# Azure Databricks

https://z.umn.edu/databricks_asr

# Getting Data

## Query ➡

- PeopleSoft
- uAchieve
- DORS
- Etc.

## Transform ➡

- Filter
- Flatten
- Convert
- Etc.

## Save

- Oracle
- Tableau
- BI
- Etc.

# Current Approaches

1. Application Development Team

2. OIT/EDMR

# Downsides

- Changes need to go back through App Dev or OIT

- The solution is maintained by people unfamiliar with the data

- Solutions are usually single use

# Example

## CSDS



Tuesday/Thursday demand

Colleges are permitted to schedule a maximum of 50% of their class hours on Tuesday and/or Thursday. ?

**College classes on T/Th: 46.3%** | **Department classes on T/Th: 47.8%**

Time period demand

Colleges are permitted to schedule up to 3% of departmental classes during any individual time period on any given weekday. ?

**Hours allowed per time period for selected departments: 0.6** ?

| Mpls | St. Paul | Mon | Tues | Weds | Thurs | Fri |
|---|---|---|---|---|---|---|
| 8:00-9:05 | 8:30-9:35 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 9:05-10:10 | 9:35-10:40 | 0.0 | 0.2 | 0.0 | 0.0 | 0.8 |
| 10:10-11:15 | 10:40-11:45 | *0.8* | *1.1* | *0.8* | *0.0* | 0.8 |
| 11:15-12:20 | 11:45-12:50 | *0.8* | *1.9* | *0.8* | *0.8* | 0.8 |
| 12:20-1:25 | 12:50-1:55 | 0.8 | 0.2 | 0.8 | 0.8 | 0.8 |

https://z.umn.edu/databricks_asr

# Example

## CSDS

- Changes to data must be done by App Dev

- App Dev are not experts in Class Scheduling data

- No one else can use this data

# Solutions

- Works with a variety of data

- Lets people manage data they know and need

- Cloud based

# Data Lake(houses)

- Lots of options in this space

- Modern way of managing and analyzing data

- Similar to a Data Warehouse, but more flexible

https://z.umn.edu/databricks_asr

# Databricks

- Cloud hosted data lakehouse

- Well integrated with Azure, OIT's preferred cloud

# Databricks

- Store data from many sources

- Send data to many targets

- Easy transformation

- Lots of automation options

- Analytics

# CSDS in Databricks

https://z.umn.edu/databricks_asr

# Goals

- Easy to load data in from a variety of sources

- Data and Business Analysts can manage data they use

- Resulting data can be used by multiple systems and people

https://z.umn.edu/databricks_asr

# Query

↻ Refresh

Description:  Created by the file upload UI
Created at:  2022-11-16 17:39:19
Last modified:  2022-11-16 17:39:35
Partition columns:
Number of files:  1
Size:  7.72 kB

## Schema:

|   | col_name | data_type | comment |
|---|----------|-----------|---------|
| 1 | INSTITUTION | string | null |
| 2 | ACAD_GROUP | string | null |
| 3 | EFFDT | string | null |
| 4 | EFF_STATUS | string | null |
| 5 | DESCR | string | null |
| 6 | DESCR100 | string | null |

## Sample Data:

|   | INSTITUTION | ACAD_GROUP | EFFDT | EFF_STATUS | DESCR | DESCR100 |
|---|-------------|------------|-------|------------|-------|----------|
| 1 | UMNTC | TALA | 01-JAN-00 | A | Arch & Landscape Arch, Coll of | College of Architecture and Landscape Architecture |
| 2 | UMNTC | TBEL | 01-JAN-00 | A | Bell Museum | Bell Museum |
| 3 | UMNTC | TCBS | 01-JAN-00 | A | Biological Sciences, Coll of | College of Biological Sciences |

# Transformation Process



```sql
1   select c.institution,
2          ffr.strm,
3          c.acad_group,
4          c.acad_org,
5          ffr.period_id,
6          ffr.day_of_week,
7          sum(ffr.seconds_used) aws
8   from ps_class_mtg_pat_fact_filtered_rollup ffr
9   left join cs_ps_class_tbl c
10   on ffr.crse_id=c.crse_id
11   and ffr.crse_offer_nbr=c.crse_offer_nbr
12   and ffr.strm=c.strm
13   and ffr.session_code=c.session_code
14   and ffr.class_section=c.class_section
15  where ffr.strm = '1229'
16  group by c.institution,
17          ffr.strm,
18          c.acad_group,
19          c.acad_org,
20          ffr.period_id,
21          ffr.day_of_week
```

▶ (3) Spark Jobs

Table ▾  +

| | institution | strm | acad_group | acad_org | period_id | day_of_week | aws |
|---|---|---|---|---|---|---|---|
| 1 | UMNTC | 1229 | TCLA | 10976 | 2 | wed | 12000 |
| 2 | UMNTC | 1229 | TIOT | 11130 | 4 | mon | 13800 |
| 3 | UMNTC | 1229 | TIOT | 11130 | 5 | thurs | 9900 |
| 4 | UMNTC | 1229 | TCLA | 10984 | 8 | thurs | 10500 |
| 5 | UMNTC | 1229 | TCLA | 10956 | 7 | wed | 6900 |
| 6 | UMNTC | 1229 | TALA | 10832 | 2 | wed | 1500 |

https://z.umn.edu/databricks_asr

# Transformation Results



https://z.umn.edu/databricks_asr

# Load

# Databricks

## Benefits

- We can store data from multiple sources in many formats

- We can let data experts transform or analyze that data

- The results of their work can be used by others

# Databricks

## Differences from other systems

- Boomi API

- Data Warehouse

# Databricks

## Next Steps

- Work with OIT to improve access to Oracle

- Investigate Automation and Operation

- More proof of concepts

# Questions

https://z.umn.edu/databricks_asr