# PROBLEM STATEMENT

Raw data contains issues:
- "N/A" values instead of numbers
- Missing values
- Different data formats

Business needs:
- Clean data
- Automated pipeline
- Easy visualization for insights

DATA PIPELINE ARCHITECTURE

CSV / Raw Dataset → Snowfiake (staging layer) → dbt Transformations → Orchestration (dbt run + test) → Visualization (dashbcards & reports)

# 🖌️ DATA CLEANING

**01** Removed incomplete rows
- Dropped records with missing values in critical fields (Name, Platform, Year, Genre, Publisher, Sales).

**02** Enforced correct data types
- Converted Rank and Year into integers.
- Converted Sales values into floats for accurate calculations.
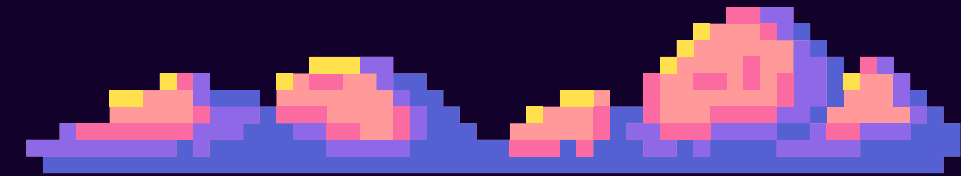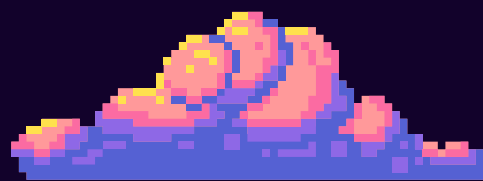
**03** Handled missing values
- Replaced any remaining null sales values with 0.0 using COALESCE.

**04** Standardized column names
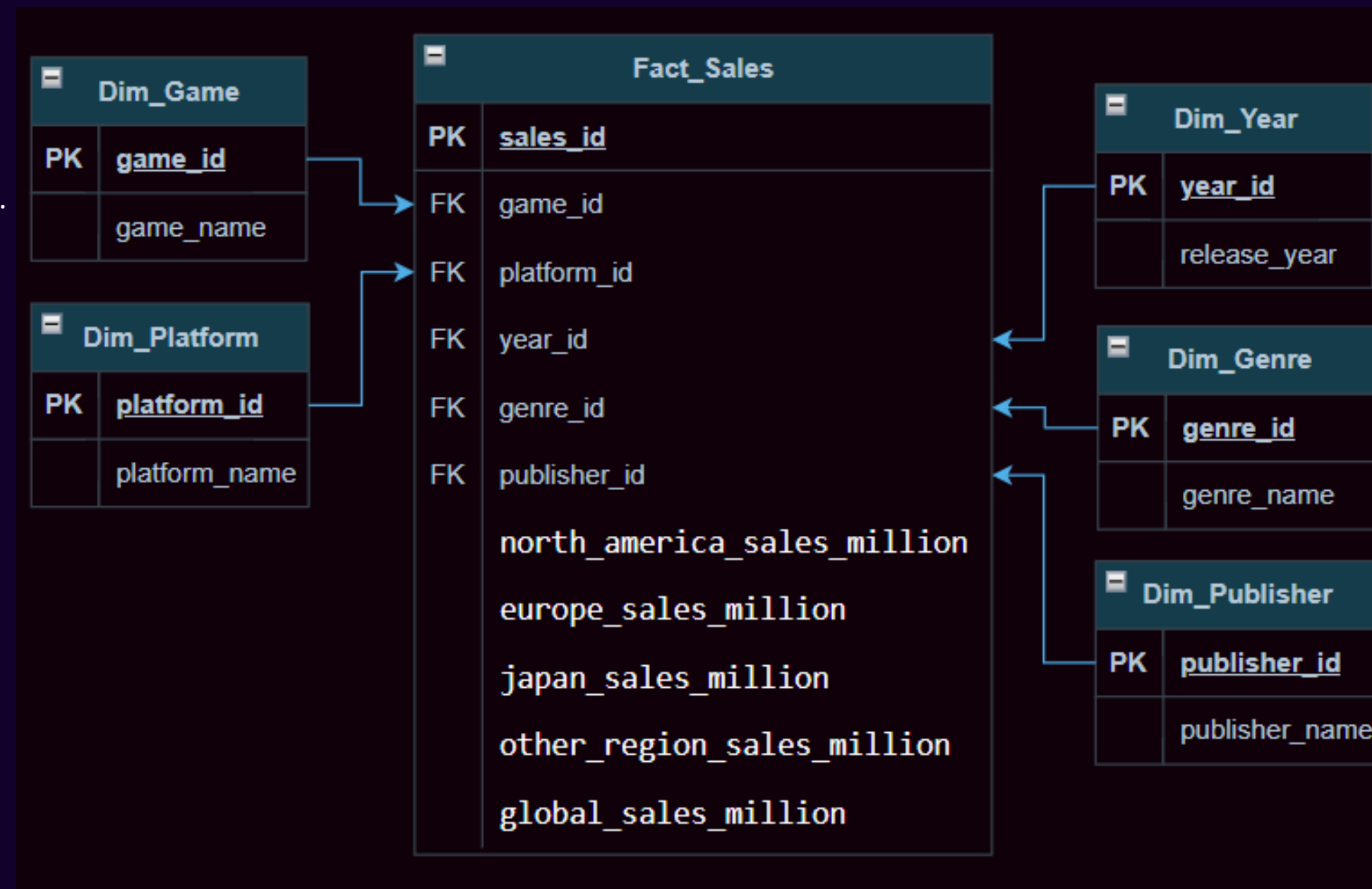- Renamed raw columns into meaningful business terms (e.g., na_sales → North America Sales (M) ).

