# A different type of DAG: Data pipeline principles for epidemiology

## Emily Riederer
*Data Science Lead at Capital One*

Disclosures: No relevant financial interests exist.

# Industry data science can learn a lot from epidemiology

Epi

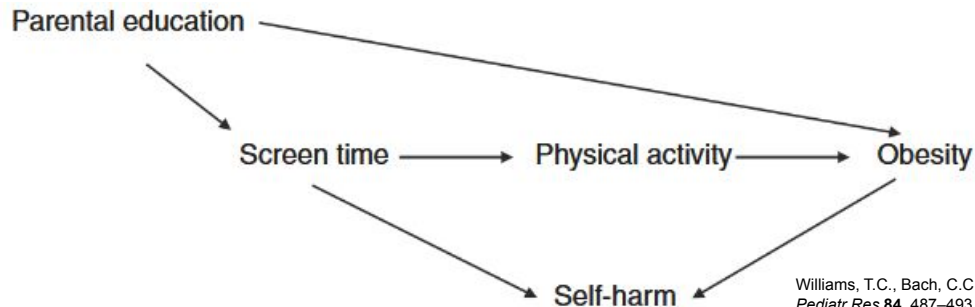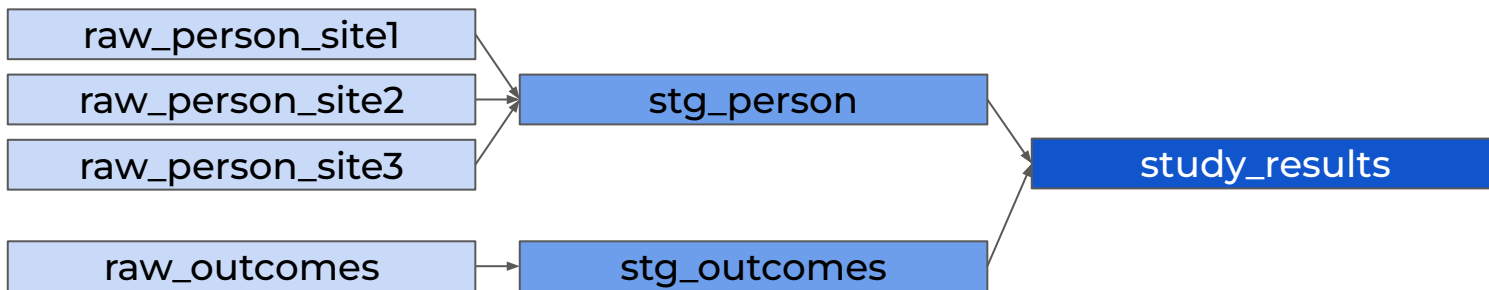| Data | Analytical Approach | Precise Communication |
|------|---------------------|-----------------------|
| Imperfection expected | Prototypical study designs | Estimands |
| Data Management | Packaged Code | Deployed Outputs |

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
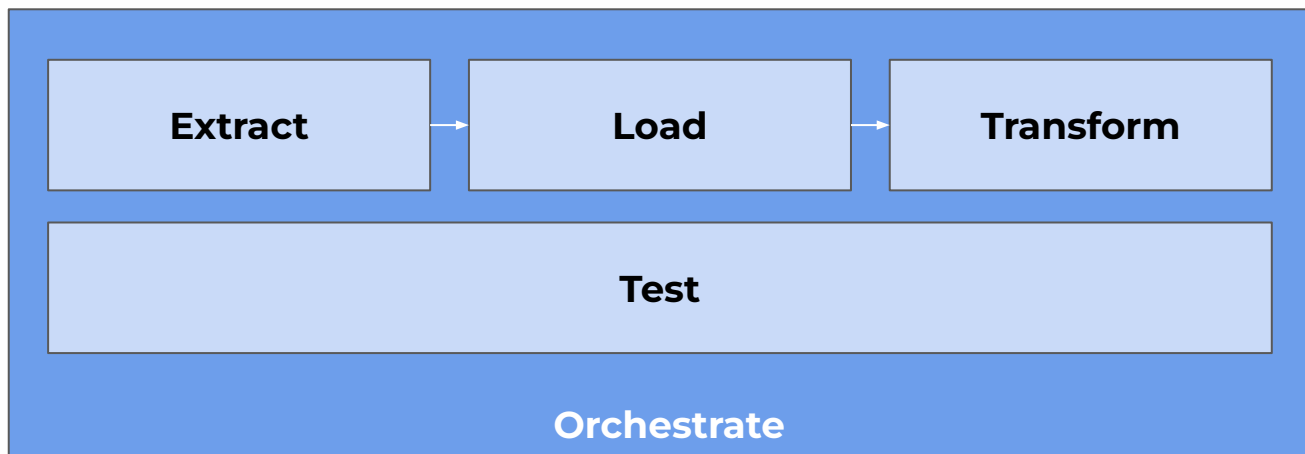
## Data Pipelines

