

# **A different type of DAG: Data pipelines principles for epidemiology**

Emily Riederer  
*Data Science Lead at Capital One*

# Disclosures

I am an employee of Capital One Financial Corporation

I have no relevant financial relationships to disclose

# Industry data science can learn a lot from epidemiology

Epi

**Data**

**Analytical Approach**

**Precise Communication**

Imperfection expected

Prototypical study designs

Estimands

Data Management

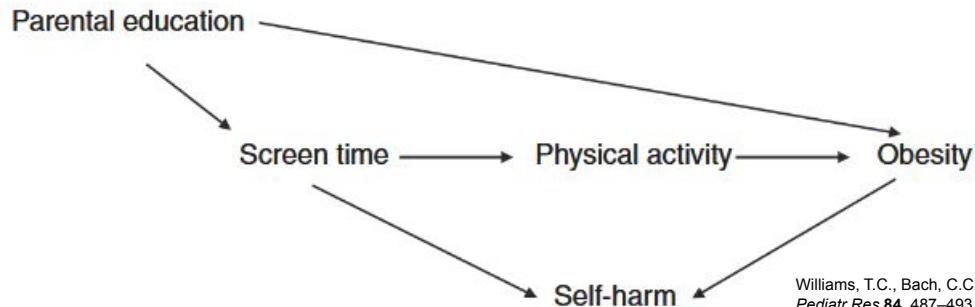
Packaged Code

Deployed Outputs

DS

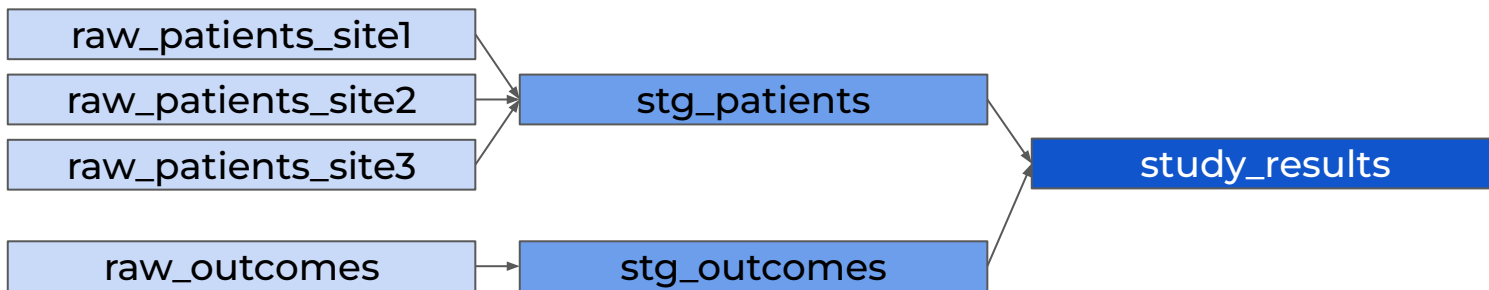
# Data pipelines are a different sort of DAG that encode the computational versus conceptual data generating process

Epi



Williams, T.C., Bach, C.C., Matthiesen, N.B. *et al.* Directed acyclic graphs: a tool for causal studies in paediatrics. *Pediatr Res* **84**, 487–493 (2018). <https://doi.org/10.1038/s41390-018-0071-3>

