California Self-Driving Vehicle Data Analysis

Kyle Markwardt, DB



The Great Seal California Department

of

Motor Vehicles

State of CA Regulations Info

Regulations v1.0 from 2014-2019 "The Test Drive Era":

- Limited approval as a developing technology
- Inconsistent data reporting standards
- Data available in scholarly research or by request from DMV
 - Not readily publicly available.

Regulations v2.0 from 2020- Today:

- Standardized reporting structure
- Data publicly available from DMV website
- Special numbered permit is issued by company, with special vehicle registration requirements
- 2x Permit Types: <u>Autonomous Testing with a Driver</u> & <u>Driverless Testing</u>
- Each company has slightly different restrictions. All require 'decent' weather and reasonable speeds



The Great Seal California Department

Motor Vehicles

About California Regulations v2.0 Raw Data

- Data is self-reported, submitted on standardized forms, and compiled by the State into .csv. Covers important categories like monthly miles, date and location of disengagement, driver present
- Yearly <u>Disengagement Reports</u> and <u>Mileage Reports</u> provided by State
- Raw data covers 18 different .csv files, ~20,000 overall rows, ~250,000 overall cells, ~2,000 vehicles
- Each individual Crash Report PDF provided by State (but uncompiled) (~700)

Processing California Regulations v2.0 Raw Data

- Varieties of highly similar terms ("Highway", "Freeway") ("Computer Disengagement" "AV System Disengagement") standardized into main categories
- When possible, VIN problems were cleaned by finding correct VIN in a different report
- Data was combined & cross-referenced to enable clear multi-year insights for VINs, Companies, Dates, etc.

Data Sources

Mileage Reports

Disengagement Reports

- Vehicle in autonomy
- Comes out
- Regardless of driver

CSVs per year

PORTS



Autonomous vehicle manufacturers that are testing vehicles in the Autonomous Vehicle Tester (AVT) Program and AVT Driverless Program are required to submit annual reports to share how often their vehicles disengaged from autonomous mode during tests (whether because of technology failure or situations requiring the test driver/operator to take manual control of the vehicle to operate safely).

To request previous disengagement reports, please email AVarchive@dmv.ca.gov.

2023 Disengagement Reports

2022 Disengagement Reports

2022 Autonomous Vehicle Disengagement Reports (CSV)

2022 Autonomous Mileage Reports (CSV)

2021-22 Autonomous Vehicle Disengagement Reports (CSV) (first-time filers)

2021-22 Autonomous Mileage Reports (CSV) (first-time filers)

Report Form

SECTION 1 — MANUFACTURER INFORMATION											
NAME OF MANU					AVT NUMBER	(8)					
BUSINESS MAILING ADDRESS			CITY	STATE	ZIP CODE	TELEPHONE NUMBER					
SECTION :	- DISENGAGEMENT E	VENT DETAIL Use	one row for each disengageme	ent event.							
DATE	VIN NUMBER	DISENGAGEMENT INITIATED BY (AV System, Test Driver, Remote Operator, or Passenger)	DISENGAGEMENT LOCATION (Interstate, Freeway, Highway, Rural Road, Street, or Parking Facility)	DESCRIPTION OF FACTS CAUSING DISENGAGEMENT *							
	VEHICLE IS CAPABLE OF OPERATING	DRIVER PRESENT									
	WITHOUT A DRIVER YES NO	YES NO									
	VEHICLE IS CAPABLE OF OPERATING	DRIVER PRESENT									
	WITHOUT A DRIVER YES NO	YES NO									
	VEHICLE IS CAPABLE OF OPERATING WITHOUT A DRIVER YES NO	DRIVER PRESENT									
	MINOUR DRIVER										
	VEHICLE IS CAPABLE OF OPERATING WITHOUT A DRIVER YES NO	DRIVER PRESENT									
	VEHICLE IS CAPABLE OF OPERATING WITHOUT A DRIVER YES NO	DRIVER PRESENT									
	VEHICLE IS CAPABLE OF OPERATING WITHOUT A DRIVER ☐ YES ☐ NO	DRIVER PRESENT									
	VEHICLE IS CAPABLE OF OPERATING WITHOUT A DRIVER YES NO	DRIVER PRESENT									
	VEHICLE IS CAPABLE OF OPERATING WITHOUT A DRIVER YES NO	DRIVER PRESENT									