# DATA MODELING

## ⭐ WHY STAR SCHEMA?

**01** Simplicity
- Easy to understand: one Fact table (sales) linked to multiple Dimension tables (games, publishers, platforms, time).

**02** Query Performance
- Optimized for analytical queries in Snowflake & Power BI (fewer joins, faster aggregations).

**03** Flexibility in Analysis
- Allows slicing and dicing data by Year, Platform, Genre, Publisher, etc.

**04** Best Practice in Data Warehousing
- Star schema is the industry standard for BI tools like Power BI, Tableau.

| Dim_Game | |
|---|---|
| PK | **game_id** |
| | game_name |

| Dim_Platform | |
|---|---|
| PK | **platform_id** |
| | platform_name |

| Fact_Sales | |
|---|---|
| PK | **sales_id** |
| FK | game_id |
| FK | platform_id |
| FK | year_id |
| FK | genre_id |
| FK | publisher_id |
| | north_america_sales_million |
| | europe_sales_million |
| | japan_sales_million |
| | other_region_sales_million |
| | global_sales_million |

| Dim_Year | |
|---|---|
| PK | **year_id** |
| | release_year |

| Dim_Genre | |
|---|---|
| PK | **genre_id** |
| | genre_name |

| Dim_Publisher | |
|---|---|
| PK | **publisher_id** |
| | publisher_name |

# AUTOMATION

## 01 DAG Definition
- Name: dbt_workflow
- Runs daily (@daily).

```python
# Define the DAG
dag = DAG(
    'dbt_workflow',
    default_args=default_args,
    description='Run dbt models using dbt Core',
    schedule_interval='@daily',
    catchup=False,
)
```
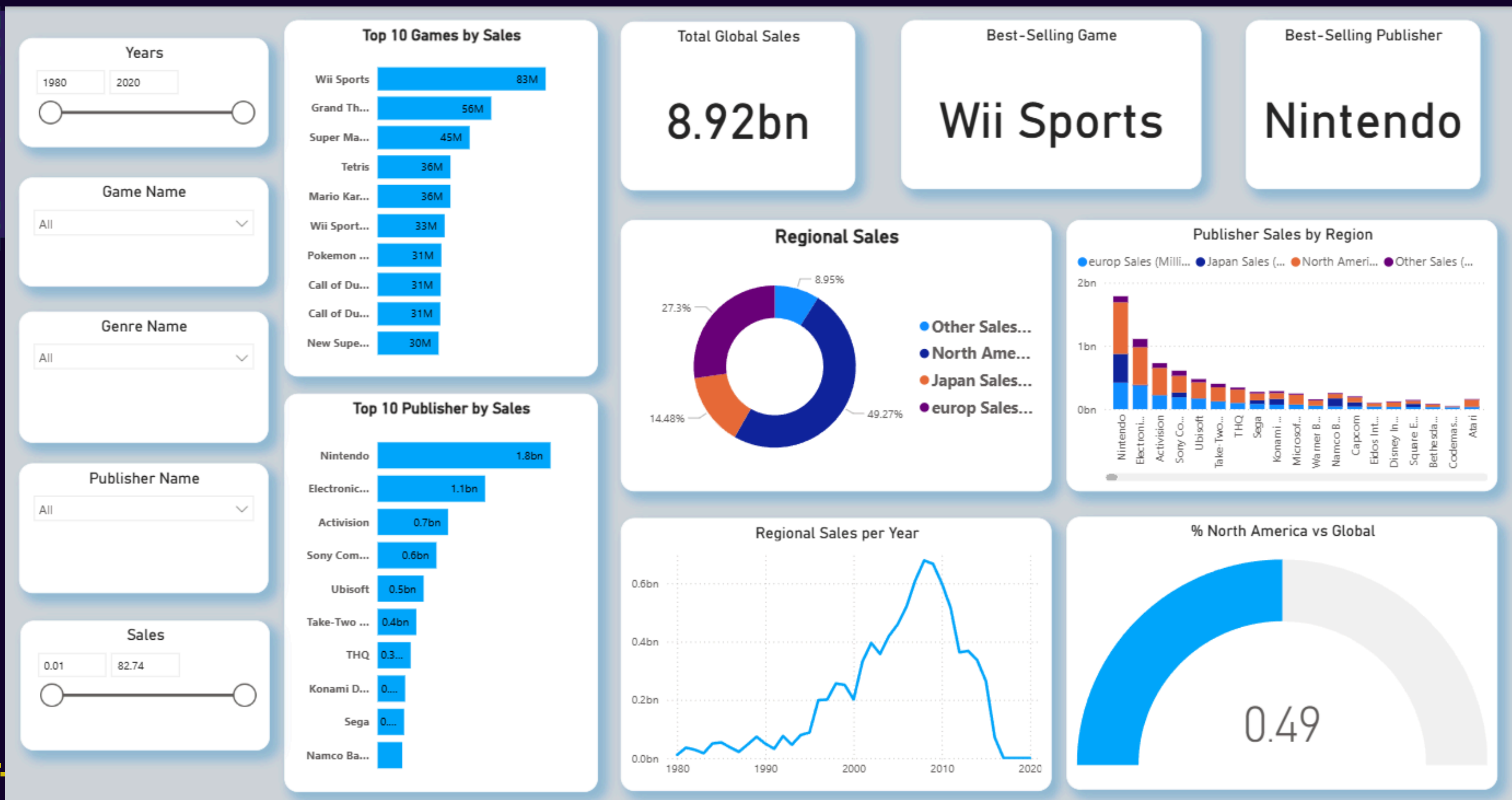
## 02 Task: Run dbt Models
- Executes dbt run command inside the project directory.
- Builds and refreshes data models in Snowflake.

```python
# Task 1: Run dbt models
dbt_run = BashOperator(
    task_id='dbt_run',
    bash_command=f'cd {DBT_PROJECT_DIR} && dbt run',
    dag=dag,
)

# Define task dependencies
dbt_run
```
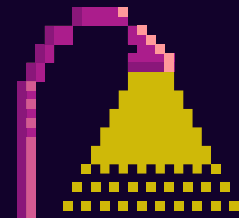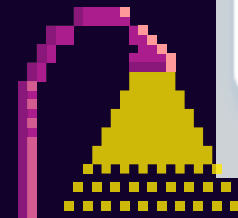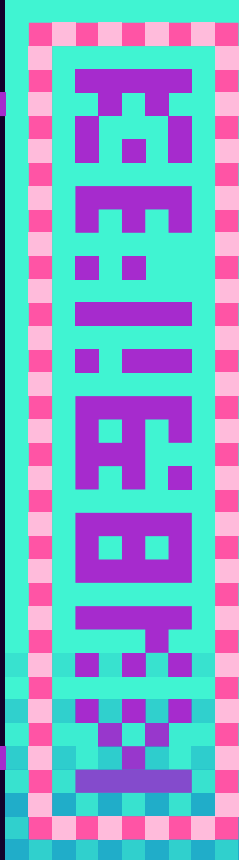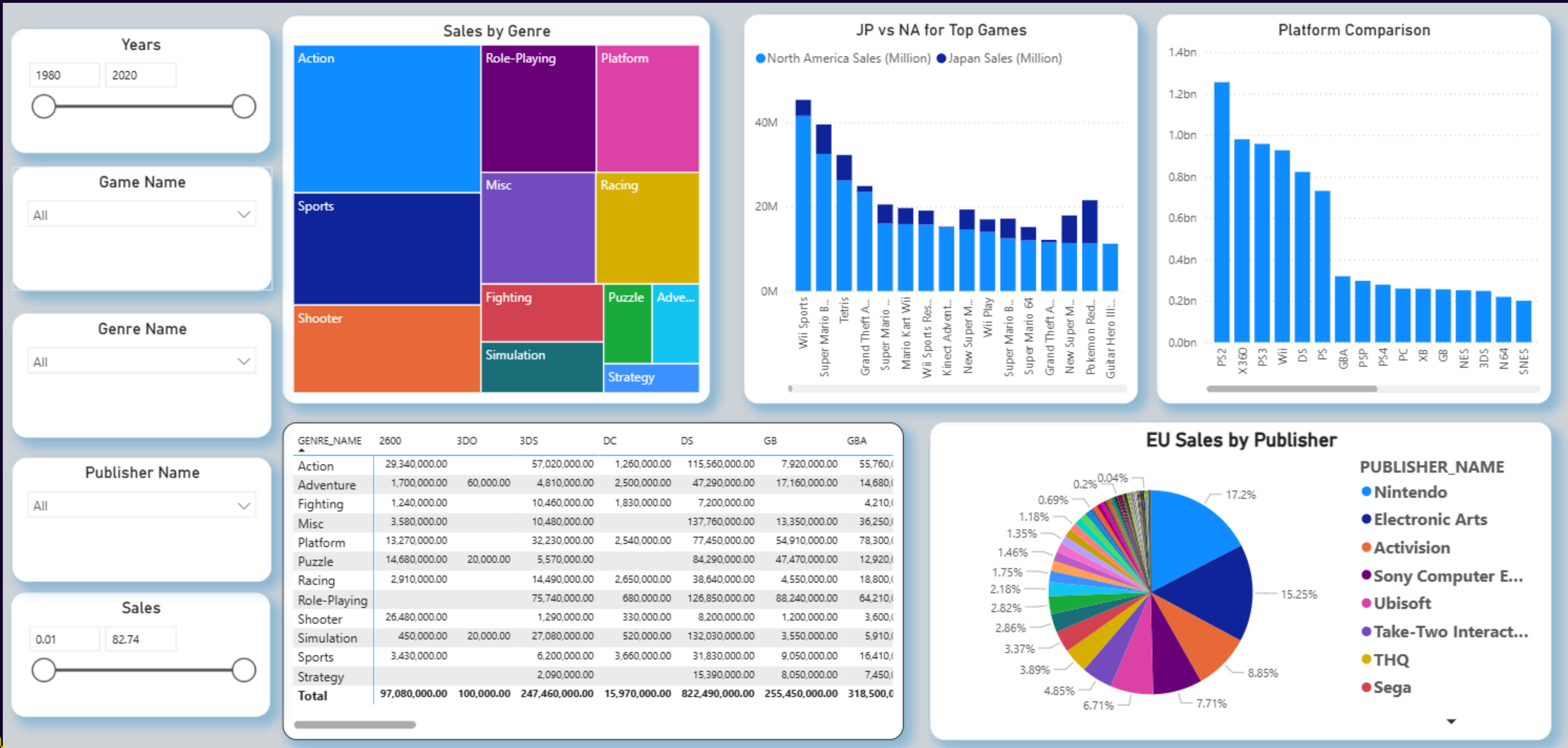
# DASHBOARD (POWER BI)

