PROJECT PHASE 1

AVIATION ACCIDENT DATA INTEGRATION



Information Integration and Analytic Data Processing

17-03-2025

Group 3:

- Tommaso Tragno [fc64699]
- Manuel Cardoso [fc56274]
- Chen Cheng [fc64872]
- Cristian Tedesco [fc65149]

Agenda

- Project Motivation & Questions
- Data Sources & Rationale
- Data Profiling & Cleaning
- Integrated Schema
- Identity Resolution & Blocking
- Preliminary Insights & Challenges
- Conclusion & Next Steps

1. Project Motivation & Questions



Explore links between accidents, passenger volume, weather, and aircraft models.

Assess if higher traffic numbers lead to more accidents.

Examine if bad weather increases the likelihood of accidents.

Identify if certain regions or airports have unique accident patterns.

Determine if the era or model of aircraft is linked to higher accident rates.

Goal of this Project

Provide actionable insights for aviation stakeholders—airlines, airports, and regulators—seeking data-driven strategies to enhance safety protocols and resource allocation under conditions of fluctuating demand and variable weather.

2. Data Sources& Rationale

2.1 NTSB Aviation Accident Database (Filtered)

 Contains detailed records of accidents, including location, date/time, flight phase, and possible contributing factor

2.2 Open-Meteo API

• Returns hourly or daily weather data (temperature, precipitation, wind speed, wind direction, etc.) for specified coordinates and date

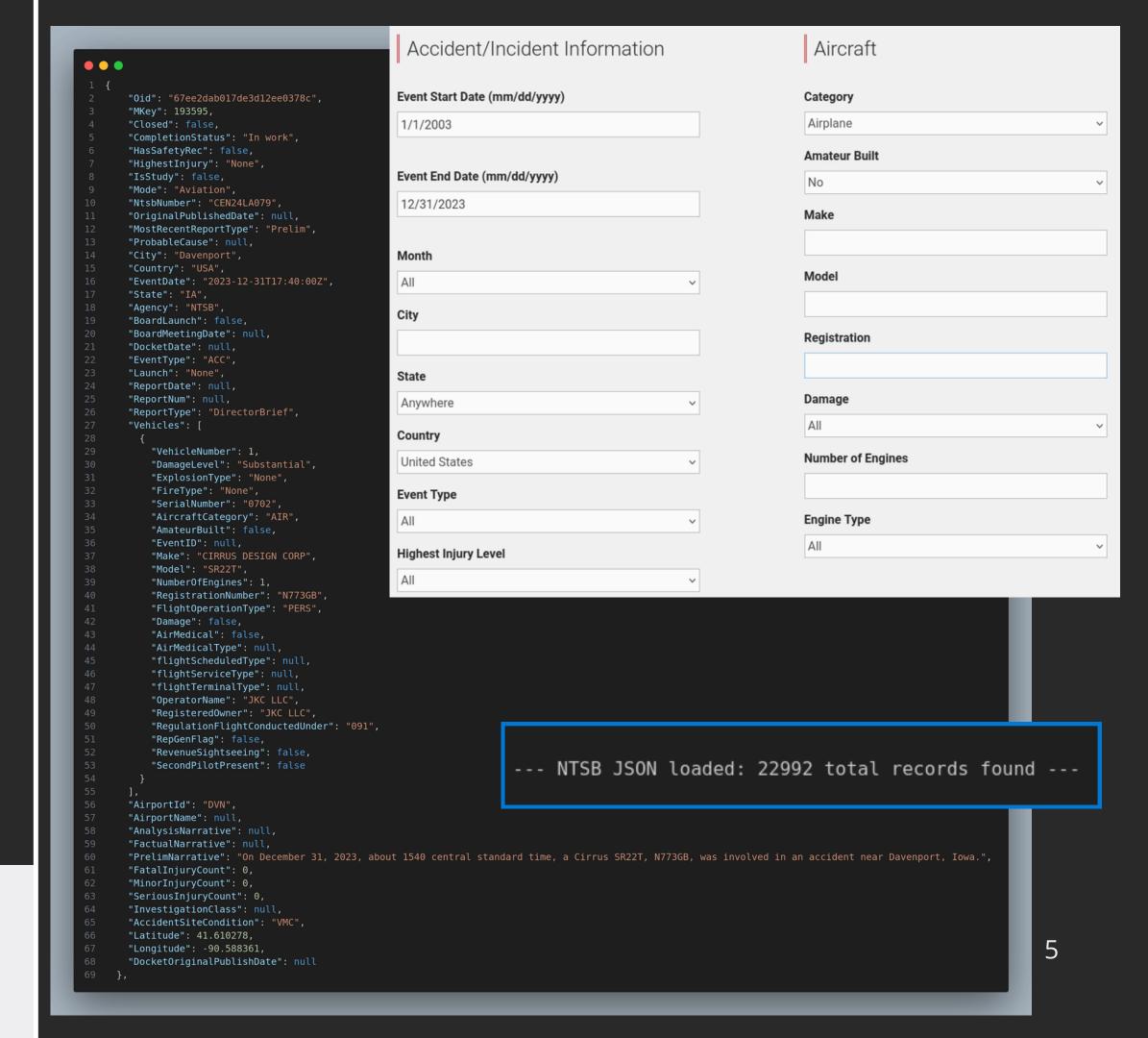
2.3 U.S. Airline Traffic Data (Kaggle)