DS



# Data pipelines and databases serve many needs in industry with parallels for individual analytical projects

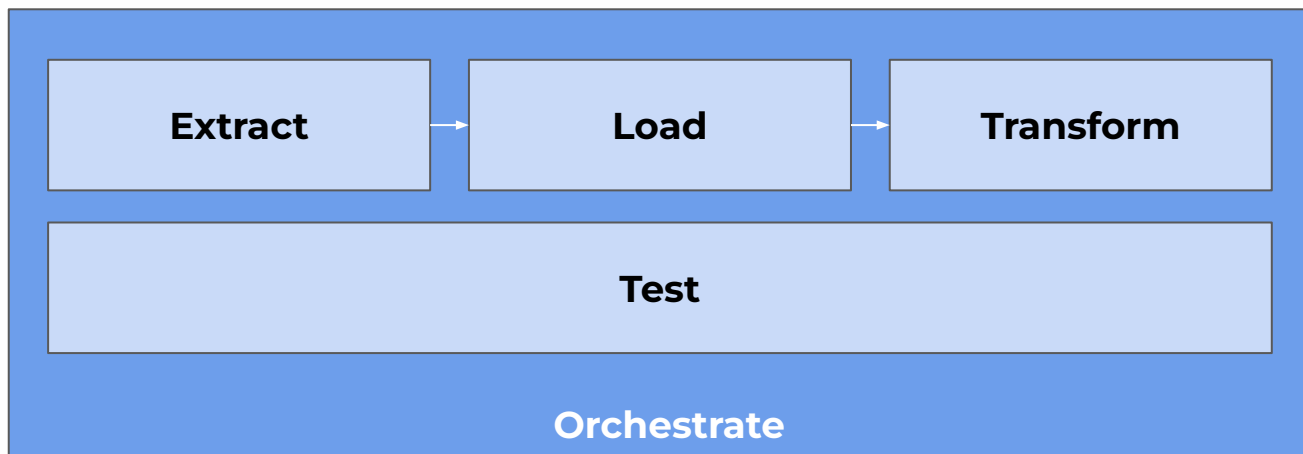
## Data Pipelines

- Data lineage and quality control
- New data ingestion, integration, and updates
- Sharing data artifacts across many users
- Sharing data artifacts across many use cases

## Project Benefit

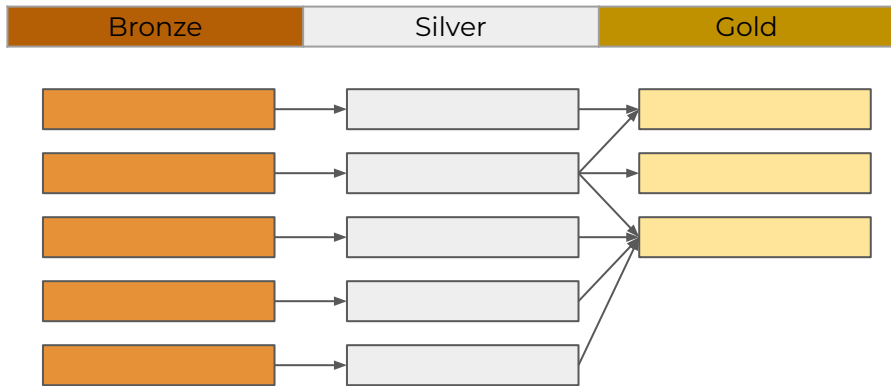
- Reproducibility
- Harmonizing disparate sources
- Collaboration
- Reuse

# Today, we will discuss some useful tools and ideas throughout the data pipeline ecosystem



# ARCHITECTURE: Pipelines preserve data artifacts at raw, standardized, and transformed stages

## Medallion Architecture



### Bronze (= Raw)

- Unmodified source data

### Silver (= Structure)

- Standardized file structure
- Cleaned column name
- Casted data types
- Light re-encoding

### Gold (= Prepared)

- Transformed and merged data
- Derived fields
- Ready for analytical use

# ARCHITECTURE: Preserving artifacts at multiple stages balances reproducibility with analytical efficiency

## In Database Schema

- data.raw
- data.structured
- data.clean

## In File System

- data/raw/
- data/structured/
- data/clean/

## Raw

- Preserves full reproducibility
- Enables testing for upstream errors pre-transformation

## Structured

- Preserves data flexible enough to transform in different ways

## Prepared

- Ready for use without rerunning expensive computations
- Ensures collaborators made same design decisions



# EXTRACT: Many different techniques can be used to collect data across disparate and sometimes inconvenient sources

## Method

**Download:** Download CSV or Parquet file from known or dynamic location

**API Call:** R {httr2} package or python {requests}

**Scraping:** Extract content from underlying webpage structure

**Headless Browsing:** Use code to navigate to URLs, click buttons, etc.

**Computer Vision:** Parse characters from an image or webpage based on visual characteristics

## When data is...

Available at a link to a standalone file, e.g. in an S3 bucket

When data is exposed through an API

When target content is rendered on the server side and in the page's source, e.g. HTML table

