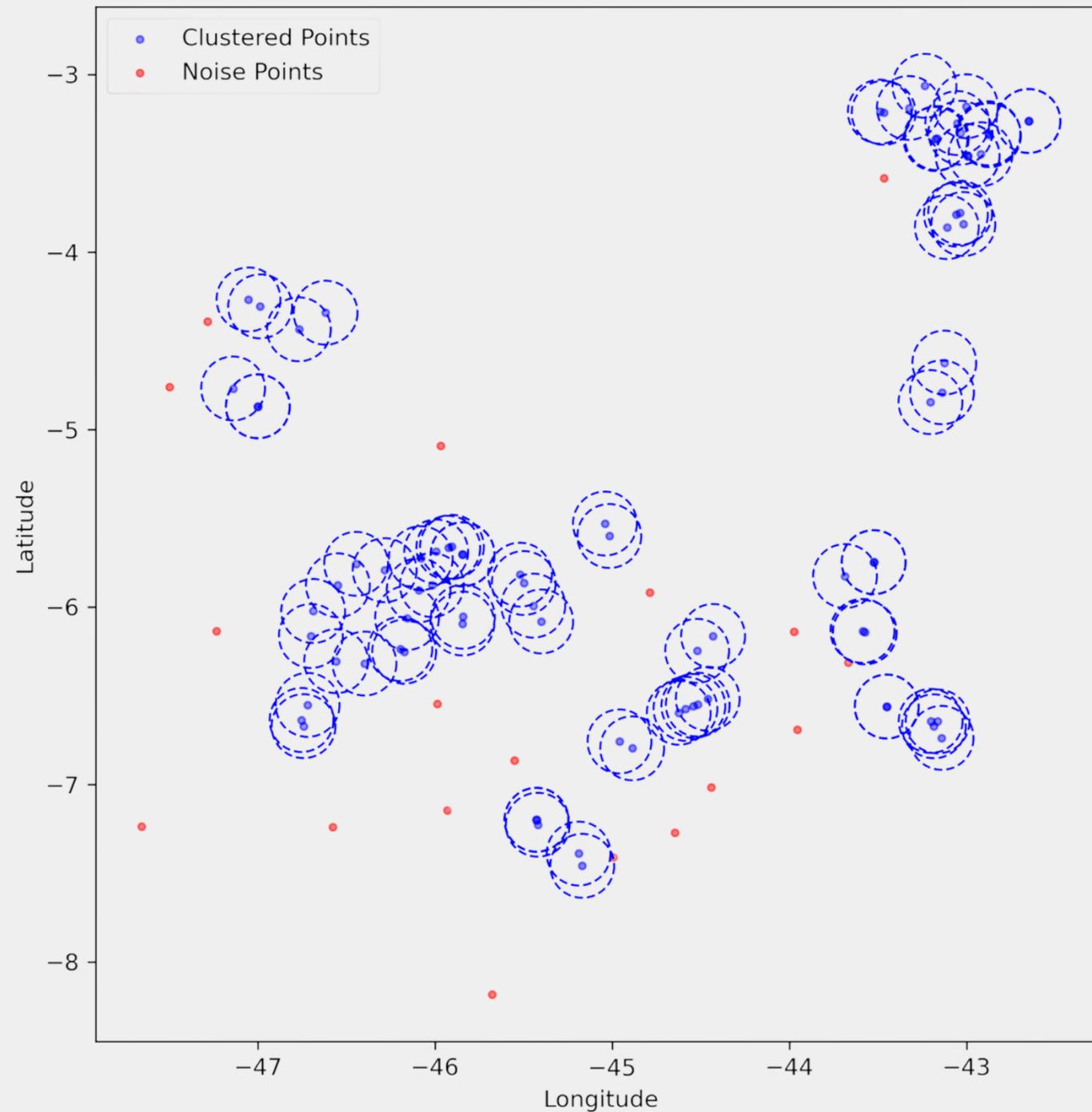


Cluster Analysis

Geospatial. Full Feature.