• Provides monthly flight volumes, passenger counts, and other tourismrelated indicators for both domestic and international flights

2.4 Aircraft Production Data (Kaggle)

Contains details on various aircraft manufacturers, models, and production volume

2.1 NTSB Aviation Accident Database



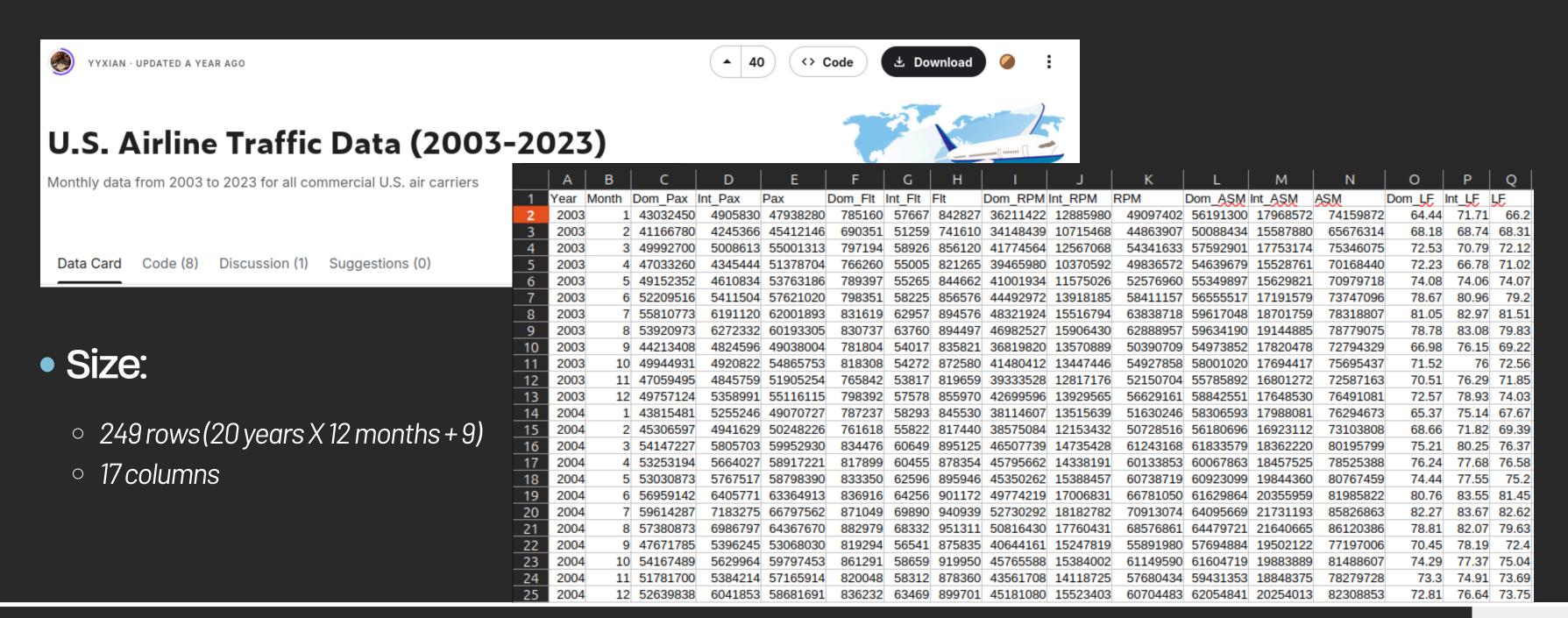
2.2 Open-Meteo API

--- Weather DataFrame sample ---<class 'pandas.core.frame.DataFrame'> RangeIndex: 491592 entries, 0 to 491591 Data columns (total 21 columns):