VIN Number	Annual Total of	Total Number of Miles Tested in Autonomous Mode (December to November)													
VIN Number	Disengagements	December	January	February	March	April	May	June	July	August	September	October	November	ANNUAL TOTAL	
														0	
														0	
														0	
														0	
														0	
														0	
														0	
														0	
														0	
														0	
														0	
														0	
														0	
														0	
														0	
														0	
ECTION 4 — ACKN	OWLEDGMENT			1										7.0	
PRINTED NAME OF AUTHORIZED REPRESENTATIVE								TITLE	TITLE						
SIGNATURE X							DATE SIGNED.								
TREET ADDRESS						CITY					STATE	Z	P CODE		
EMAILADDRESS							FAX NU	MBER		TELEPHO	TELEPHONE NUMBER				

Data Processing

Data Cleaning and QA

Drop Columns

Locations: 15 - ['Express Way', 'Freeway', 'Freeway', 'HIGHWAY', 'HIghway', 'Highway', 'Interstate', 'Parking Facility', 'Parking facility', 'Rural Roac

```
Operators: 51 - ['AIMOTIVE INC.', 'APOLLO AUTONOMOUS DRIVING USA LLC', 'APPLE INC.', 'ARGO AI, LLC', 'AURORA OPERATIONS, INC.', 'AUTOX TECHNOLOGIES, INC'

Cleaned locations: 8 - ['EXPRESS WAY', 'FREEWAY', 'HIGHWAY', 'INTERSTATE', 'PARKING FACILITY', 'RURAL ROAD', 'STREET', 'URBAN']

Cleaned operators: 38 - ['AIMOTIVE', 'APOLLO', 'APOLLO AUTONOMOUS DRIVING USA', 'APPLE', 'ARGO AI', 'AURORA OPERATIONS', 'AUTOX TECHNOLOGIES', 'BOSCH', '

Original number of operators: 50

Actual number of unique operators: 32

'AIMOTIVE'

'APOLLO'

'APPLE'

'ARGO AI'

'AURORA OPERATIONS'

'AUTOX'

'BOSCH'

'BOSCH'
```

'CRUISE'

Data Sources

Collision reports

Individual PDFs per incident

Access problem:

- No API
- Parse failed (href)
- Clicker failed (JS)

Manual download

TONOMOUS HICLE COLLISION PORTS



Autonomous Vehicle Regulations
us Vehicle Definitions
us Vehicles Tests without a Driver
us Vehicles Testing with a Driver
us Vehicle Deployment Program
us Vehicle Testing Permit Holders
ment Reports

Manufacturers who are testing autonomous vehicles need to report any collision that resulted in property damage, bodily injury, or death within 10 days of the incident.

As of March 22, 2024, the DMV has received 698 Autonomous Vehicle Collision Reports.

Collision reports prior to January 1, 2019 have been archived by DMV and are available upon request. Please email Marchive@dmv.ca.gov to request a digital copy of an archived report. Requests must include the manufacturer and the date of the collision. Please do not include any sensitive personal information such as your social security number, driver license number, or financial account number on the request.

2024 +

- Mercedes Benz, November 29, 2023 (PDF)
- Nuro November 27, 2023 (PDF)

2023

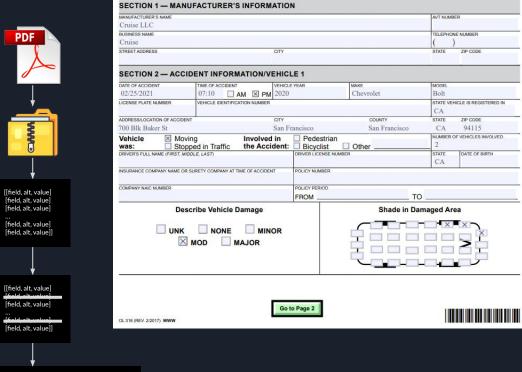
- Waymo November 22, 2023 (PDF)
 - Waymo November 3, 2023 (PDF)