Geospatial Cluster Analysis

DBSCAN Clustering with Buffer Zones



➤ Total True Charcoal Sites: 86

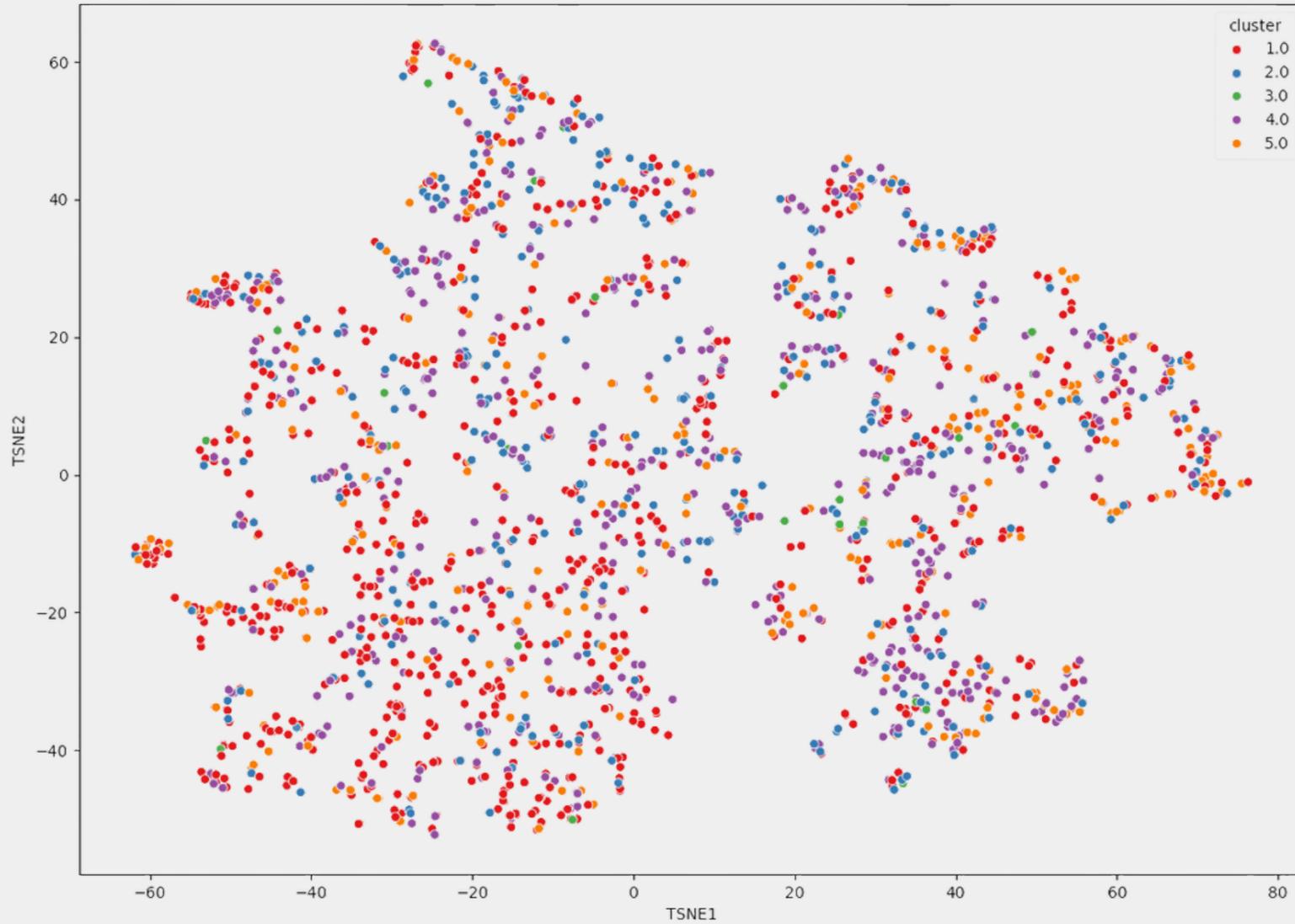
➤ Unique Clusters Identified: 26

➤ Max Clustering Distance: 20 km

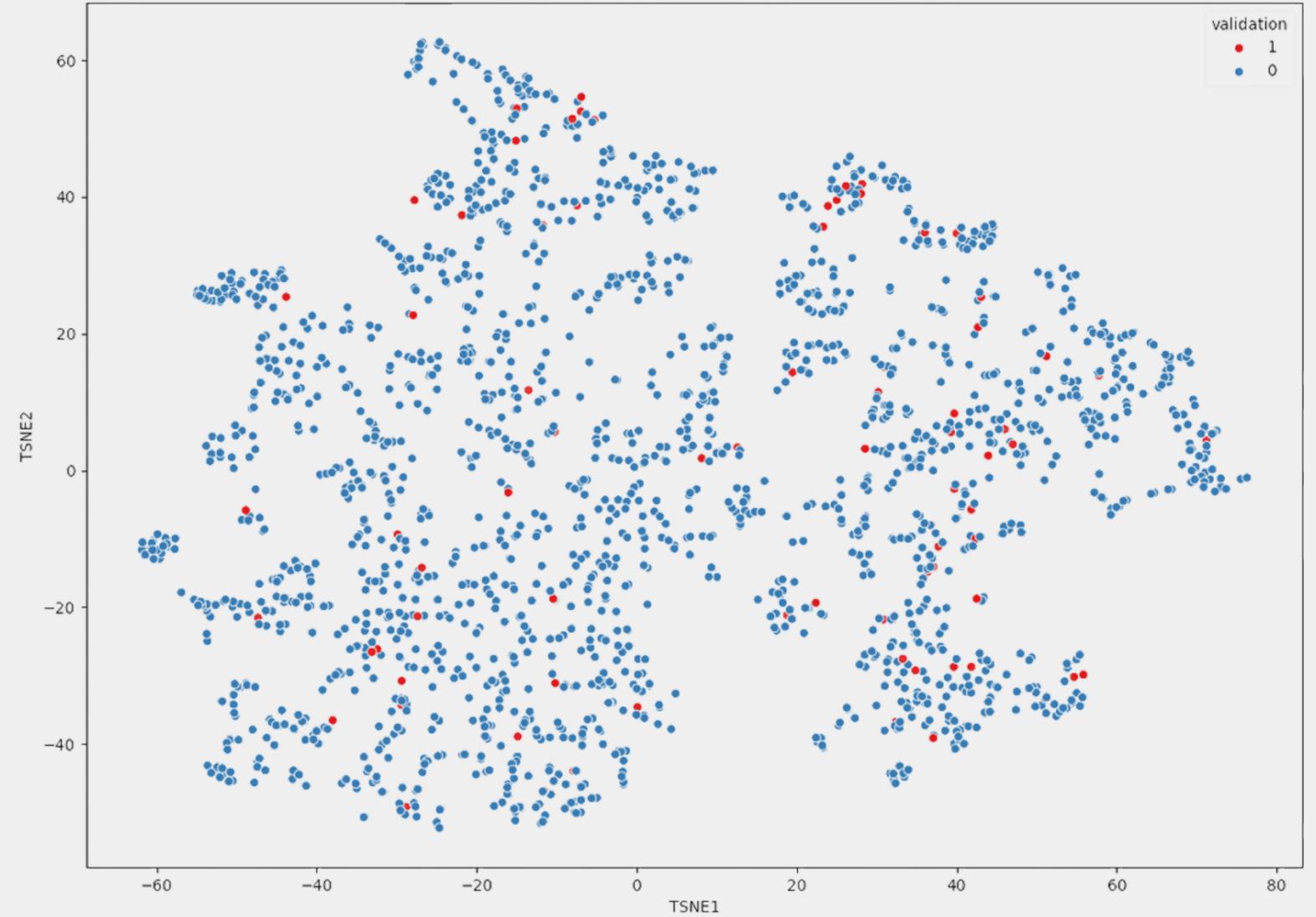
➤ Min Sites per Cluster: 2

Feature Cluster Analysis

t-SNE Visualization of Top Features by Cluster



t-SNE Visualization of Top Features by Validation





Conclusion

Future Work. Acknowledgements.

Future Work

Image feature embeddings

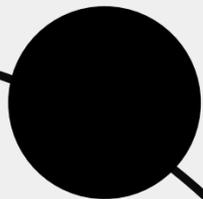
Enrichen the information from the first stage of the model.

Improved time-series modeling

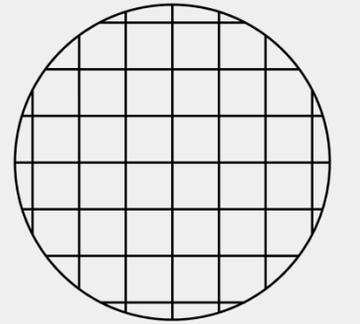
“Hotspot” feature, deforestation chronology

Feedback from fieldwork

Brazil FLPO task force deployment this August



Acknowledgements



Dr. Ben Seiler

Dr. Mike Baiocchi

Dr. Kim Babiarz, Jonas Junnior, and the
HTDL Lab

Shilaan Alzahawi, Dr. Balasubramanian
Narasimhan, Dr. Annie Lamar, and the
whole DSSG team