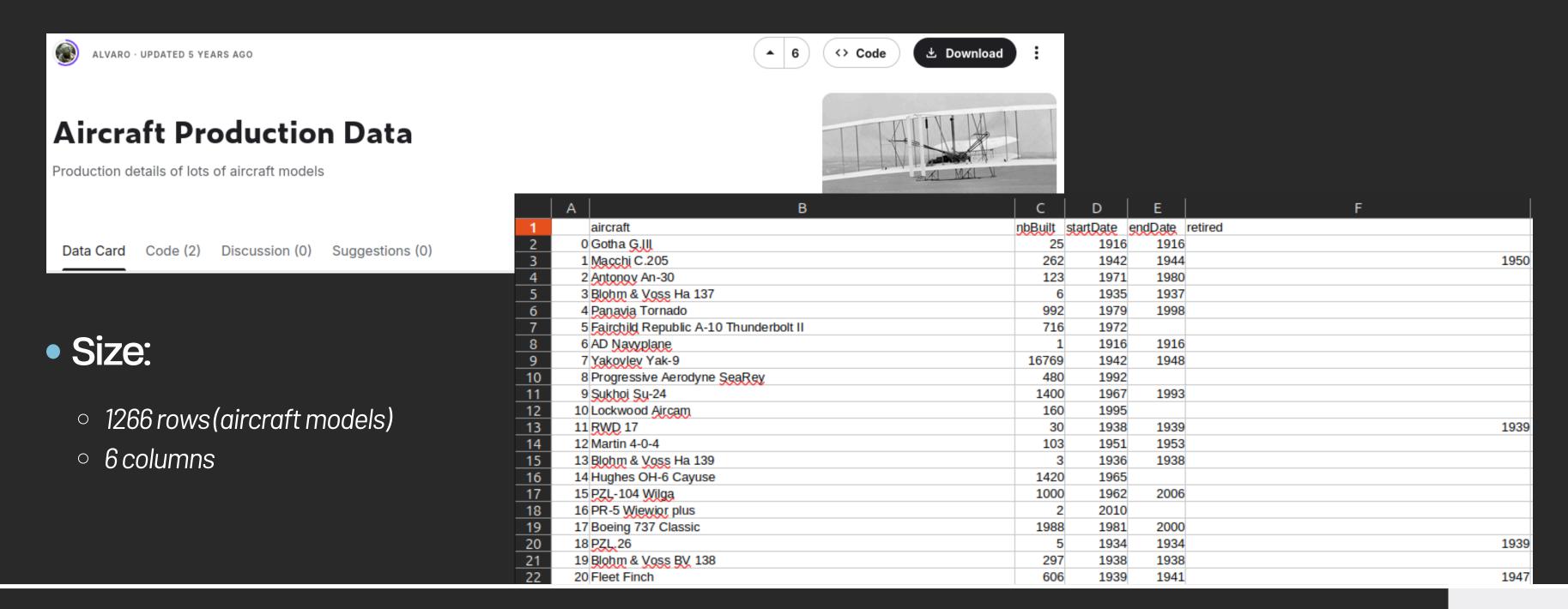
```
def fetch weather data(lat, lon, date str, cache key)
       endpoint = "https://archive-api.open-meteo.com/v1/archive"
       params = {
            "longitude": lon,
            "start date": date str,
            "end date": date str.
            "hourly": ",".join([
                "temperature_2m", "relative_humidity_2m", "dew_point_2m", "pressure_msl",
                "wind speed 100m", "wind direction 10m", "wind direction 100m",
                "wind gusts 10m", "weather code", "snow depth"
       timeout attempts = 0
        while timeout attempts < MAX RETRIES TIMEOUT:</pre>
           status429 attempts = 0
                    response = requests.get(endpoint, params=params, timeout=REQUEST_TIMEOUT)
                    # Check HTTP status
                    if response.status code == 200:
                       return response.json().get("hourly", {})
                    elif response.status code == 429:
                       status429 attempts += 1
                       if status429 attempts <= RETRIES ON 429:</pre>
                           print(f"[{cache_key}] Got 429. Waiting {SLEEP_ON_429_SECS}s then retrying (attempt {status429_attempts}/{RETRIES_ON_429}).")
                           time.sleep(SLEEP_ON_429_SECS)
                           print(f"[{cache_key}] Too many 429 responses. Giving up this accident.")
                       print(f"[{cache_key}] Request failed: {response.status_code}. Skipping.")
                except requests.exceptions.Timeout:
```

```
1 "cen24la079 2023-12-31 41.610278 -90.588361": {"time": ["2023-12-31T00:00", "2023-12-31T01:00", "2023-12-31T02:00", "2023-12-31T03:00", "2023-12-31T04:00", "2023-12-31T05:00", "2023-12
 2 -31T06:00", "2023-12-31T07:00", "2023-12-31T08:00", "2023-12-31T09:00", "2023-12-31T10:00", "2023-12-31T11:00", "2023-12-31T12:00", "2023-12-31T13:00", "2023-12-31T14:00", "2023-12-31T1