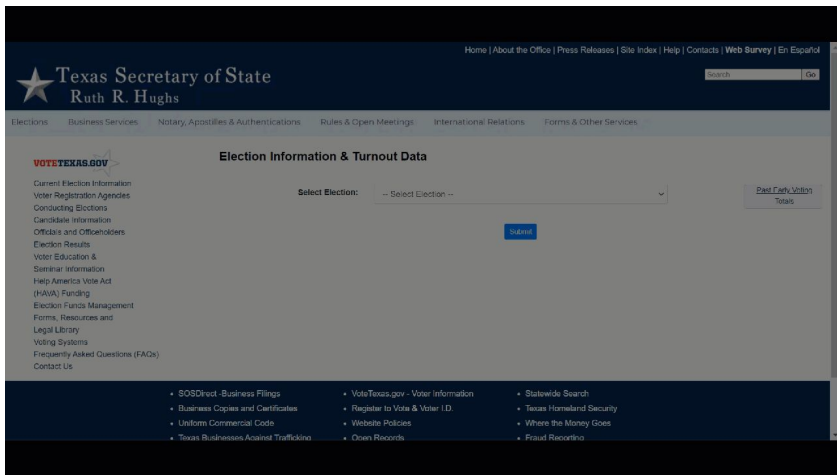
When target content is rendered client-side or you need to navigate through a UI

When data is trapped within a complex UI or dashboard with no access to the source

*Increasing Complexity*

# EXTRACT: Automation can save manual data collection effort

## Example: Headless Browsing with Playwright



```
from playwright.sync_api import sync_playwright
```

```
with sync_playwright() as p:
```

```
    browser = p.firefox.launch()
    context = browser.new_context(accept_downloads = True)
    page = context.new_page()
```

```
# pick selections from menu
```

```
page.goto(url)
page.select_option('#idElection',
    label = "2020 NOVEMBER 3RD GENERAL ELECTION")
page.click('#electionsInfoForm button')
page.wait_for_selector('#selectedDate')
```

```
page.select_option('#selectedDate', value = target_date)
page.click('#electionsInfoForm button:nth-child(2)')
page.wait_for_selector('"Generate Statewide Report"')
```

```
# download report
```

```
with page.expect_download() as download_info:
    page.click('"Generate Statewide Report"')
download = download_info.value
download.save_as(target_file)
```

# LOAD: Despite being the most common data storage formats, CSVs and Excel files can be brittle for computing

id_patient <str>	treated <int>	name <str>
'01'	1	Emily
'02'	1	Malcolm,
'03'	1	Travis

→  
`write_csv('patients.csv')`  
`read_csv('patients.csv')`

Wrong type inferred,  
leading zeros ignored

id_patient <int>	treated <int>	name <str>
1	1	Emily
2	1	Malcolm,
3	1	Travis

Could break file or  
skip row if not  
escaped

## BONUS:




DuckDB's [CSV Sniffer](#) can help identify and quarantine bad rows that keep you from opening a file

# LOAD: Parquet has numerous advantages over CSV as a standard data format




- 1 Similar scripting read/write** Easily handled by {readr} or {arrow}
- 2 Lower risk of accidental editing** Not easily opened in UI or modified without explicit read/write
- 3 Preserves data types** Remembers column types to avoid data corruption
- 4 Storing strategies allow efficient reads** Allows for partitioned read/writes for pre-filtering of large files
- 5 Smaller file sizes** Compressed on disk

# LOAD: Partitioning improves data organization and aids in efficient reading of large files




```
arrow::write_dataset(  
  data,  
  path = 'data/',  
  format = "parquet",  
  partitioning = "year"  
)
```

 data/year=2021/  
 data/year=2022/  
 data/year=2023/

```
arrow::open_dataset(  
  'data'  
)
```

 data/year=2021/  
 data/year=2022/  
 data/year=2023/

```
arrow::open_dataset(  
  'data/year=2022/'  
)
```

 data/year=2021/  
 data/year=2022/  
 data/year=2023/

# TRANSFORM: Data transformation choices help us organize our data for maximum success across projects

## Data Modeling

Organizing the relationship between types of information in your warehouse

## Table Design

Designing table contents, column names, documentation, etc. for maximum usability