PDF Mining

- 1. Download PDFs
- 2. Zip for memory (568 files)
- 3. Read fields pdfplumber

4. Select wanted fields

5. Create DF





Scraping User-Entered data from PDF

```
# This function read the filed name, alternate field name, and field values
def parse field helper(form data, field, prefix=None):
    """ appends any PDF AcroForm field/value pairs in `field` to provided `form data` list
        if `field` has child fields, those will be parsed recursively.
    resolved field = field.resolve()
    field name = '.'.join(filter(lambda x: x, [prefix, resolve and decode(resolved field.get("T"))]))
   if "Kids" in resolved field:
        for kid field in resolved field["Kids"]:
            parse field helper(form data, kid field, prefix=field name)
    if "T" in resolved field or "TU" in resolved field:
        # "T" is a field-name, but it's sometimes absent.
        # "TU" is the "alternate field name" and is often more human-readable
        # your PDF may have one, the other, or both.
        alternate field name = resolve and decode(resolved field.get("TU")) if resolved field.get("TU") else None
        field value = resolve and decode(resolved field["V"]) if 'V' in resolved field else None
        # Remove non-printable characters and trailing spaces - This affects every Cruise file.
        field name = ''.join(char for char in field name if char.isprintable()).strip()
        alternate field name = ''.join(char for char in alternate field name if char.isprintable()).strip() if alternate field name else None
        field value = ''.join(char for char in field value if char.isprintable()).strip() if field value else None
        form data.append([field name, alternate field name, field value])
```

```
Returns a list of tuples, e.g: [[field_name_1, alt_name_1, value_1], ... [field_name_n, alt_name_n, value_n]]
```

Filter for what is interesting

```
# Define filtering criteria
def find_important_tuples(tuple, search_condition):
    # Check if the first element of the tuple is in the search condition list
    return tuple[0] in search_condition
```

Manually checked whole list for fields we wanted (\sim ½)

