

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)

Data Process

```
Locations: 15 - ['Express Way', 'Freeway', 'Freeway', 'HIGHWAY', '
Operators: 51 - ['AIMOTIVE INC.', 'APOLLO AUTONOMOUS DRIVING USA LLC', 'APPLE INC
Cleaned locations: 8 - ['EXPRESS WAY', 'FREEWAY', 'HIGHWAY', 'INTERSTATE', 'PARKI
Cleaned operators: 38 - ['AIMOTIVE', 'APOLLO', 'APOLLO AUTONOMOUS DRIVING USA', '
Original number of operators: 50
Actual number of unique operators: 32
  'AIMOTIVE'
  'APOLLO'
  'APPLE'
  'ARGO AI'
  'AURORA OPERATIONS'
  'AUTOX'
  'BOSCH'
  'CRUISE'
                  # What are "Unnamed" colums? Probably useless "reserved"
                  unnamed_cols = ['Unnamed: 17', 'Unnamed: 18', 'Unnamed: 19', 'Unnamed: 20', 'Unnamed: 21']
                  for check in unnamed cols:
                             non null count = mileage df[check].notna().sum()
                             print(f'Number of non values in {check}: {non null count}')
                  # Empty. Drop 'em
                 mileage_df = mileage_df.drop(columns=unnamed_cols)
                  mileage df.columns

√ 0.0s

         Number of non values in Unnamed: 17: 0
         Number of non values in Unnamed: 18: 0
         Number of non values in Unnamed: 19: 0
```

Number of non values in Unnamed: 20: 0 Number of non values in Unnamed: 21: 0

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

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 AVarchive@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);



	SECTION 1 — MANUFACTURER'S INFORMATION MANUFACTURER'S NAME					LAVT NUMBER	
Cruise LLC					AVINOMBER		
BUSINESS NAME					TELEPHONE NUMBER		
Cruise					()		
STREET ADDRESS		CITY	СПУ		STATE	ZIP CODE	
SECTION 2 - ACC	DENT INFORMATION/\	/EHICL	.E 1				
ATE OF ACCIDENT TIME OF ACCIDENT		VEHICLE	YEAR	MAKE	MODEL		
02/25/2021 07:10 AM X PM		M 2020		Chevrolet	Bolt		
LICENSE PLATE NUMBER VEHICLE IDENTIFICATION NUMBER		3			STATE VEHICLE IS REGISTERED IN CA		
ADDRESS/LOCATION OF ACCIDENT		CITY		COUNTY	STATE	ZIP CODE	
700 Blk Baker St	00 Blk Baker St San		ancisco	San Francisco	CA	94115	
ehicle Moving Involved in ras: Stopped in Traffic the Acciden					NUMBER OF VEHICLES INVOLVED 2		
DRIVER'S FULL NAME (FIRST, MIDDLE, LAST)			DRIVER LICENSE NUMBER		CA	DATE OF BIRTH	
INSURANCE COMPANY NAME OF	SURETY COMPANY AT TIME OF ACCIDE	ENT	POLICY NUMBER			-	
COMPANY NAIC NUMBER			POLICY PERIOD				
			FROM	TC			
Describe Vehicle Damage UNK NONE MINOR MOD MAJOR				Shade in Damaged Area			
	MOD MAJOR			Z			

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"

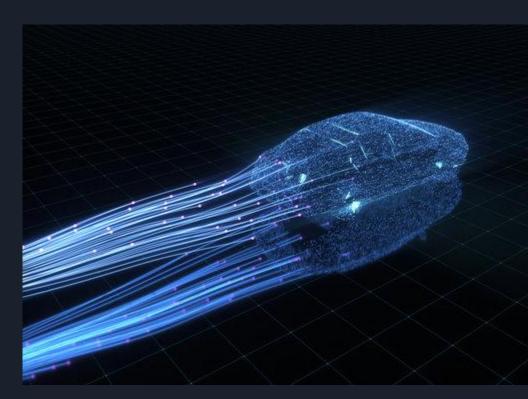
read in all the data once.

Self-Driving Trends in California

Maybe total mileage instead of annual

How many US-20 trips (3300m)





Comparisons twixt companies

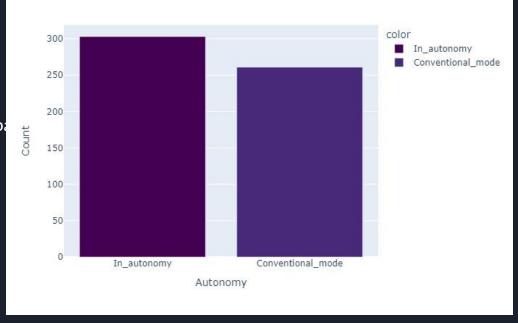


Things break c

What happened?