 3 5:00", "2023-12-31T16:00", "2023-12-31T17:00", "2023-12-31T18:00", "2023-12-31T19:00", "2023-12-31T20:00", "2023-12-31T21:00", "2023-12-31T22:00", "2023-12-31T23:00"], "temperature_2m":
 4 [1.1, -0.6, -1.0, -0.4, -1.3, -1.8, -2.2, -2.7, -2.7, -2.7, -2.9, -3.2, -3.1, -2.9, -2.6, -2.4, -1.6, -0.7, -0.1, 0.6, 1.2, 1.0, 1.1, 0.8, 0.1], "relative humidity 2m": [85, 89, 89, 80, 71,
 5 69, 70, 74, 73, 71, 72, 76, 78, 77, 77, 74, 70, 73, 74, 73, 72, 68, 67, 70], "dew_point_2m": [-1.1, -2.2, -2.6, -3.5, -5.8, -6.7, -6.8, -6.7, -6.8, -7.4, -7.5, -6.8, -6.2, -6.1, -5.8,
 6 5.7, -5.6, -4.4, -3.4, -3.1, -3.5, -4.2, -4.5, -4.8], "pressure msl": [1011.4, 1012.5, 1012.8, 1013.3, 1013.5, 1013.9, 1014.2, 1014.0, 1014.1, 1014.6, 1014.6, 1014.8, 1015.1, 1015.8, 10
     16.1, 1016.6, 1017.3, 1017.7, 1017.8, 1018.2, 1018.8, 1019.7, 1020.8, 1021.8], "surface pressure": [983.2, 984.0, 984.3, 984.8, 984.9, 985.3, 985.5, 985.3, 985.4, 985.9, 985.8, 986.0, 9
 13 0, 0, 0, 0, 0, 0, 50, 99, 18, 0, 0, 0, 0, 0, 0], "wind speed 10m": [12.0, 11.2, 13.2, 16.6, 17.4, 16.7, 15.8, 16.7, 17.8, 18.1, 18.5, 17.7, 17.4, 16.9, 15.9, 16.8, 18.4, 20.4, 22.6, 25
14 .1, 21.5, 20.1, 18.1, 18.0], "wind speed 100m": [26.0, 25.5, 27.7, 27.6, 27.5, 26.1, 25.6, 27.9, 28.5, 28.7, 29.3, 27.5, 27.4, 26.7, 24.8, 23.8, 25.5, 28.0, 31.5, 35.4, 30.8, 29.6, 28.6
15 , 30.1], "wind_direction_10m": [261, 272, 292, 319, 310, 303, 297, 284, 288, 293, 294, 294, 292, 294, 295, 305, 310, 315, 326, 331, 329, 327, 325, 323], "wind_direction_100m": [266, 279]
16 , 300, 320, 312, 305, 299, 287, 291, 297, 298, 297, 295, 297, 300, 308, 312, 315, 327, 332, 330, 328, 326, 326], "wind gusts 10m": [19.1, 19.1, 21.2, 26.6, 29.2, 28.8, 28.1, 27.4, 29.9,
18 "snow depth": [0.02, 0.02, 0.02, 0.02, 0.02, 0.02, 0.02, 0.02, 0.02, 0.02, 0.02, 0.02, 0.02, 0.02, 0.02, 0.02, 0.02, 0.02, 0.02, 0.02, 0.02, 0.02, 0.02, 0.02]}
```