Created list to filter and some sub-lists

```
all fields = ['MANufACTURERS NAME', BUSINESS NAME',
              'DATE Of ACCIDENT', 'Time of Accident', 'AM', 'PM',
              'VEHICLE YEAR', 'MAKE', 'MODEL',
               'section 2 accident infoRmation.0','section 2 accident infoRmation.1.0','section 2 accident infoRmation
              'Moving', 'Stopped in Traffic', 'Pedestrian', 'Bicyclist', 'undefined', 'Other',
              'NuMBER OF VEHICLES INVOLVED',
              'Unknown', 'None', 'minor', 'Moderate', 'major',
              Left Rear 1', 'Rear Bumper', 'Right Rear 1', 'Left Rear 2', 'Left Rear 3', 'Right Rear 2', 'Right Rear 3',
              'Left Rear Passenger 1', 'Left Rear Passenger 2', 'Right Rear Passenger 1', 'Right Rear Passenger 2',
              'Left Rear Passenger 3', Left Rear Passenger 4', Right Rear Passenger 3', Right Rear Passenger 4',
              'Front Driver Side 1', 'Front Driver Side 2', 'Front Passenger Side 1', 'Front Passenger Side 2',
              'Front Driver Side 3', 'Front Driver Side 4', 'Front Passenger Side 3', 'Front Passenger Side 4',
              'Left Front Corner 1', 'Left Front Corner 2', 'Right Front Corner 1', 'Right Front Corner 2',
              'Left Front Corner 3', 'Front Bumper', 'Right Front Corner 3',
              'Moving_2', 'Stopped in Traffic_2', 'Pedestrian_2', 'Bicyclist_2', 'undefined_2', 'Other_2',
               'ADDRESS 2.1.0.1', 'Autonomous Mode', 'Conventional Mode',
              'WEATHER A 1', 'WEATHER A 2', 'WEATHER B 1', 'WEATHER B 2', 'WEATHER C 1', 'WEATHER C 2',
              'WEATHER D 1', 'WEATHER D 2', 'WEATHER E 1', 'WEATHER E 2', 'WEATHER F 1', 'WEATHER F 2', 'WEATHER G 1', 'WEATHER
              'LIGHTING A 1', 'LIGHTING A 2', 'LIGHTING B 1', 'LIGHTING B 2', 'LIGHTING C 1', 'LIGHTING C 2',
              'LIGHTING D 1', 'LIGHTING D 2', 'LIGHTING E 1', 'LIGHTING E 2',
               'ROADWAY A 1','ROADWAY A 2','ROADWAY B 1','ROADWAY B 2','ROADWAY C 1', 'ROADWAY C 2','ROADWAY D 1','ROADWA
              'ROAD CONDITIONS A 1', 'ROAD CONDITIONS A 2', 'ROAD CONDITIONS B 1', 'ROAD CONDITIONS B 2', 'ROAD CONDITIONS
              'ROAD CONDITIONS D 1', 'ROAD CONDITIONS D 2', 'ROAD CONDITIONS E 1', 'ROAD CONDITIONS E 2', 'ROAD CONDITIONS
              'ROAD CONDITIONS G 1', 'ROAD CONDITIONS G 2', 'ROAD CONDITIONS H 1', 'ROAD CONDITIONS H 2',
              'MOVEMENT A 1', 'MOVEMENT A 2', 'MOVEMENT B 1', 'MOVEMENT B 2', 'MOVEMENT C 1', 'MOVEMENT C 2', 'MOVEMENT D
              'MOVEMENT I 1', 'MOVEMENT I 2', 'MOVEMENT J 1', 'MOVEMENT J 2', 'MOVEMENT K 1', 'MOVEMENT K 2', 'MOVEM
              'MOVEMENT M 1', 'MOVEMENT M 2', 'MOVEMENT N 1', 'MOVEMENT N 2', 'MOVEMENT O 1', 'MOVEMENT O 2',
              'MOVEMENT P 1', 'MOVEMENT P 2', 'MOVEMENT Q 1', 'MOVEMENT Q 2', 'MOVEMENT R 1', 'MOVEMENT R 2',
              'TYPE A 1', 'TYPE A 2', 'TYPE B 1', 'TYPE B 2', 'TYPE C 1', 'TYPE C 2', 'TYPE D 1', 'TYPE D 2', 'TYPE E 1', 'TYP
              'OTHER A YES', 'OTHER A NO', 'OTHER B', 'OTHER C', 'OTHER D', 'OTHER E', 'OTHER F', 'OTHER G',
              'OTHER H YES', OTHER H NO', OTHER I', OTHER J', OTHER K', OTHER L']
```

PDF Megafunction

- Loop through all files in location (zip)
- 2. Call the scrape and filter functions
- 3. Make each pdf into one row
- 4. Concatenate

```
# Loop over all pdfs and add entries. Keyword can search foles for specific name
def extract from zip(zip file path, list of pdf fields, keyword=None):
   collisions = []
   counter = 0
    with zipfile.ZipFile(zip file path, 'r') as zip ref:
        for filename in zip ref.namelist():
           # Check if the current item is a pdf file
           if filename.endswith('.pdf') and (keyword is None or keyword in filename):
                print(f'Extracting: {filename}...')
               # Read PDF from zip file
               with zip_ref.open(filename, 'r') as pdf_file:
                   pdf_data = io.BytesIO(pdf_file.read())
               # Open each pdf
                pdf = pdfplumber.open(pdf data)
                form data = []
                fields = resolve(pdf.doc.catalog["AcroForm"])["Fiel
                # For each field, run the pdf parsing funct
                                                                                    add it to form data list
                                                                o extract adta an
                for field in fields:
                    parse field helper(form data, field)
                # Filter the long list of tuples [all_fields, geo_only, cond lions_only, damage_only]
               filtered list = [tuple for tuple in form data if find important tuples(tuple, list of pdf fields)]
                data dict = {}
                # Set df sturcture so each pdf is one row - alt text is column name, value is vlaue]
                # Populate the dictionary with values from filtered list
                for tuple in filtered list:
                    column name = tuple[1]
                    row value = tuple[2]
                    data dict[column name] = row value
               # Create DataFrame from the dictionary
               collision_report = pd.DataFrame([data_dict])
               collisions.append(collision report)
               print('Done.')
               counter += 1
   print(f"Extracted data from {counter} collision reports.")
   df = pd.concat(collisions)
   df.reset index(drop=True, inplace=True)
   return df
```

Result!