## Years
1980 — 2020

## Game Name
All

## Genre Name
All

## Publisher Name
All

## Sales
0.01 — 82.74

### Top 10 Games by Sales
| Game | Sales |
|------|-------|
| Wii Sports | 83M |
| Grand Th... | 56M |
| Super Ma... | 45M |
| Tetris | 36M |
| Mario Kar... | 36M |
| Wii Sport... | 33M |
| Pokemon ... | 31M |
| Call of Du... | 31M |
| Call of Du... | 31M |
| New Supe... | 30M |

### Top 10 Publisher by Sales
| Publisher | Sales |
|-----------|-------|
| Nintendo | 1.8bn |
| Electronic... | 1.1bn |
| Activision | 0.7bn |
| Sony Com... | 0.6bn |
| Ubisoft | 0.5bn |
| Take-Two ... | 0.4bn |
| THQ | 0.3... |
| Konami D... | 0.... |
| Sega | 0.... |
| Namco Ba... | |

### Total Global Sales
8.92bn

### Best-Selling Game
Wii Sports

### Best-Selling Publisher
Nintendo

### Regional Sales
- 8.95%
- 27.3%
- 14.48%
- 49.27%
- Other Sales...
- North Ame...
- Japan Sales...
- europ Sales...

### Publisher Sales by Region
- europ Sales (Milli... • Japan Sales (... • North Ameri... • Other Sales (...

2bn
1bn
0bn

Nintendo, Electroni..., Activision, Sony Co..., Ubisoft, Take-Two..., THQ, Sega, Konami..., Microsof..., Warner B..., Namco B..., Capcom, Eidos Int..., Disney In..., Square E..., Bethesda..., Codemas..., Atari

### Regional Sales per Year
0.6bn
0.4bn
0.2bn
0.0bn

1980  1990  2000  2010  2020

### % North America vs Global
0.49

# CHALLENGES & SOLUTIONS

**01** Data quality issues → Solved by cleaning (null removal, type casting, handling N/A).

**02** Large raw data → Automated with Airflow + dbt.

**03** Business understanding → Visualized in Power BI.

# CONCLUSION

**01** Built an end-to-end data pipeline: raw → cleaned → Snowflake → Power BI.

**02** Automated pipeline with Airflow.

**03** Designed Star Schema for analytics.

THANK YOU