## Materialization

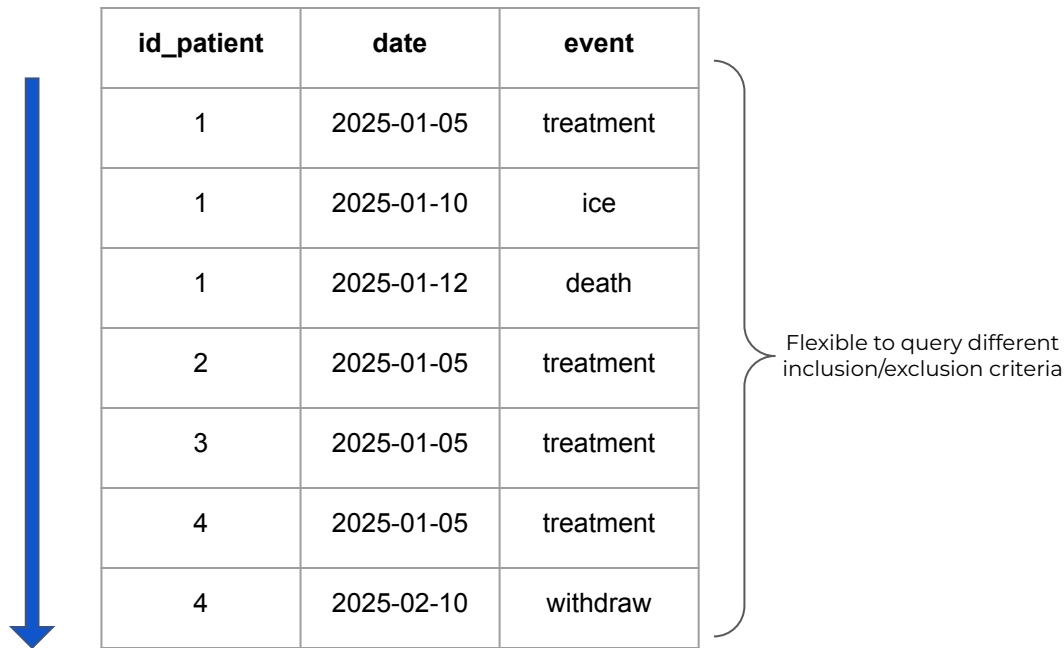
Deciding which intermediates you save as tables versus what exists only in transient code

# TRANSFORM: Keep the simplest data structures as long as possible (“shift right”) to preserve flexibility

id_patient	date	event
1	2025-01-05	treatment
1	2025-01-10	ice
2	2025-01-05	treatment
3	2025-01-05	treatment
3	2025-03-14	death
4	2025-01-05	treatment
4	2025-02-10	withdraw

id_patient	treated	withdrew	ice	death
1	1		1	
2	1			
3	1			1
4	1	1		

# TRANSFORM: Keep the simplest data structures as long as possible (“shift right”) to preserve flexibility



id_patient	date	event
1	2025-01-05	treatment
1	2025-01-10	ice
1	2025-01-12	death
2	2025-01-05	treatment
3	2025-01-05	treatment
4	2025-01-05	treatment
4	2025-02-10	withdraw

Easily updated with new records  
(favor appends to joins)



