

**FusionSite**

# Driver Safety Model

Trang Hoang and Isabel Arvelo  
04-14-25



# Agenda

- 1 Problem Description
- 2 Data Description
- 3 Model Development
- 4 Results
- 5 Highlights
- 6 Challenges



# What is FusionSite?

FusionSite is a high-growth premier waste management provider with services across 15 states.

# Problem Description

## PROBLEM

### Driver Safety

FusionSite needs proactive driver safety protocols across its fleet operations

Current approach is reactive (waiting for incidents)

Industry leader in portable sanitation, waste management & luxury trailer rentals  
Reputation for reliability at stake

## SOLUTION

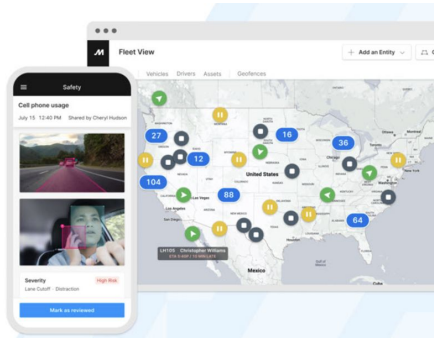
### Risk Prediction

Developing a Daily Driver Risk Model  
Leveraging GoMotive Fleet Management & Driver Safety Platform to analyze driving events data

Integrating driver historical driving patterns and pre-trip inspections

Using location-based data to incorporate state crash data and precipitation data in the radius of each site

# Data Sources



Motive API



Internal Company Data



Climate Data



Geospatial Data



Transportation Data

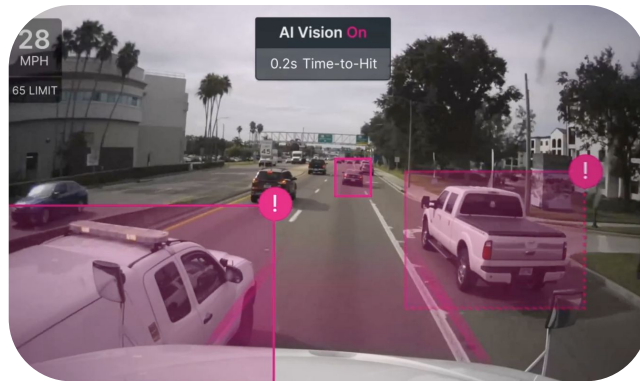
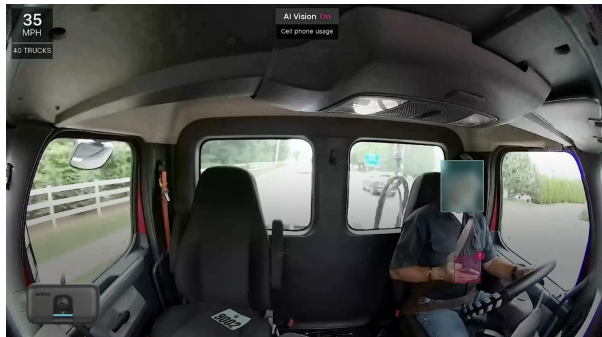


Insurance Claims Data

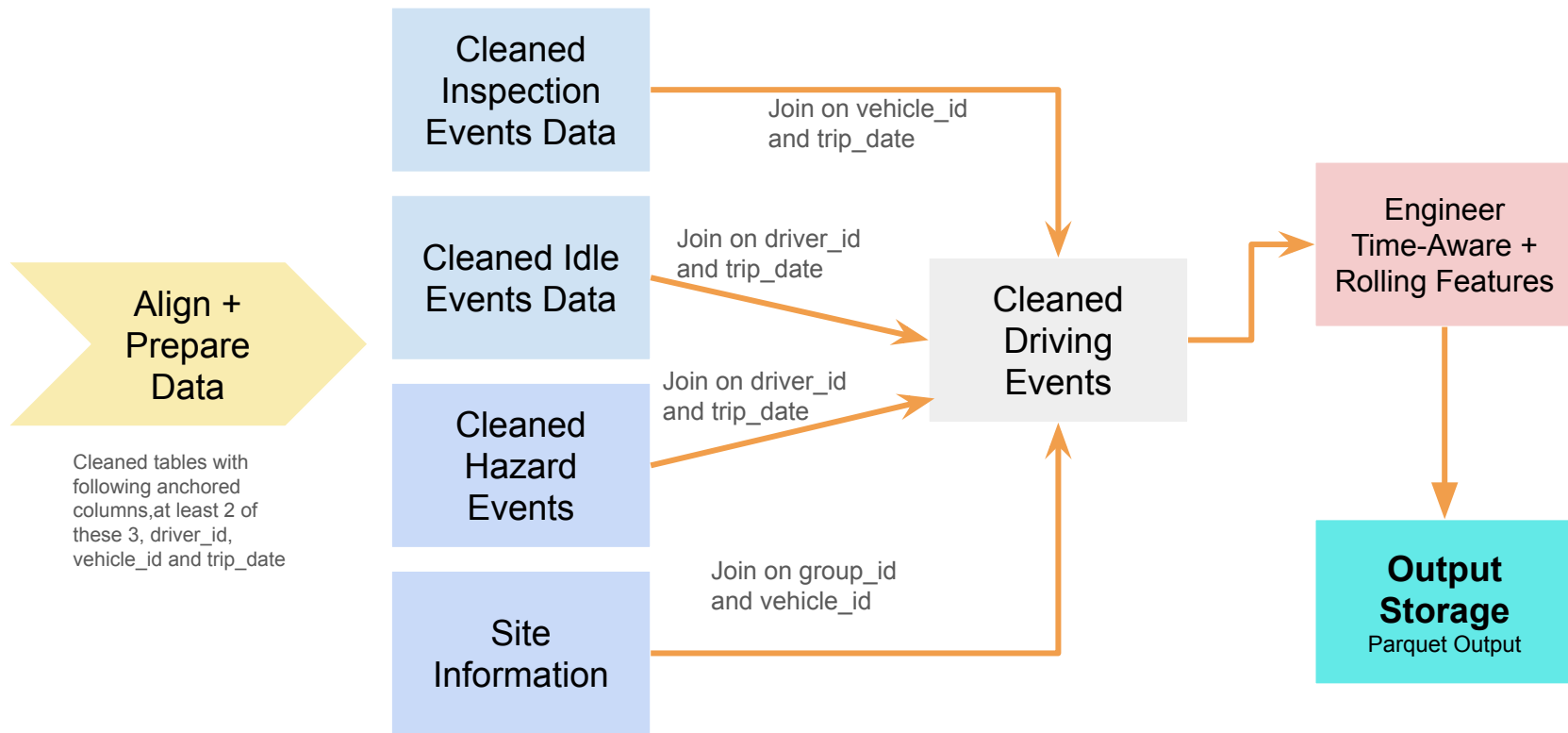
# Motive Data Preprocessing

Data about driver and vehicle activity from January 2023 to February 2025:

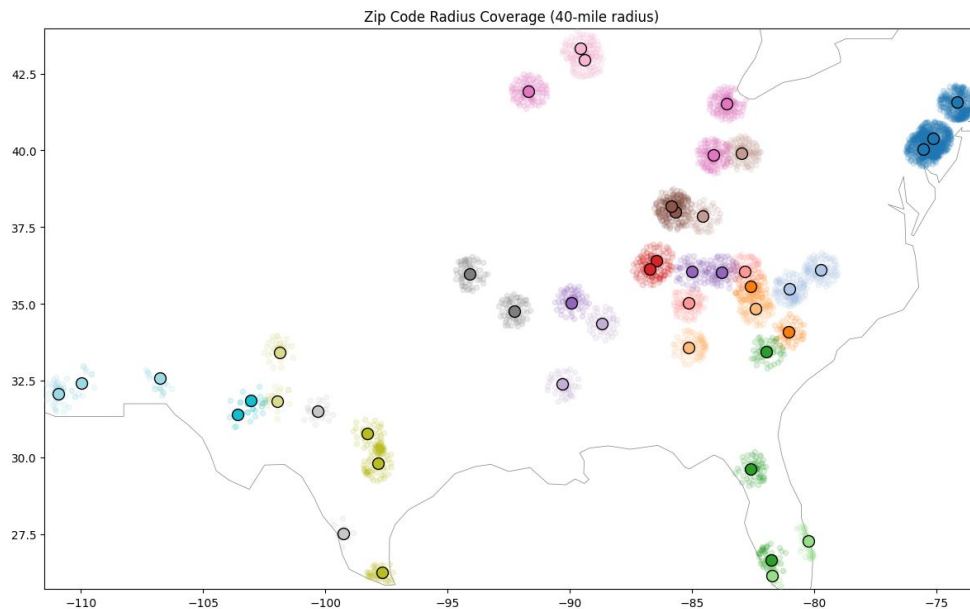
1. **Aggregated Driving Events**  
Captures each driver's behavior, history, and context on a per-trip basis, forming a comprehensive feature set.
2. **Aggregated Combined Events**  
Tracks the volume and types of safety-related events (e.g., seatbelt violations, distractions, speeding), including severity levels of speeding (low/mid/high). It helps identify risky behaviors and monitor unresolved coaching issues.
3. **Aggregated Inspection Events**  
Provides a trip-level summary of inspection activities per vehicle by counting total inspections and flagging open or resolved issues.
4. **Aggregated Idle Events**  
Summarizes idle behavior per driver per day to gauge frequency and duration of truck idling, useful for assessing driver efficiency, fuel consumption, and potential safety concerns.



# Joining FusionSite trips data



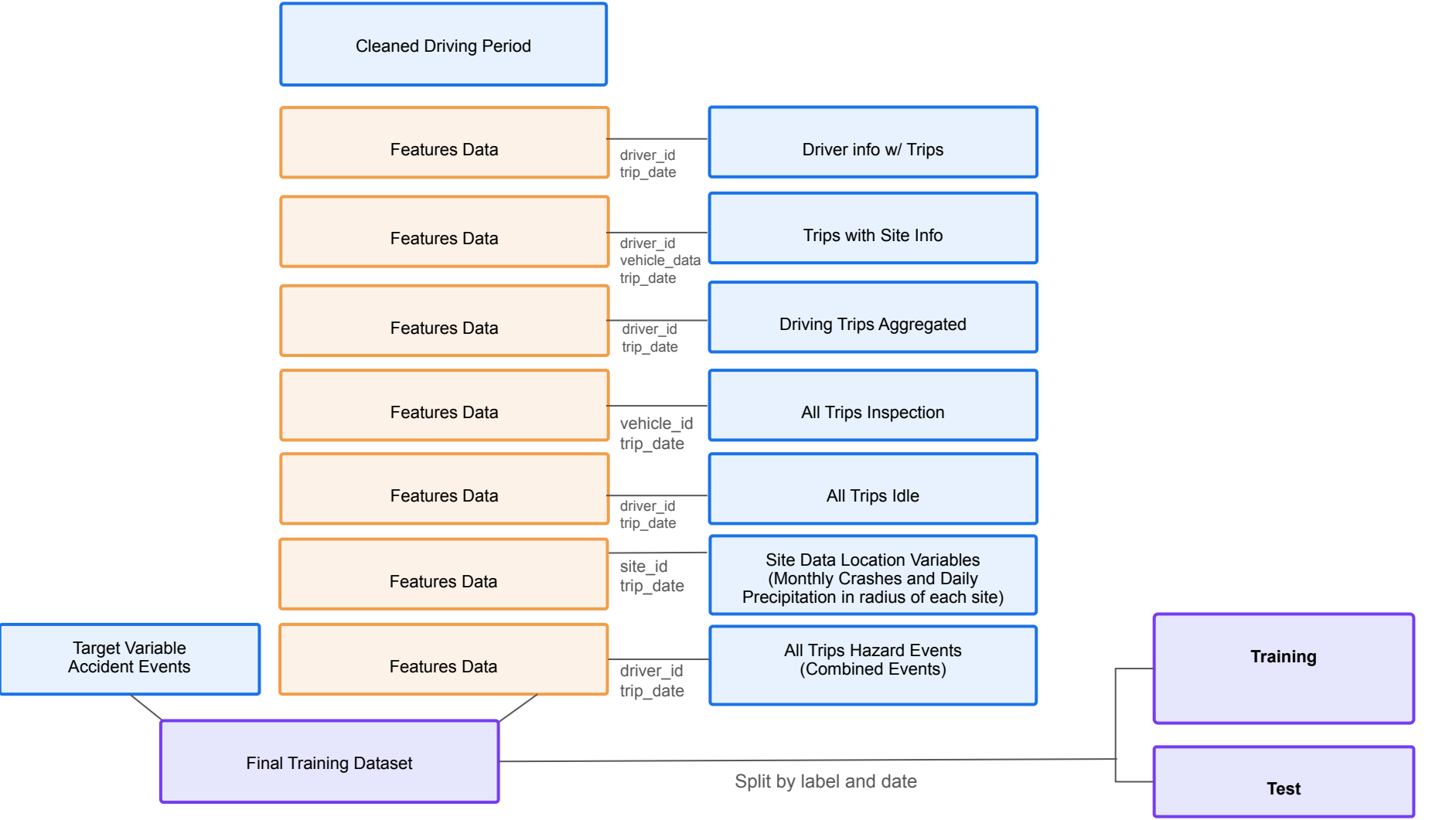
# Joining External Data to Site Radius



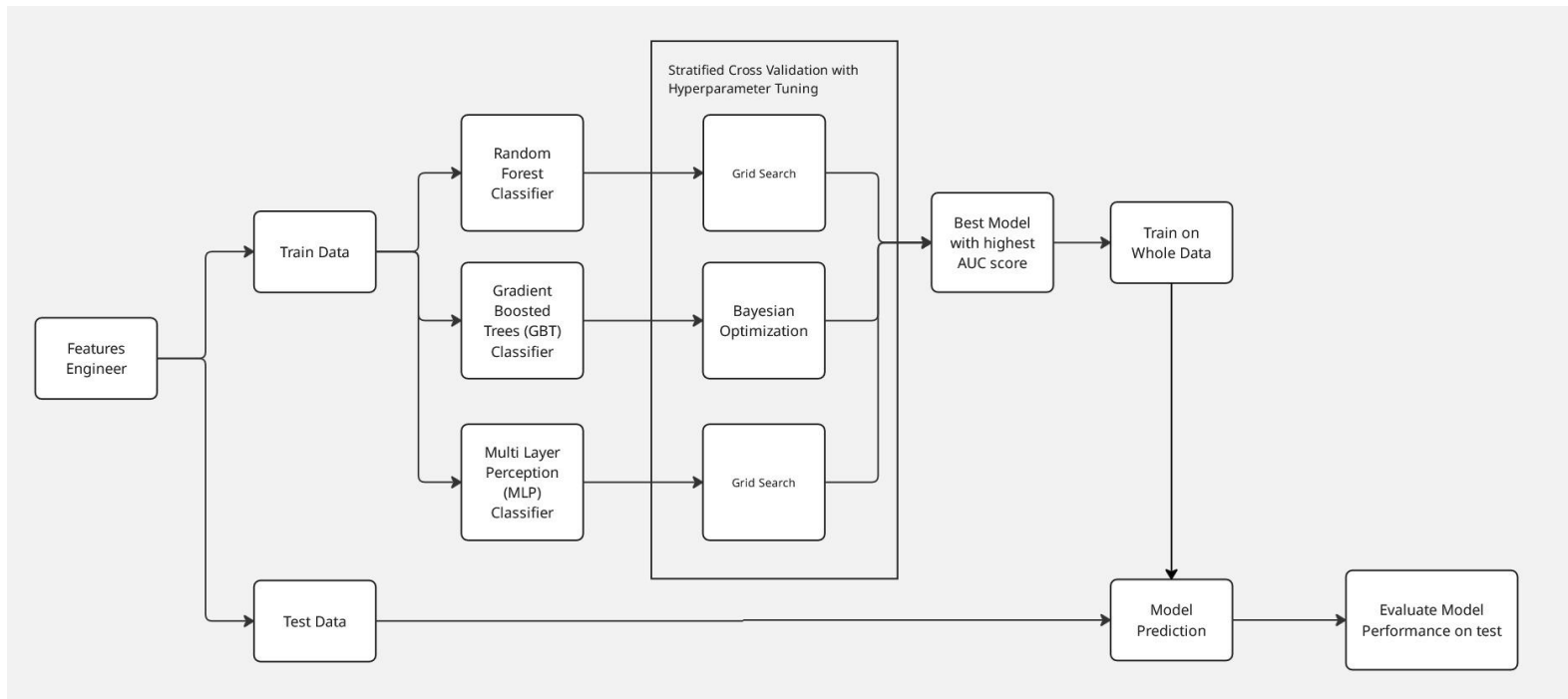
For each site radius we:

- Aggregated details on individual crashes to **monthly crash statistics** across all counties that intersect with the site's service area zip codes
- Aggregated **spatial climate data** to summary statistics for daily precipitation in the service area
- Created **lookback windows** for moving averages of each variable





# Hyperparameter Tuning

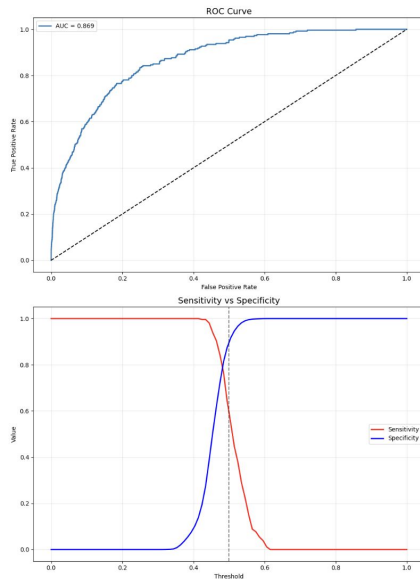


# Results

## Training Data

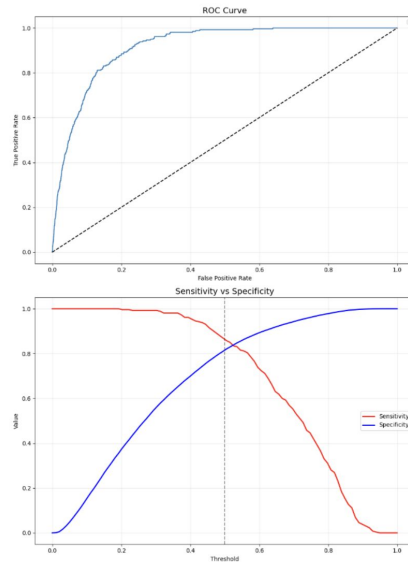
- Our models learned something with AUC values  $> 0.90$  after training, but did they learn the relationship between the features and driver accidents?

Model Evaluation



Random Forest

Model Evaluation

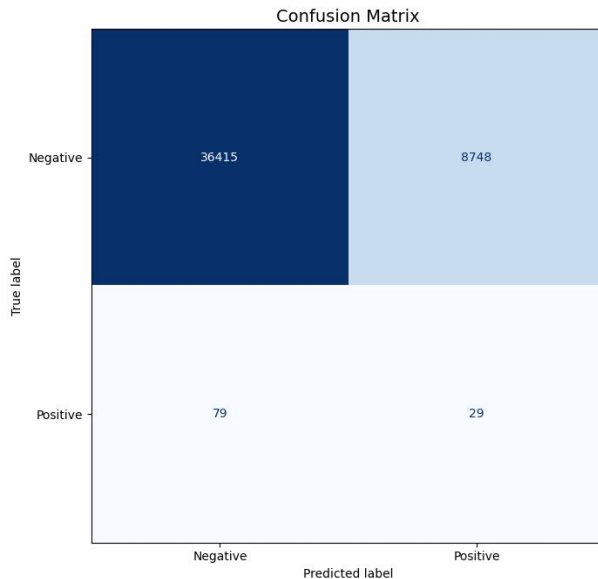


Gradient Boosted Tree

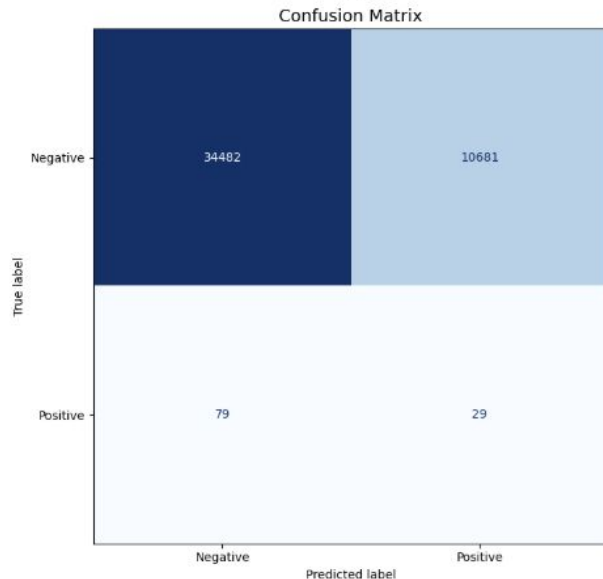
# Results

## Test Data

- Both models performed just marginally better than random with AUCs of 0.57
- The positive predictive value was extremely low at just 0.32% for RF and 0.30% for GBT
- Less than half a percent of positive predictions were actual accidents.
- Only about 18-28% of true accidents while generating thousands of false positives.



Random Forest



Gradient Boosted Tree

# Success

- Served as an informal data audit, uncovering key data inconsistencies
- Provided actionable insights into current data quality and structure
- Helped lay the foundation for future in-house predictive modeling efforts
- Acquired extensive knowledge of the Spark package, recognizing some Spark model limitations like non-stratified cross-validation.
- Fostered smooth collaboration among team members throughout the process.

# Challenges

- Date Imputation
- No guarantee drivers will stay within 40 miles of site for any given trip
- Data not missing completely at random (MCAR)
- Data Completeness in insurance claims
- Extreme Class Imbalance.
- Scalability

# Where to go from here

**Data**

A predictive model can only ever be as good as the data that goes into it.



## **COLLECTION**

Focus on systematic data collection protocols.



## **START SIMPLE**

Start with simple model based on known safety hazards.



## **ITERATE**

Analyze patterns, profiles and behaviors over time.

# Appendix



# Motive Data Preprocessing - Data Driving Events

**Description:** The final table generate a comprehensive feature set that captures each driver's behavior, history, and context on a per-trip basis, The date range of this data is from Jan 2023 to Feb 2025.