2.3 U.S. Airline Traffic Data



2.4 Aircraft Production Data



3. Data Profiling& Cleaning

- 3.1 Table Flattening & Duplicates
- 3.2 Type Conversion
- 3.3 Distribution, Missingness, Uniqueness
- 3.4 Data Issues

3.1 Table Flattening & Duplicates

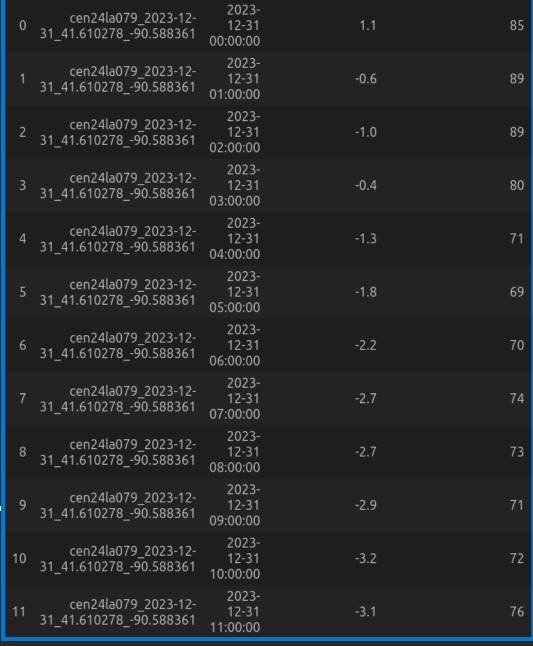
- NTSB Database
- Flattening Nested Information
 - Convert hierarchical data (like JSON) into a flat, tabular structure for analysis.
 - Record will appear multiple times in the resulting DataFrame. The top-level fields get repeated in each row, while any fields from the nested array become columns.
 - May introduce key duplication and data redundancy.

```
"DamageLevel": "None"
"ExplosionType": "None",
"FireType": "None",
"SerialNumber": "C0218",
"AircraftCategory": "AIR"
"AmateurBuilt": false,
"EventID": null,
"Make": "DIAMOND AIRCRAFT IND INC",
"Model": "DA20-C1".
"NumberOfEngines": 1,
"RegistrationNumber": "N857PA".
"FlightOperationType": null,
"Damage": false,
"AirMedical": false,
"AirMedicalType": null,
"flightScheduledType": null,
"flightServiceType": null,
"flightTerminalType": null,
"OperatorName": "DIAMOND AIRCRAFT SALES OF KENTUCKY LLC",
"RegisteredOwner": "DIAMOND AIRCRAFT SALES OF KENTUCKY LLC"
"RegulationFlightConductedUnder": "UNK",
"RevenueSightseeing": false,
"SecondPilotPresent": false
"VehicleNumber": 2,
"DamageLevel": "None",
"ExplosionType": "None",
"FireType": "None",
"SerialNumber": "1955",
"AircraftCategory": "HELI",
"AmateurBuilt": false,
"EventID": null,
"Make": "ROBINSON HELICOPTER",
"Model": "R44",
"NumberOfEngines": 1,
"RegistrationNumber": "N744AF"
"FlightOperationType": null,
"Damage": false,
"AirMedical": false,
"AirMedicalType": null,
"flightScheduledType": null,
"flightServiceType": null,
"flightTerminalType": null,
"OperatorName": "SKYLINE HELICOPTER TOURS LLC",
```

П	Vehicles.VehicleNumber	Vehicles.SerialNumber	Vehicles.Make	Vehicles.Model	NtsbNumber	ProbableCause	City	Country	EventDate	State	Agency	EventType	AirportId	AirportName	Latitude	Longitude
3	9 1	c0218	diamond aircraft ind inc	da20-c1	ops24la011	None	north las vegas	usa	2023-12- 09 13:06:00		ntsb	осс	vgt	north las vegas	36.211268	-115.19968
4	0 2	1955	robinson helicopter	г44	ops24la011	None	north las vegas	usa	2023-12- 09 13:06:00		ntsb	осс	vgt	north las vegas	36.211268	-115.19968