Additional challenges:

CRUISE!! - non-printable characters EVERY field

Checkboxes:

- X = "", "Yes"
 - o Change to 1
- Blank = None, NaN, "Off"
 - o Change to 0

Timestamps - 12hr/24hr

```
all collisions df = extract from zip(path to zip, all fields)
   all collisions df.head()

√ 4m 16.6s

Extracting: Aimotive 091619.pdf...
Extracting: Apollo-OL316-062923-Redacted.pdf...
Done.
Extracting: Apollo-OL316-090623-Redacted.pdf...
Done.
Extracting: Apollo-OL316-101623-Redacted.pdf...
Done.
Extracting: Apple 012624.pdf...
Done.
Extracting: Apple 022123.pdf...
Extracting: apple 081921.pdf...
Done.
Extracting: apple 082321.pdf...
Done.
Extracting: Apple 091919.pdf...
Done.
Extracting: Apple 101022.pdf...
Extracting: Apple 120621.pdf...
Done.
Extracting: Apple 122022.pdf...
Done.
Extracting: Apple OL316 061422 Redacted.pdf...
Done.
Extracting: Zoox-OL316-101223-Redacted.pdf...
Done.
Extracted data from 568 collision reports.
```

path to zip = "data/collisions/Collision PDFs.zip"

read in all the data once.

Self-Driving Trends in California

2021-2023 Total Autonomous Miles in CA:

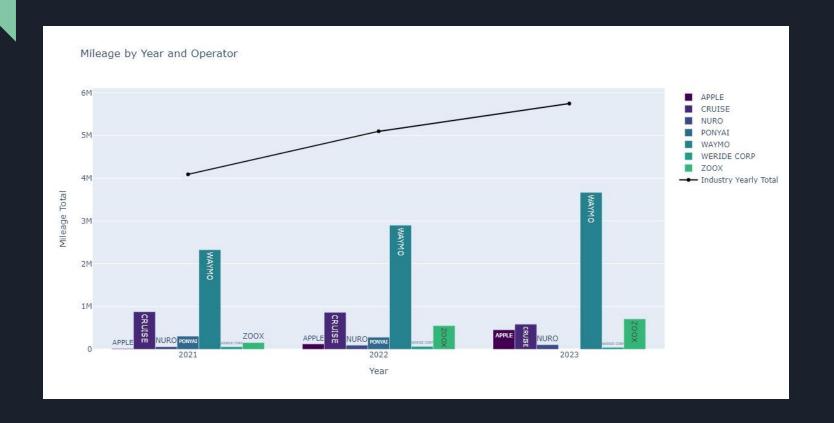
14,938,551

US-20 trips (3300m) equivalent:

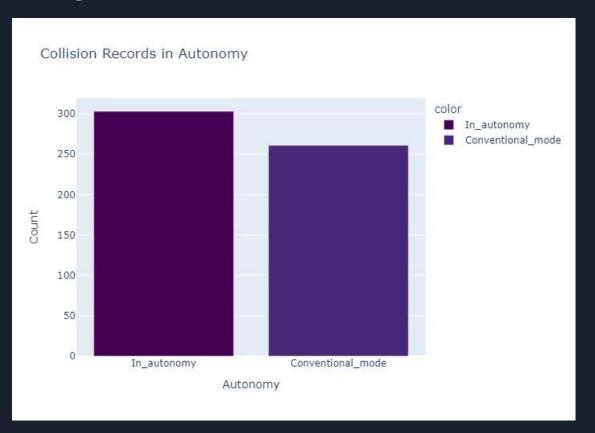
Over 4,500



Top Company Annual Mileage



Things break down



Average Rate of Autonomous Collisions

All Company aggregate:

Autonomous collisions:

~25/mo

2023 autonomous mileage:

5.75M mi

Per 100M mi:

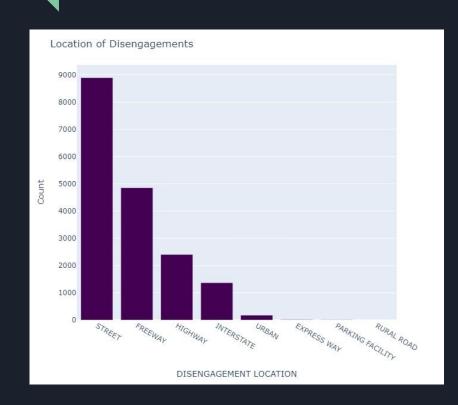
~5200

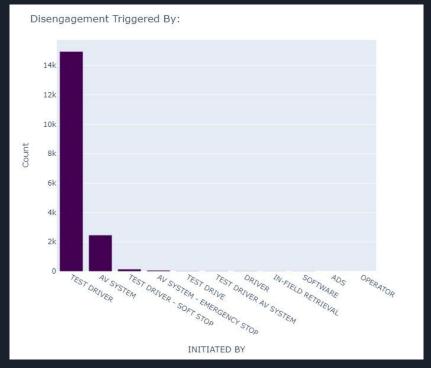


What Time Did Collisions Occur?

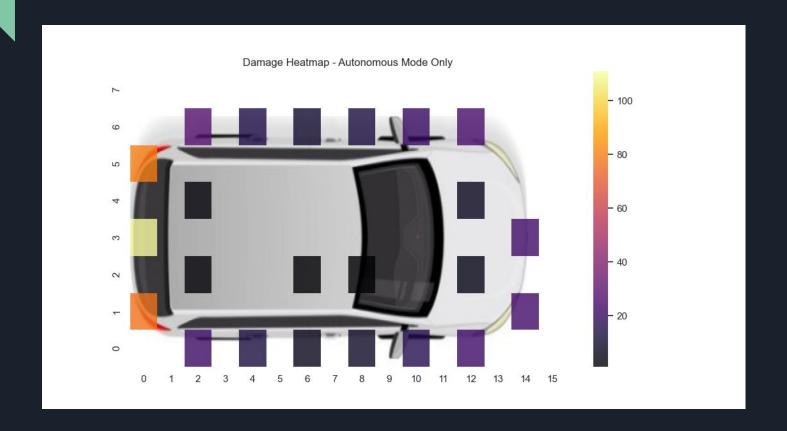


Disengagement Insights





Collision heatmap



Waymo Gets There*

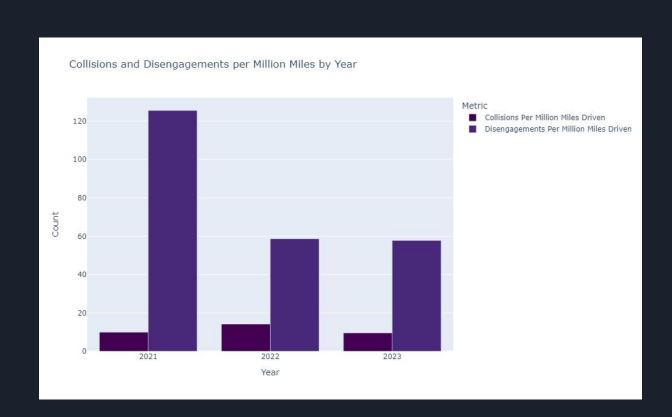
NHTSA - 2021

Crashes/100M miles:

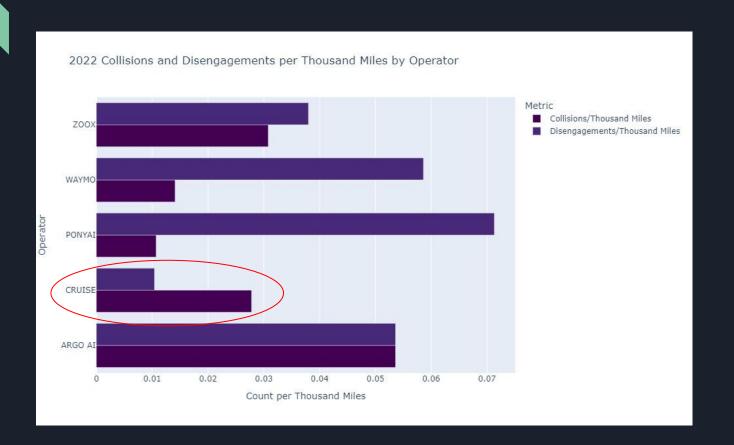
195

Waymo:

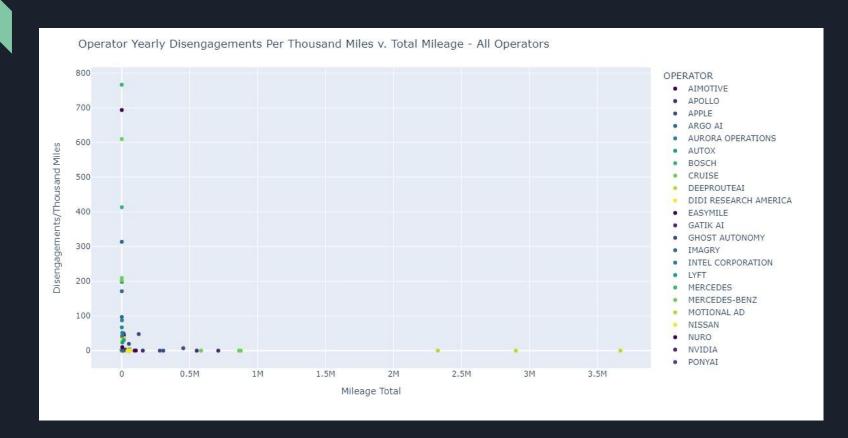
~1400



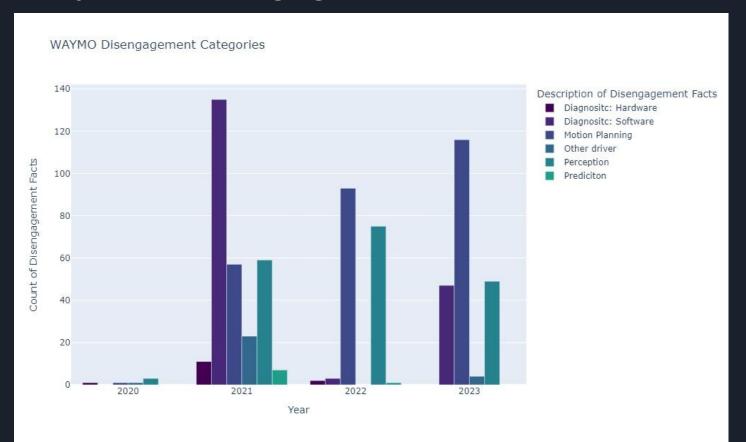
Broader View



Disengagement Rates vs Miles Logged



Waymo - Disengagement Causes



Disengagement Insights



State of California

BAD:

- No one appears to be checking this data
 - Duplicate VINs, wrong format
 - Entries like "same"
 - Total lack of data validation
- Unlikely outliers:
 - E.g., Cruise lack of disengagement reports. Eye opening delta.
 - Could it be real?
- What role do crashes play in relationship to the processed dataset?
 - The State hasn't compiled crashes in an accessible raw data format.

GOOD:

- Allowed Waymo to get there
- Other states more friendly to self-driving because of CA, but are not publishing this data
- Surprisingly accessible data for the DMV

Wish List: Zippy & Pokey

Data is granular enough (VIN level) we would have liked to tell the story of the highest and lowest mileage vehicles. Several vehicles have interesting stories:

Zippy VIN is SADHW2513M1616427 Zippy is a WAYMO vehicle, which drove 55273.1 miles. Pokey VIN is ~5071 Pokey is a ZOOX vehicle, which drove 0.0 miles.





Other Directions We Could Go



- Other states
 - Some operators (e.g., Aurora, do not primarily work in CA.)
- Geo-map the collision locations
 - We were unable due to the nature of the locations given (x street near y blvd, for example)
- Other environments (non-urban)
- NLP and Disengagement Reports:
 - Each disengagement report (16,000+) has a short written narrative from the company. What would NLP analysis of this reveal?