## Raw Features

- Id
- Event\_id
- driver\_id
- driver\_first\_name
- driver\_last\_name
- vehicle\_id
- start\_date
- end\_date
- driving\_distance
- Driving\_period\_type
- driver\_company\_id
- minutes\_driving
- created\_at/updated\_at
- number
- status-2
- make
- model



Aggregate by driver and trip date

## Derived Features

- driver\_id
- Trip Date
- driver\_total\_trip\_count
- first\_driving\_date
- driver\_log\_trip\_count
- rolling\_xday\_trip\_count
- previous\_trip\_date
- days\_since\_last\_trip
- driving\_year\_since\_first\_trip

# Motive Data Preprocessing - Data Combined Events

**Description:** This aggregation tables table captures the volume and types of safety-related events (e.g., crashes, distractions, speeding), their severity in term of speeding (low/mid/high). It helps assess driving behavior, identify risky patterns, and monitor unresolved coaching interventions

## Raw Features

- Id
- Event\_id
- driver\_id
- driver\_first\_name
- driver\_last\_name
- vehicle\_id
- coaching\_status
- start\_date
- severity
- group\_id
- group\_name
- month
- created\_at
- updated\_at
- max\_over\_speed\_in\_kph
- max\_over\_speed\_in\_mph



Aggregate by driver and  
start/trip date

## Derived Features

- driver\_id
- Trip Date
- Total\_events\_per\_trip
- total\_<event\_type>\_count
- speeding\_bin
- cnt\_max\_over\_speed\_in\_mph\_per\_speed\_bin
- pending\_review\_count

# Motive Data Preprocessing - Data Vehicle Inspection Events

**Description:** This aggregation creates a trip-level summary of inspection activity per vehicle by counting total inspections and identifying trips with open or resolved issues.

## Raw Features

- id
- event\_id
- driver\_id
- driver\_first\_name
- driver\_last\_name
- vehicle\_id
- start\_date
- end\_date
- driving\_distance
- driving\_period\_type
- driver\_company\_id
- minutes\_driving



Aggregate by vehicle and trip date

## Derived Features

- vehicle\_id
- trip Date
- num\_inspections\_per\_trip
- num\_issues\_per\_trip
- inspection\_date
- last\_inspection\_date
- last\_inspection\_status

# Motive Data Preprocessing - Data Idle Events

**Description:** This aggregation summarizes idle behavior per driver per day. It helps quantify how often and how long a truck is idling — key indicators of driver efficiency, fuel usage, and potential safety concerns.

## Raw Features

- Id
- Event\_id
- driver\_id
- vehicle\_id
- start\_time
- end\_time
- driver\_company\_id
- minutes\_idling

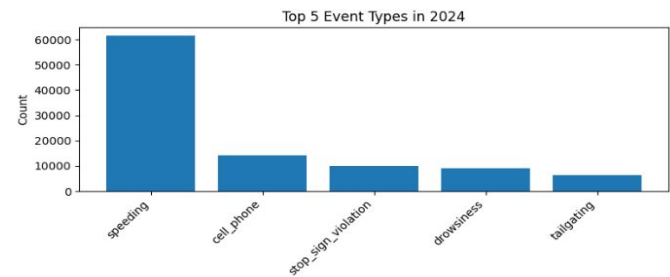
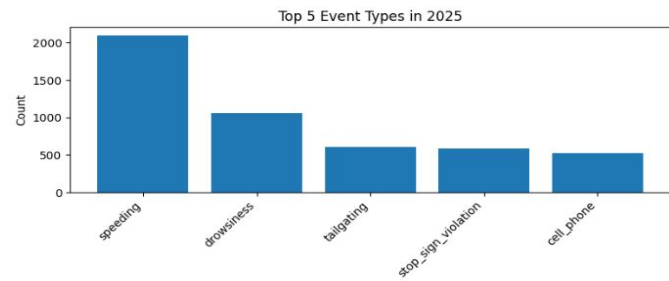
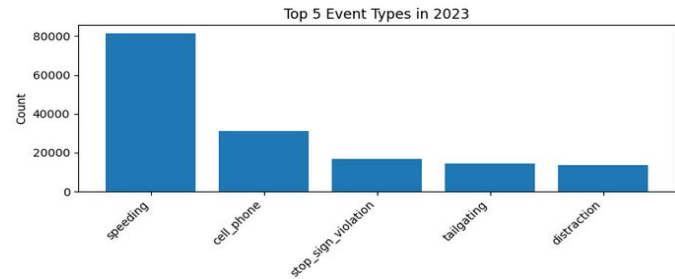
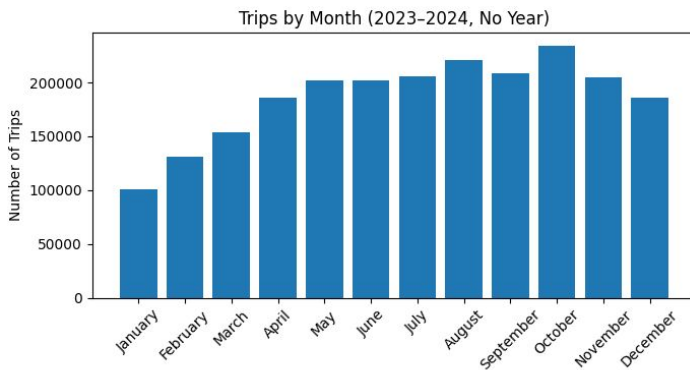
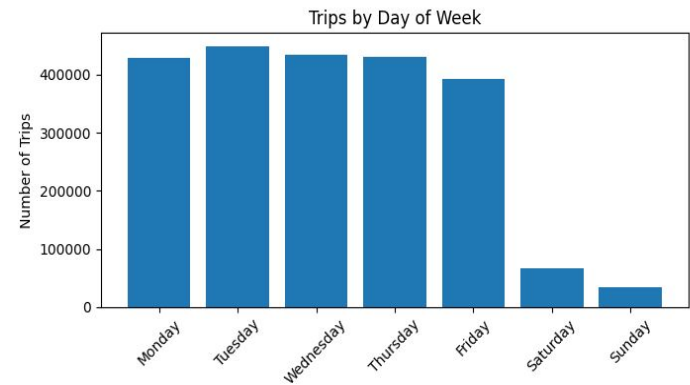


Aggregate by driver and trip date

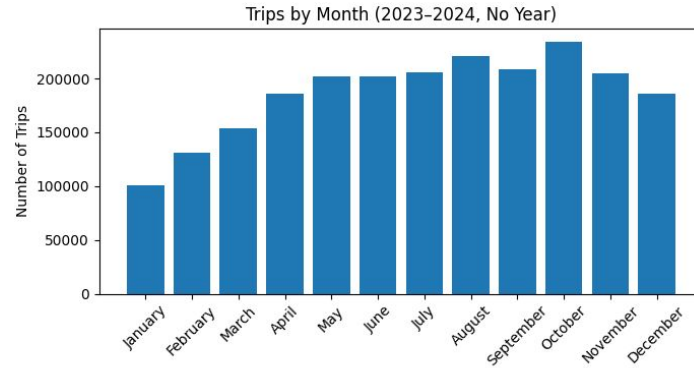
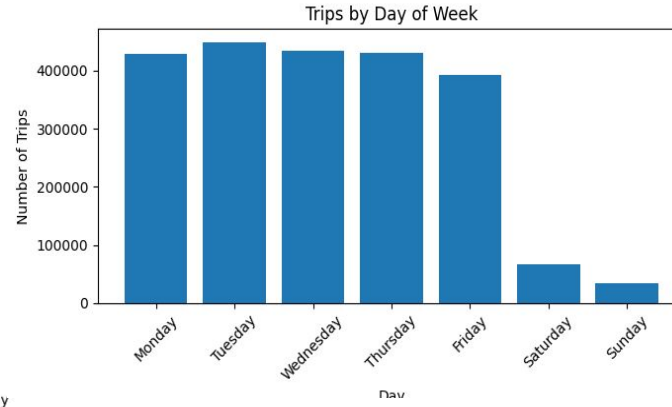
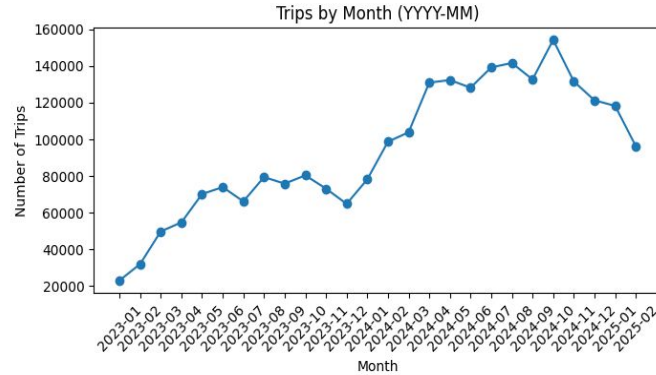
## Derived Features