3.1 Table Flattening & Duplicates

			3 5:00",	"2023-12-31T16:00	", "2023-12	?-31T17:00",	"2023-12	2-31T18:00	", "2023-1	12-31T19:00'	", "2023-1	12-31T20:	00", "202	3-12-31T21	:00", "2023	-12-31T22	:00", "2023-	12-31T23:00"],	"temperature	e_2m":
			4 [1.1,	-0.6, -1.0, -0.4,	-1.3, -1.8	3, -2.2, -2.	7, -2.7,	-2.9, -3.2	2, -3.1,	-2.9, -2.6,	-2.4, -1	.6, -0.7,	-0.1, 0.	6, 1.2, 1.	0, 1.1, 0.8	, 0.1], "	relative_hum	idity_2m": [85	, 89, 89, 80,	71,
			5 69 70	74, 73, 71, 72,	76, 78, 77,	77, 74, 76), 73, 74,	, 73, 72, (58, 67, 70	<code>9], "dew_po:</code>	int_2m":	[-1.1, -2]	.2, -2.6,	-3.5, -5.	8, -6.7, -6	.8, -6.7,	-6.9, -7.4,	-7.5, -6.8, -	6.2, -6.1, -5	5.8, -
AccidentID	time	temperature_2m	relative_humidity_2m	6, -4.4, -3.4, -	3.1, -3.5,	-4.2, -4.5,	-4.8], '	pressure_r	msl": [10]	11.4, 1012.	5, 1012.8	, 1013.3,	1013.5,	1013.9, 10	14.2, 1014.	0, 1014.1	, 1014.6, 10	14.6, 1014.8,	1015.1, 1015.	8, 10
	2023-			16.6, 1017.3, 16																
cen24la079_2023-12-	12-31	11	85	7.1, 987.4, 987.	9, 988.7, 9	89.2, 989.3	3, 989.8,	990.3, 99	1.2, 992.3	3, 993.2], '	"precipita	ation": [0.0, 0.0,	0.0, 0.0,	0.0, 0.0,	0.0, 0.0,	0.0, 0.0, 0	.0, 0.0, 0.0,	0.0, 0.0, 0.0	0.0
_41.61027890.588361	00:00:00		03	.1, 0.1, 0.0, 0.	0, 0.0, 0.0)], "rain":	[0.0, 0.6	9, 0.0, 0.0	9, 0.0, 0	.0, 0.0, 0.0	0.0, 0	.0, 0.0,	0.0, 0.0,	0.0, 0.0,	0.0, 0.0,	0.0, 0.0,	0.0, 0.0, 0	.0, 0.0, 0.0],	"snowfall":	[0.0,
				0, 0.0, 0.0, 0.6	, 0.0, 0.0,	0.0, 0.0,	0.0, 0.0,	, 0.0, 0.0	, 0.0, 0.0	9, 0.0, 0.2	1, 0.07, 0	0.07, 0.0	, 0.0, 0.	0, 0.0], "	cloud_cover	": [100,	87, 100, 100	, 100, 100, 10	0, 100, 100,	100,
cen24la079 2023-12-	2023-	0.6	80	, 100, 100, 100,	99, 100, 1	100, 100, 10	00, 100, 1	100, 99, 10	90], "clou	ud_cover_lo	v": [0, 46	6, 100, 1	.00, 100,	100, 100,	100, 100, 1	.00, 100,	100, 100, 10	0, 98, 96, 100	, 100, 97, 10	00, 10
41.61027890.588361	12-31 01:00:00	-0.6	89	8, 100], "cloud	cover_mid":	[88, 0, 0,	0, 0, 0,	, 0, 25, 1	2, 37, 45	, 98, 100, 3	100, 96, 8	85, 100,	100, 100,	100, 100,	100, 42, 6], "cloud	_cover_high"	: [100, 76, 0,	0, 0, 0, 0,	0, 0,
	01.00.00			, 0, 0, 0, 50, 9	9, 18, 0, 6), 0, 0, 0,	0], "wind	d_speed_10	n": [12.0	, 11.2, 13.	2, <mark>16.6,</mark> 3	17.4, 16.	7, 15.8,	16.7, 17.8	, 18.1, 18.	5, 17.7,	17.4, 16.9,	15.9, 16.8, 18	.4, 20.4, 22.	6, 25
con24la079 2023-12-	2023-			, 20.1, 18.1, 18	.0], "wind_	speed_100m'	': [26.0,	25.5, 27.	7, 27.6, 2	27.5, 26.1,	25.6, 27	.9, 28.5,	28.7, 29	.3, 27.5,	27.4, 26.7,	24.8, 23	.8, 25.5, 28	.0, 31.5, 35.4	, 30.8, 29.6,	28.6
cen24la079_2023-12- 41.61027890.588361	12-31	-1.0	89	"wind_direction	10m": [261	, 272, 292,	319, 316	9, 303, 29 ³	7, 284, 28	88, 293, 29 ⁴	4, 294, 29	92, 294,	295, 305,	310, 315,	326, 331,	329, 327,	325, 323],	"wind_direction	n_100m": [266	5, 279
41.010276_50.500501	02:00:00			20, 312, 305, 29	_ 9, 287, 291	, 297, 298,	297, 295	5, 297, 30	9, 308, 3	12, 315, 32 ¹	7, 332, 33	30, 328,	326, 326]	, "wind_gu	sts_10m": [19.1, 19.	1, 21.2, 26.	6, 29.2, 28.8	28.1, 27.4,	29.9,
0.41 0.70 0.000 4.0	2023-			1.0, 31.7, 29.5,	28.8, 27.7	7, 29.2, 32.	8, 33.8,	39.6, 44.3	3, 43.2, 3	37.1, 33.5,	29.5], "v	weather_c	ode": [3,	3, 3, 3,	3, 3, 3, 3,	3, 3, 3,	3, 3, 3, 3,	3, 3, 73, 71,	71, 3, 3, 3,	3],
cen24la079_2023-12-	12-31	-0.4	80	pth": [0.02, 0.6	2, 0.02, 0.	02, 0.02, 6	0.02, 0.02	2, 0.02, 0	.02, 0.02	, 0.02, 0.02	2, 0.02, 0	0.02, 0.0	2, 0.02,	0.02, 0.02	, 0.02, 0.0	2, 0.02,	0.02, 0.02,	0.02]}		