# TRANSFORM: Consistent column naming conventions across tables aids comprehension and downstream analysis code

Stub	Semantics
ID	Unique entity identifier
IND / IS	Binary 0/1 indicator; rest of name describes 1 condition
BIN	Binary 0/1 indicator; rest of name describes 1 condition
N	Count of quantity or event occurrences
AMT	Sum-able real number amount ("denominator free")
...	...

```
library(dplyr)

data %>%

  group_by(NM_STATE) %>%

  summarize(

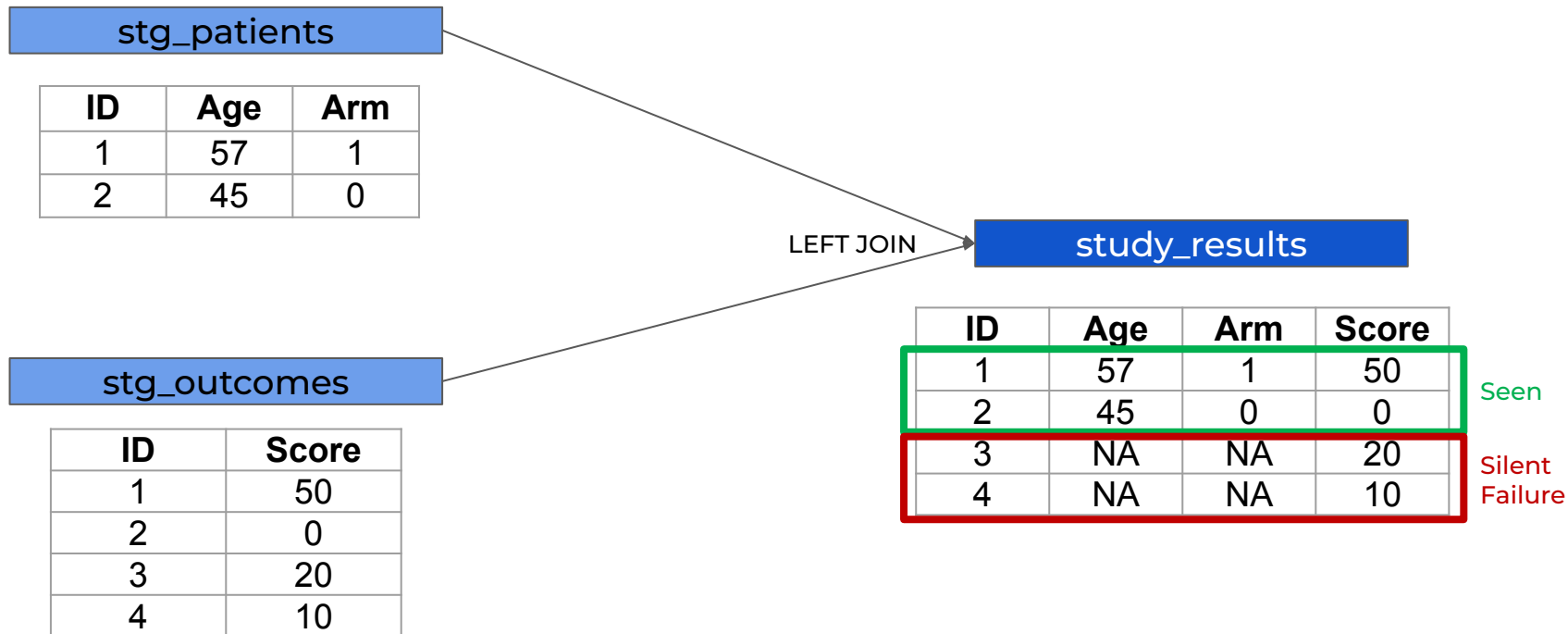
    across(starts_with("IND"), mean),

    across(contains("_ACTL_"), sum)

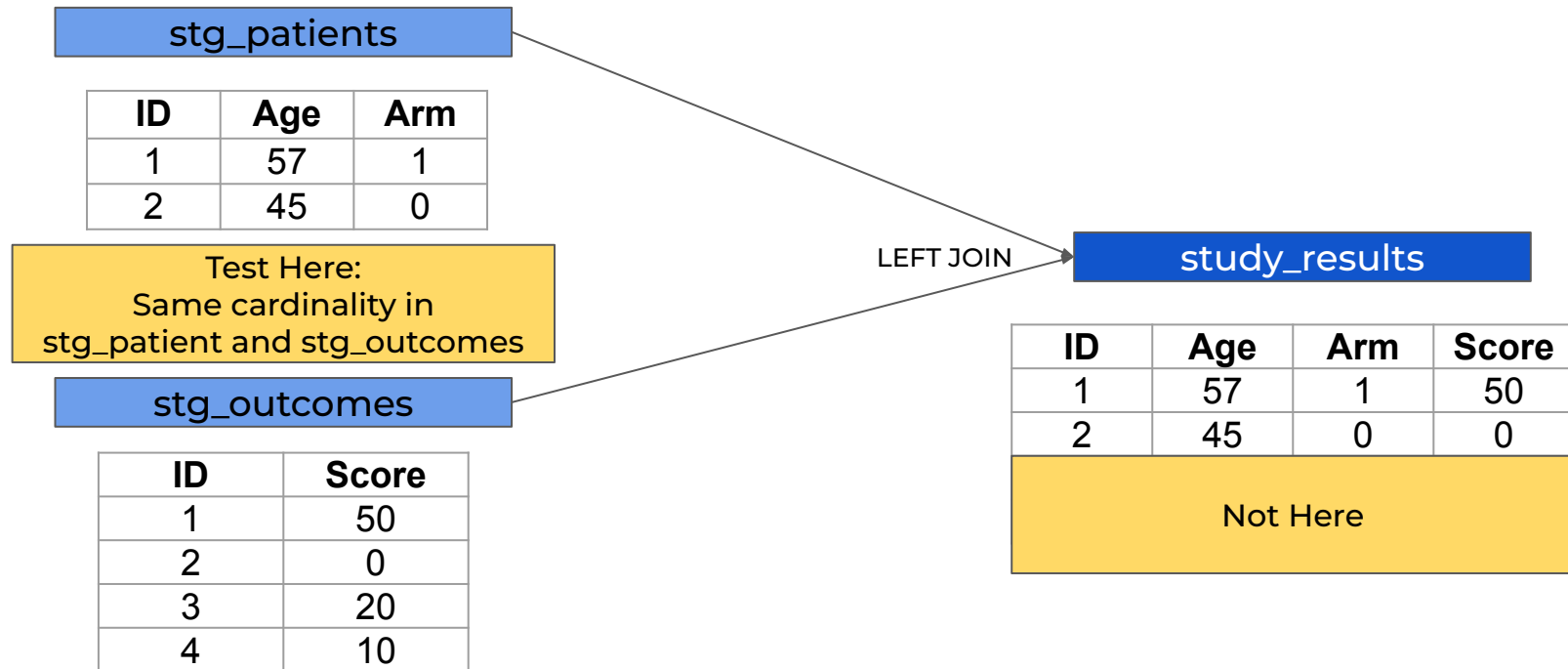
  )

#> # A tibble: 51 x 4
#>   NM_STATE      IND_COUNTY_HPSA  N_CASE_ACTL  N_DEATH_ACTL
#>   <chr>                <dbl>        <dbl>        <dbl>
#> 1 Alabama              0.149        455582        7566
#> 2 Alaska               0.235         51338         250
#> 3 Arizona              0            753379       13098
```

# TEST: A multistage pipeline allows us to conduct hypothesis-drive upstream data quality checks (“shift left”)



# TEST: A multistage pipeline allows us to conduct hypothesis-drive upstream data quality checks (“shift left”)



# ORCHESTRATE: Orchestration tools help us run many steps of an analysis pipeline in order to keep dependencies consistent

## Define Order of Operation

Define the dependency graph to ensure all steps are rerun in order for consistent results

## Restarts & Error Handling

Be able to arbitrarily initiate pipeline start based on changes and/or specific node

## Execution & Logging

Run pipeline end-to-end and track signs of success

## ORCHESTRATE: There are many tools available for orchestration depending on your project needs

	<b>dbt</b>	<b>make</b>	<b>{targets} R package</b>	<b>Dagster</b>
Works with	SQL/python	Arbitrary command line operations that read/write files	R	Python
Use cases	Pure data pipeline	End-to-end project but fewer bells and whistles	End-to-end automation of data prep through modeling; fuller featured but higher learning curve	

# Thank you!