- driver\_id
- idle\_date
- idle\_event\_count\_per\_trip
- avg\_idle\_duration\_per\_trip
- total\_idle\_minutes\_per\_trip

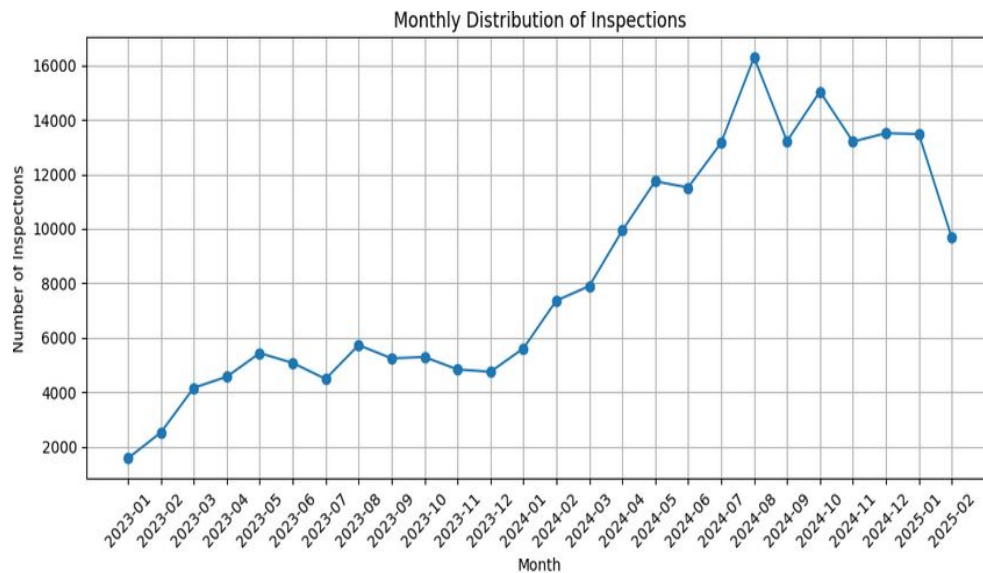
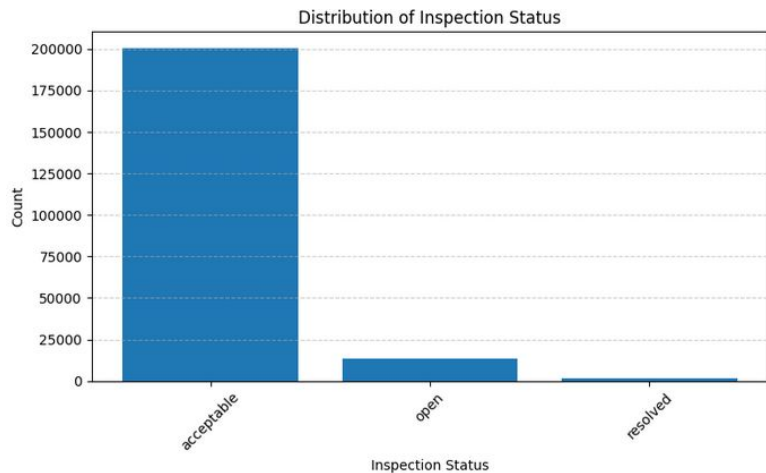
# Motive EDA



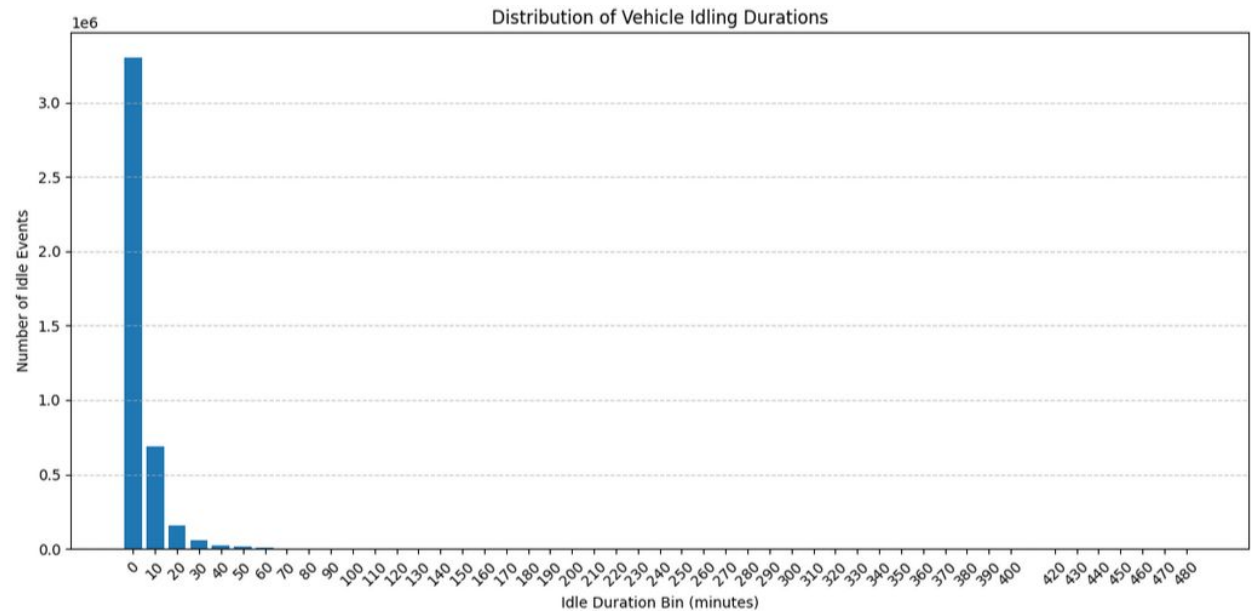
# Motive Data Preprocessing - Data Driving Events



# Motive Data Preprocessing - Data Vehicle Inspection Events



# Motive Data Preprocessing - Data Idle Events



summary	minutes_idling
count	4264634
mean	8.2234
stddev	8.981
min	2.0
max	485.2



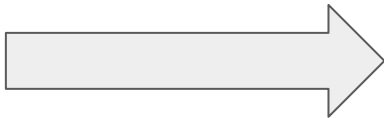
# Site Radius Data

## Raw Features

- ID
- Brand
- State
- City
- Latitude
- Longitude
- County
- Zip

## US Census Zip Code Data

- Zip
- Latitude
- Longitude



- Cross join site and target ZIP codes to create all possible combinations
- Add target coordinates (latitude/longitude) as columns to the DataFrame
- Calculate distances between locations using the Haversine distance
- Filter results to only include locations within 40 miles

## Derived Features

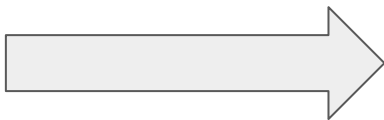
- One DataFrame for each site that contains
  - Zip code
  - Distance in miles
- For all zip codes within a 40 mile radius of each site

# State Crash Data

**Description:** This aggregation summarizes provides a monthly summary of crash-related incidents around each site to help encode external risk factors for accidents based on where drivers will spend most of time on the road.