Open-Meteo Database

"cen24la079_2023-12-31_41.610278_-90.588361": {"time": ["2023-12-31T00:00", "2023-12-31T01:00", "2023-12-31T02:00", "2023-12-31T03:00", "2023-12-31T04:00", "2023-12-31T05:00", "2023-12-31T06:00", "2023-12-31T07:00", "2023-12-31T13:00", "2023-12-31T14:00", "2023-12-31T12:00", "2023-12-31T12:00", "2023-12-31T13:00", "2023-12-31T14:00", "2023-12-31T14:00", "2023-12-31T12:00", "2023-12-31T13:00", "2023-12-31T14:00", "2023-12-3

3.2 Type Conversion

Ensuring Consistent Data Formats

- Convert columns to appropriate data types (e.g., strings to integers, dates to datetime).
- Standardized dates using ISO 8601 format (YYYY-MM-DD) for consistency and compatibility.
- Enables accurate sorting, filtering, and merging across datasets.

```
# Type Conversion

df_ntsb['EventDate'] = pd.to_datetime(df_ntsb['EventDate']).dt.tz_localize(None)

df_ntsb['Vehicles.VehicleNumber'] = pd.to_numeric(df_ntsb['Vehicles.VehicleNumber'], errors='coerce').astype(int)

df_ntsb['MKey'] = pd.to_numeric(df_ntsb['MKey'], errors='coerce').astype(int)

df_ntsb['Vehicles.NumberOfEngines'] = pd.to_numeric(df_ntsb['Vehicles.NumberOfEngines'], errors='coerce').fillna(0).astype(int)

df_ntsb['Latitude'] = pd.to_numeric(df_ntsb['Latitude'], errors='coerce').astype(float)

df_ntsb['Longitude'] = pd.to_numeric(df_ntsb['Longitude'], errors='coerce').astype(float)

df_ntsb['TotalInjuryCount'] = pd.to_numeric(df_ntsb['TotalInjuryCount'], errors='coerce').astype(int)
```

3.3 Distribution, Missingness, Uniqueness

Data Profiling – Understanding the Dataset

- Distribution: Understand how data is spread look at min, max, mean, standard deviation, and quartiles to detect skewness or outliers.
- Missingness: Identify how many values are missing per column (count and %) to assess data quality and plan cleaning or imputation.
- Uniqueness: Measure cardinality (number of unique values) to spot identifiers, repeated patterns, or categorical vs.
 continuous variables.
- Mode & Frequency: Reveal dominant values, helpful for detecting defaults or common entries.

Column	DataType	TotalCount	NonNullCount	NumMissing	MissingPerc	Cardinality	Mode	ModeFreq	Mean	Min	Q25	Q50	Q75	Max	Std

3.4 Data Issues

Invalid Values in Aircraft Models Data Columns

- Some entries in startDate contain values unrelated to actual years.
- Affects temporal analysis and integration with time-based datasets.
- Manually corrected a few anomalies using reliable sources from the web.

	Column	DataType	Min	Q25	Q50	Q75	Max	Std
0	aircraft	object	NaN	NaN	NaN	NaN	NaN	NaN
1	nbBuilt	int64	0.0	32.25	185.0	703.00	43400.0	3618.899938
2	startDate	int64	1.0	1937.00	1951.0	1974.75	2015.0	224.918816
3	endDate	Int64	1.0	1938.00	1949.0	1979.00	2016.0	227.826756

```
df_filtered = df_aircraft[(df_aircraft['startDate'] < 1000) | (df_aircraft['endDate'] < 1000)]
df_filtered.style.map(
    lambda val: 'background-color: red' if val < 1000 else '',
    subset=['startDate', 'endDate']
)</pre>
```

	aircraft	nbBuilt	startDate	endDate
82	lockheed c-5 galaxy	131	5	5
86	british aerospace nimrod aew3	8	11	11
171	schneider es-57 kingfisher	11	2	
190	bell 222	230	222	1991
284	flitfire	49	10	10
308	grumman c-2 greyhound	58	2	2
498	chu hummingbird	2	2	2
514	embraer legacy 500	500	500	
518	lockheed martin f-22 raptor	195	22	22
536	gallaudet d-4	2	2	2
637	fleet canuck	225	198	198
668	bell ah-1 supercobra	1271	1	1
688	chu cjc-3	1	1	1
896	dallach sunrise	0	5	5
951	sukhoi su-30mki	200	30	30
1049	boeing kb-29 superfortress	282	92	
1089	myasishchev m-4	2	93	93
1200	yakovlev yak-100	2	115	115