What were they doing when b

Collision Records in Autonomy



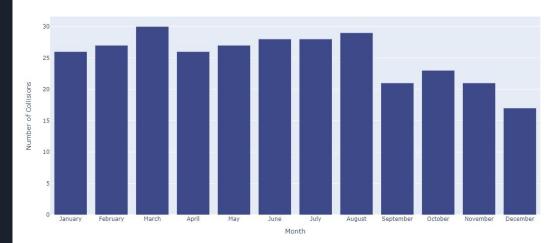
<u>Autonomous collisions:</u>

~25/ month, ~300/year

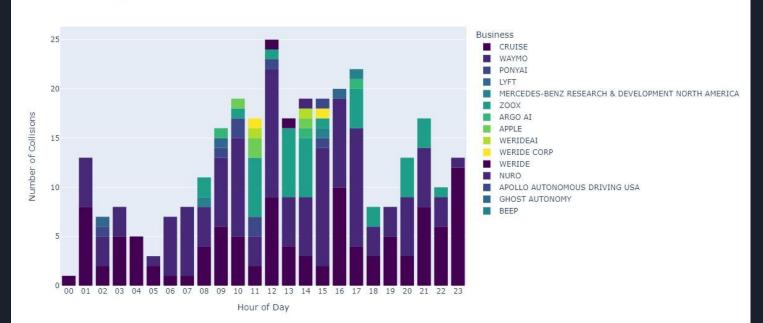
2023 autonomous mileage:

5.75M mi

Autonomous Collisions by Month



Autonomy Collisions by Time of Day

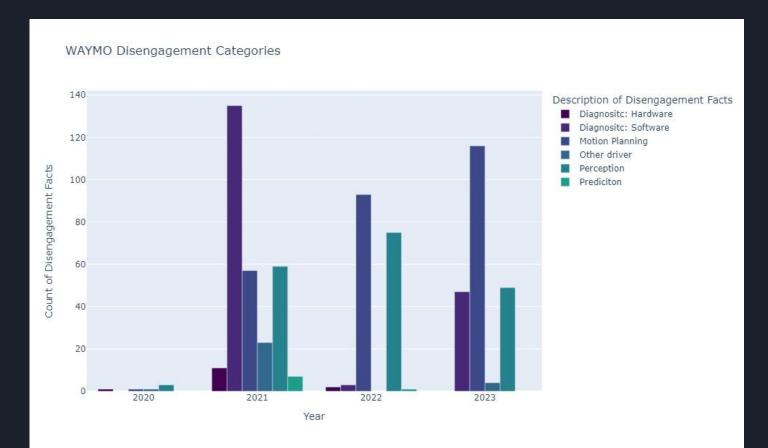


Collision heatmap

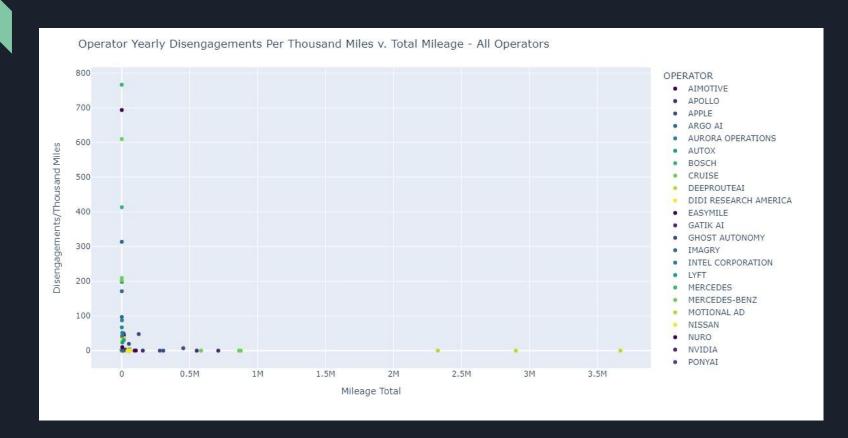
Waymo Gets there

NHTSA - 2021 crashes per 100M miles: 195

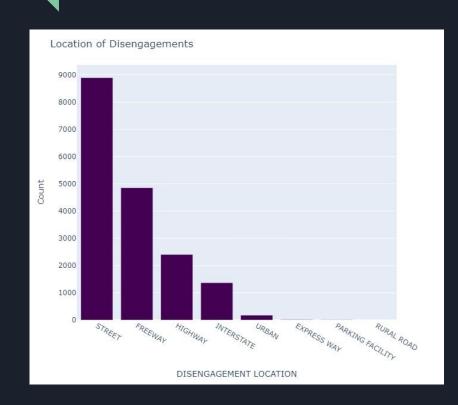
Waymo - disengagement rates

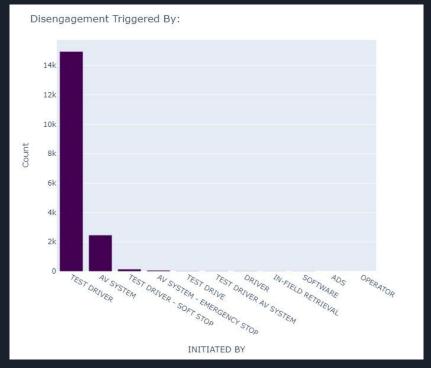


Disengagement Rates vs Miles Logged



Disengagement Insights





Disengagement Insights



Waymo - collision rates

- Compare to NHTSA avg for urban



California

- VINs CA not checking
 - o Dupes
 - Errant (not 17 digits)
- Outliers:
 - Cruise lack of disengagement reports
- Allowed Waymo to get there
- Other states which allow are not publishing this data

Zippy & Pokey

Focus on granularity of data





Conclusions

Things we could do

Waymo - vs all

General break down

Directions we could go in

Other states

Other environments (non-urban)