## Raw Features

- Crash ID
- Crash Year
- Report State
- Fatal Count
- Fatalities
- Injuries
- Vehicles In Accident
- City
- City Code
- State
- Location
- Crash Date
- Year



Aggregate by Year, Month,  
State, County and count and  
sum over columns

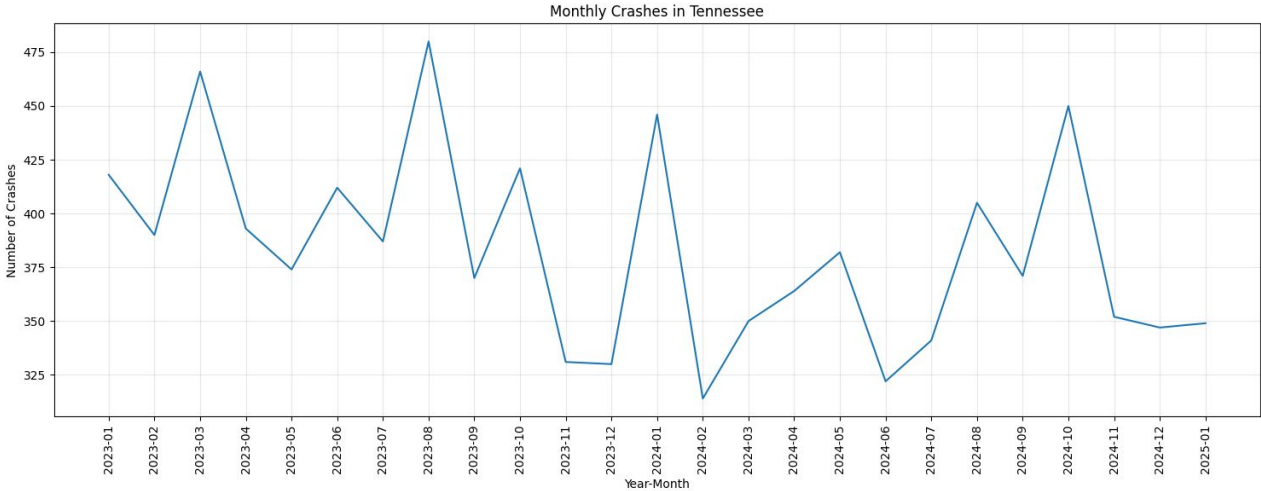
## Derived Features

- Crash Year
- Crash Month
- State
- County Code
- Crash Count
- Total Fatalities
- Total Injuries
- Total Vehicles

# State Crash Data EDA

Correlation matrix for key metrics:

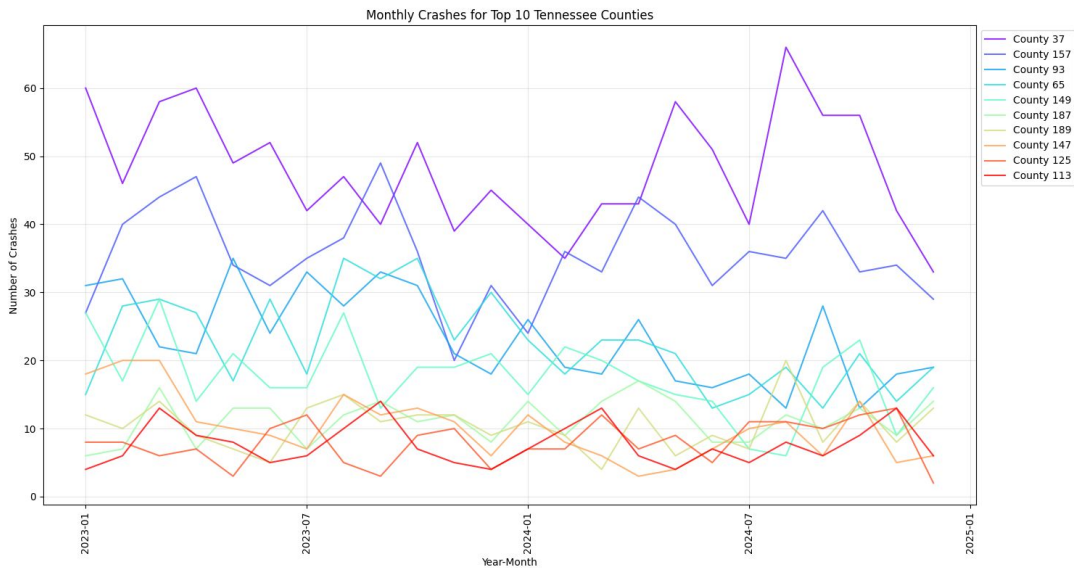
	crash_count	total_fatalities	total_injuries	total_vehicles
crash_count	1.000	0.437	0.851	0.932
total_fatalities	0.437	1.000	0.440	0.481
total_injuries	0.851	0.440	1.000	0.752
total_vehicles	0.932	0.481	0.752	1.000



# State Crash Data EDA

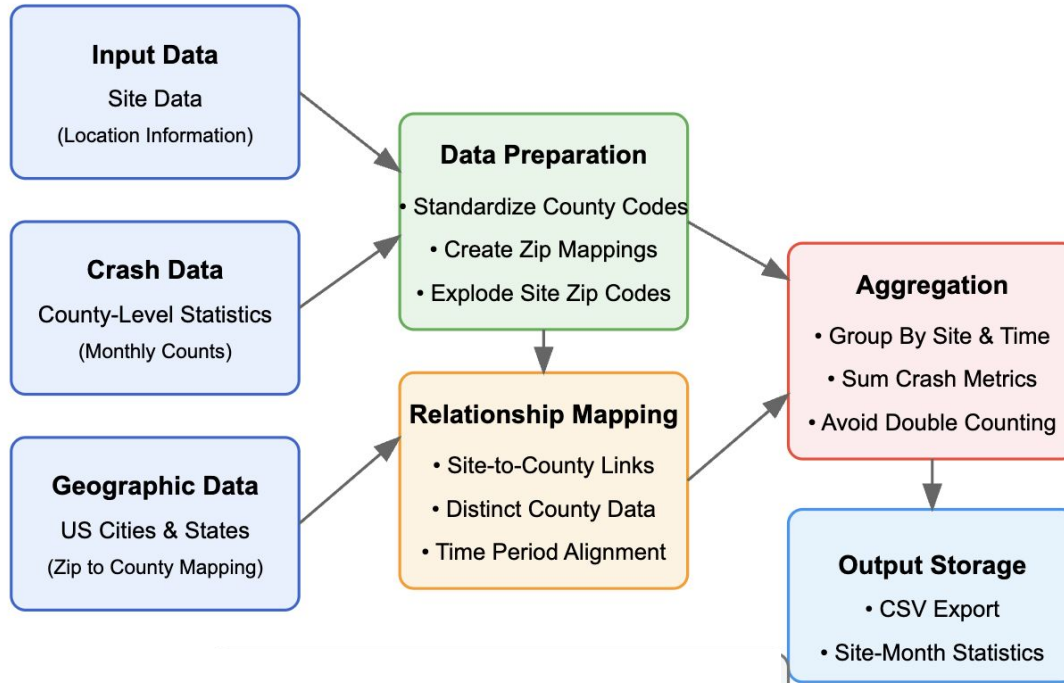
Correlation matrix for key metrics:

	crash_count	total_fatalities	total_injuries	total_vehicles
crash_count	1.000	0.437	0.851	0.932
total_fatalities	0.437	1.000	0.440	0.481
total_injuries	0.851	0.440	1.000	0.752
total_vehicles	0.932	0.481	0.752	1.000