3.4 Data Issues

- Removing Formatting for Type
 Conversion in Airline Traffic Data
 - Numbers stored as strings with commas (e.g., "1,200") can't be converted to integers.
 - Strip commas before converting to numeric types.
 - Ensures correct data types for computation and analysis.

```
1 Year, Month, Dom Pax, Int Pax, Pax, Dom Flt,
2 2003,1 "43,032,450", "4,905,830", "47,938,280", "785,160",
3 2003,2 "41,166,780", "4,245,366", "45,412,146", "690,351",
4 2003,3 "49,992,700", "5,008,613", "55,001,313", "797,194",
```

```
# Remove commas from all columns and then convert
df_airline_traffic = df_airline_traffic.replace(',', '', regex=True)

# Now convert each column to numeric. If everything converts well, no rows become NaN.
df_airline_traffic = df_airline_traffic.apply(pd.to_numeric, errors='coerce').astype(int)
```

4. Integrated Schema

4.1 Data Selection

- Focused on relevant attributes to answer the research question
- Removed unnecessary data for clarity and efficiency

4.2 Data Integration

- Combined multiple tables into a single unified schema
- Maintained data consistency, integrity, and efficiency

4.3 Schema Design

- Used high-cardinality relationships for optimal data linkage
- Supported different types of correspondences (e.g., many-to-one)

4.4 Binding Methods

Applied techniques discussed in "Identity Resolution"

4. Integrated Schema

Origin Dataset	From	Target	Type Corresp.	Description
NTSB	NtsbNumber	AccidentNumber	1-1	
NTSB	EventDate	DateTime	1-1	
NTSB	City	City	1-1	
NTSB	State	State	1-1	
NTSB	Longitude	Longitude	1-1	
NTSB	Latitude	Latitude	1-1	
NTSB	AirportName	AirportName	1-1	
NTSB	Operator	Operator	1-1	
NTSB	Aircraft Damage	Aircraft Damage	1-1	
NTSB	FatalInjuryCount; SeriousInjuryCount; Mi- norInjuryCount	TotalInjuryCount	N-1	Sum all the value
NTSB	HighestInjury	HighestInjury	1-1	
NTSB	Model, Make	Aircraft	N-1	Blocking with Aircraft Data
Aircraft Data	StartDate	ProductionStartDate	1-1	
Aircraft Data	EndDate	ProductionEndDate	1-1	
Weather API	Temperature_2m	Temperature	1-1	
Weather API	Precipitation	Precipitation	1-1	
Weather API	Wind_Speed_10m	WindSpeed	1-1	
Weather API	Weather code	Weather code	1-1	
Weather API	other weather info	other weather info	1-1	
Airline Traffic	Pax	PassengersPerMonth	1-1	
Airline Traffic	Flt	FlightsPerMonth	1-1	
Airline Traffic	LF	LoadFactorPerMonth	1-1	

INTEGRATED MODEL

5. Identity Resolution & Blocking

Blocking with Q-grams (q=3) and Substring Matching

Conditional Numeric Filter

Combined Textual Similarity (Threshold 0.75)

sta	artDate	Vehicles.Model	Matched_Aircraftl	JW	LEV	JAC Fi	nalScore
2.	1965	concorde	concorde	1.0	1.0	1.0	1.0
3.	1938	quad city challenger	quad city challenger	1.0	1.0	1.0	1.0
8.	1967	bac 167 strikemaster	bac strikemaster	0.96	0.8	0.66	0.82
9.	1962	pzl 104m wilga	pzl 104 wilga	0.98	0.92	0.5	0.82
33.	1985	gulfstream	gulfstream	0.95	0.76	0.5	0.76

Final Score = $0.4 \times Jaro-Winkler + 0.3 \times Levenshtein + 0.3 \times Jaccard$

6. Preliminary Insights & Challenges

Passenger Volume Limitation:

Passenger data is monthly, while accidents occur daily—this 1-to-many cardinality reduces correlation accuracy

Aircraft Data Limitation:

The aircraft dataset lacks detailed specs restricting deeper analysis.

Low Match Rate:

From ~23,000 accident records, only 38 aircraft matched due to identity resolution limits

7. Conclusion

Conclusion

- We performed data cleaning & schema integration for four data sources
- Implemented identity resolution to match aircraft models
- o Preliminary results show plausible correlations but more analysis needed

8. References

- National Transportation Safety Board (NTSB). url: https://www.ntsb.gov/
- Yyxian. U.S. Airline Traffic Data. Kaggle Dataset. url: https://www.kaggle.com/datasets/yyxian/u-s-airline-traffic-data/
- Alvaroibrain. Aircraft Production Data. Kaggle Dataset. url: https://www.kaggle.com/datasets/alvaroibrain/aircraft-production-data
- Open-Meteo Historical Weather API Documentation. url: https://open-meteo.com/en/docs/historical-weather-api?

Thanks for the attention!