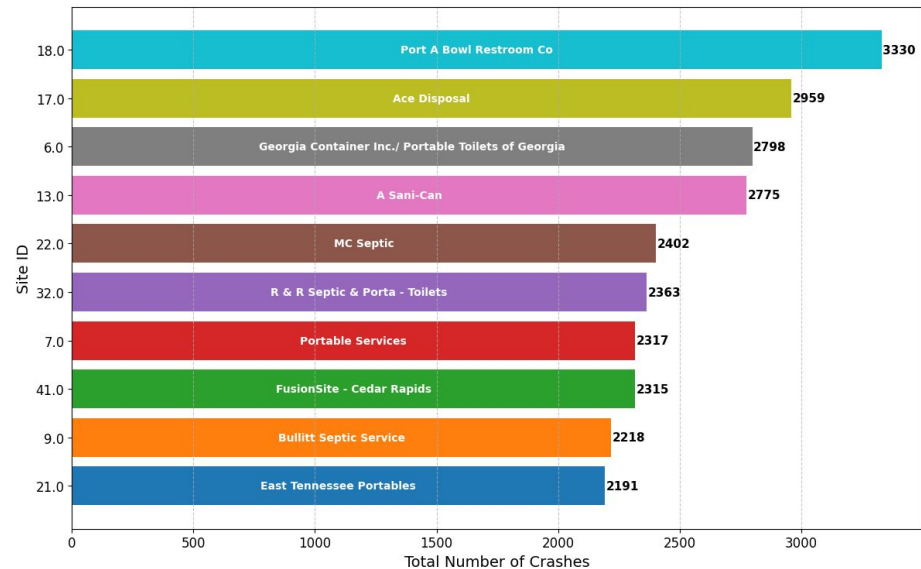
# Joining State Crash Data to Site Data

## Crash Data Integration Schema

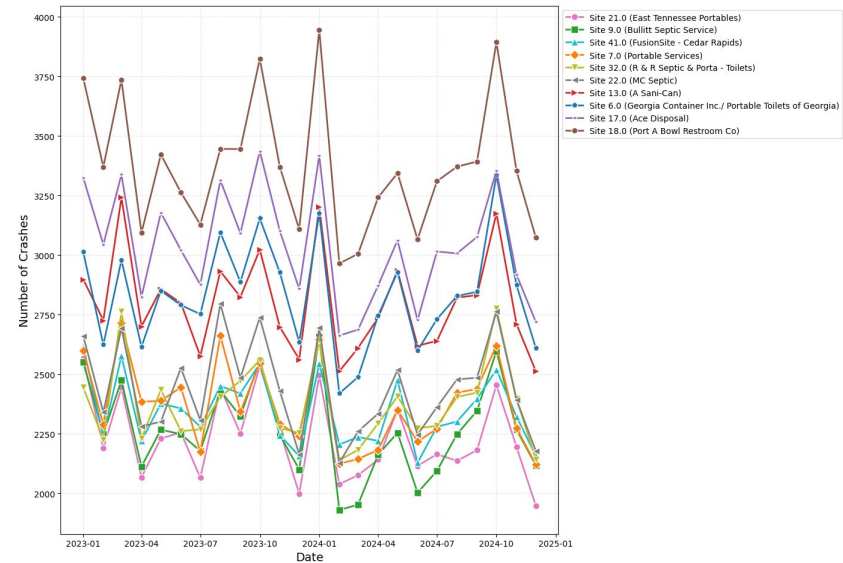


# Site Radius Crash Data EDA

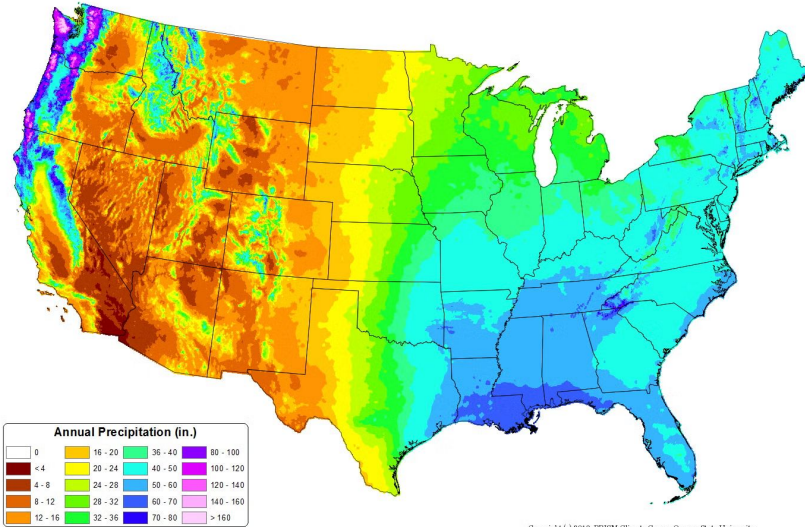
Top 10 Site Radiuses with Highest Average Monthly Crash Count in 2024



Monthly Crash Counts Over Time for Top 10 Sites



# Precipitation Data



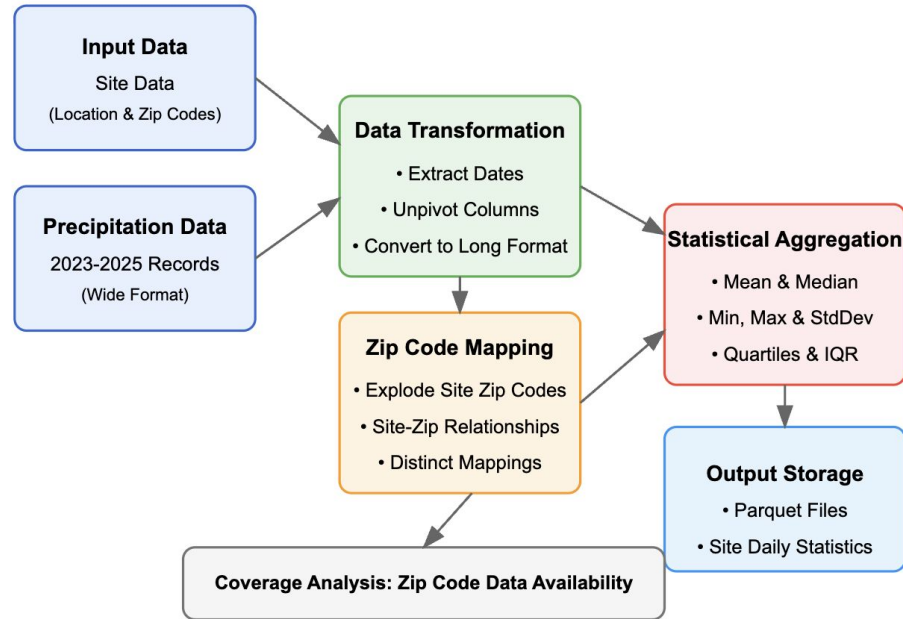
## Raw Features

- Zip Code
- Latitude
- Longitude
- Total precipitation (rain+melted snow) for each Date in mm

*\*See Weather Data Processing for details*

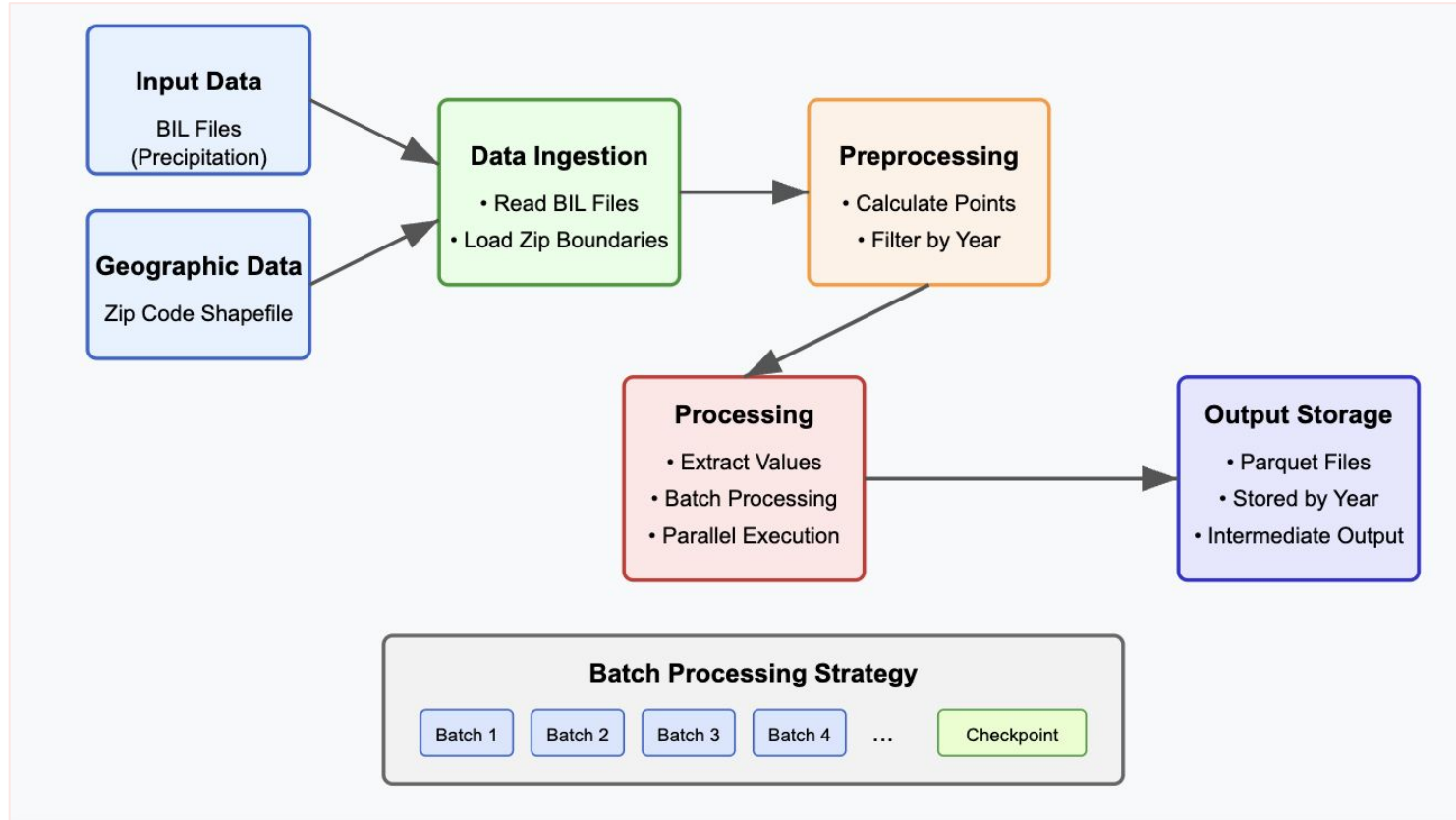
# Joining Precipitation Data to Site Data

## Precipitation Data Integration Schema





# Weather Data Preprocessing



# Weather Data Preprocessing

## **Precipitation Data (BIL Files)**

- Primary raster values
- Raster Metadata
  - Width and height
  - Coordinate reference system (CRS)
  - Transform parameters
  - Bounds
  - Nodata values

## **Zip Code Shapefile**

- Latitude
- Longitude
- Zip Code Identifiers
  - GEOID20
  - ZCTA5CE20

# State Crash Data Preprocessing

## Preprocessing

- Generate representative points for each zip code polygon
- These points are used as sampling locations where values from the raster data will be extracted
- The code converts zip code polygons to specific point coordinates that fall within each zip code area

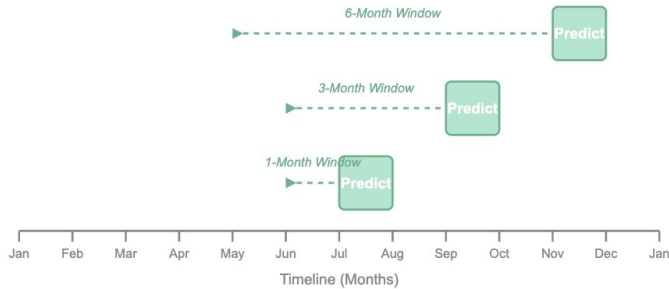
## Processing

- Raster values are sampled at the representative points
- Converts geographic coordinates to pixel coordinates in the raster
- Checks if points fall within the raster bounds and handles nodata values appropriately
- Associates specific BIL data values with each zip code location

# Location Data Feature Engineering

- Created lookback windows at different time scales for each variable

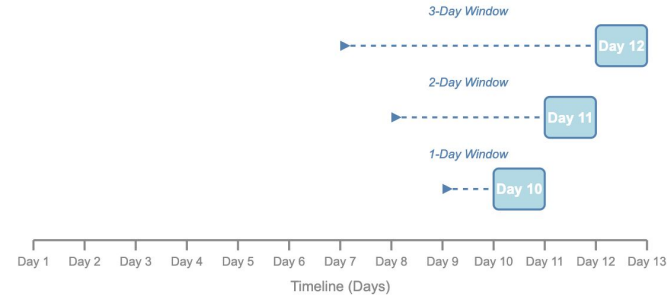
## State Crash Data



### Accident Data Moving Averages

- Total Monthly Crashes
- Total Monthly Fatalities
- Total Monthly Injuries
- Total Monthly Vehicles

## Precipitation Data



### Precipitation Data Moving Averages

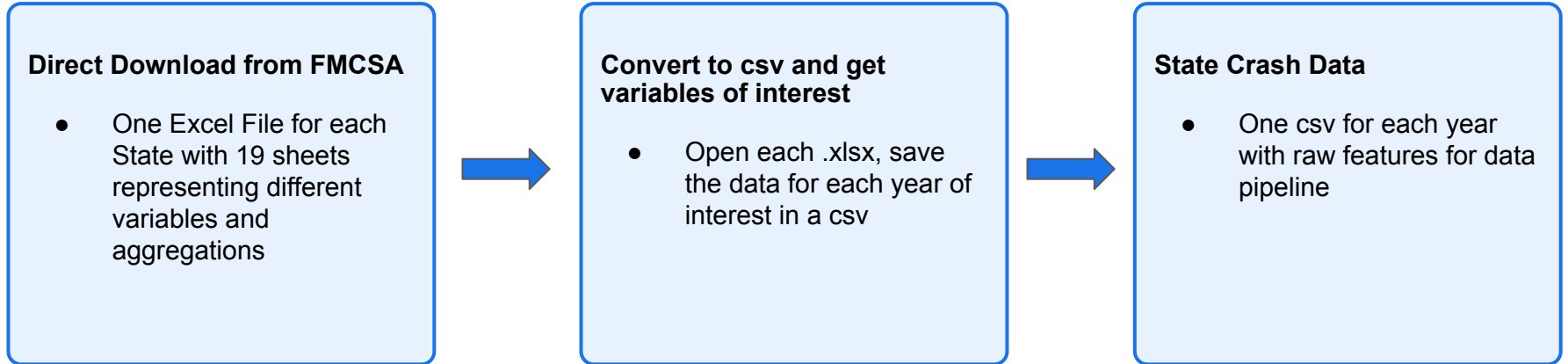
- Summary statistics (min, max, mean, median, iqr) for rainfall over previous days

# Weather Data Preprocessing

## Output

- Zip Code
- Latitude
- Longitude
- Total precipitation (rain+melted snow) for each Date in mm

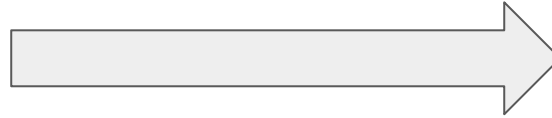
# State Crash Data Data Collection



# Target Variable

Accident Data

Insurance Claims Data



- Identify common columns
- Concatenate all rows
- Delete any duplicates across the driver id and date columns

## All Accidents

- Driver id
- Insured driver name
- Insured driver first name
- Insured driver last name
- Date
- Year