- Data lineage and quality control

- New data ingestion, integration, and updates

- Sharing data artifacts across many users

- Sharing data artifacts across many use cases

- Produce complex outputs (models, reports)

## Project Benefit

- Reproducibility

- Harmonizing disparate sources

- Collaboration

- Reuse across projects

- In-sync analysis

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

# ARCHITECTURE: The same principles apply across platforms – all you need to get started is a hard drive

|  | **Local** | **Shared** | **Cloud** |
|---|---|---|---|
| **Storage** | 📁 Local Drive | 📁 Network Drive<br>🗄 Google Cloud Storage<br>🗄 Amazon S3 | 🗄 Google BigQuery DB |
| **Compute** | R, python, SQL + dbt, or any scripting language<br>🗄 duckdb | | 🗄 Google BigQuery<br>SQL + dbt (data pipeline)<br>Dagster (orchestration) |

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

## Medallion Architecture



| Bronze | Silver | Gold |

**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

ELT TO

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

## In Database Schema

- ▤ `data.raw`
- ▤ `data.structured`
- ▤ `data.prepared`

## In File System

- 📁 `data/raw/`
- 📁 `data/structured/`
- 📁 `data/prepared/`

**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

ELT
TO

# ARCHITECTURE: Read/write patterns make the pipelines acyclic, idempotent, resilient to re-runs

```
df = readr::read_csv('data/raw/persons.csv')

df_cleaned = dplyr::filter(df, dt_enroll < cutoff)

readr::write_csv(df_cleaned, 'data/structured/persons.csv')
```

```
df = readr::read_csv('data/raw/persons.csv')

df = dplyr::filter(df, dt_enroll < cutoff)

readr::write_csv(df, 'data/raw/persons.csv')
```

```
df = readr::read_csv('data/raw/persons.csv')

df_cleaned = dplyr::filter(df, dt_enroll < cutoff)

readr::write_csv(df_cleaned, 'data/raw/persons.csv')
```

ELT TO

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

## When data is available as...

Link to a standalone file, e.g. in an S3 bucket

Exposed through a REST API

Rendered server-side in the webpage's source, e.g. a static HTML table

Rendered dynamically in a web-based UI

Trapped within a complex UI or report

## Method

**Download:** R {readr} or python {urllib.requests}

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

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

**Headless Browsing:** Use code to navigate a UI

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

*Increasing Complexity*

ELT TO

# EXTRACT: Automation can save manual data collection effort

**Example:**
**Headless Browsing with Playwright**

```
# 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)
```

ELT
TO

# EXTRACT: Automation can save manual data collection effort

**Example:
Optical Character Recognition with Tesseract**



```
# take screenshot with Playwright
page.goto(url)
page.wait_for_load_state(state = 'networkidle')
time.sleep(30)
page.screenshot(path = 'ri.png')
```

```
# parse image to data with Tesseract
img = cv2.imread('ri.png')
text = pytesseract.image_to_string(img)
n_tot = re.search('Turnout\n\n(\d+)', text).group(1)
n_mail = re.search('Mail Ballots', text).group(1)
```

ELT
TO

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

Wrong type inferred,
leading zeros ignored

| zip_code<br><str> | treated<br><int> | name<br><str> |
|---|---|---|
| '06830' | 1 | Emily |
| '27514' | 1 | Barret,M. |

write_csv('persons.csv')
read_csv('persons.csv')

| zip_code<br><int> | treated<br><int> | name<br><str> |
|---|---|---|
| 6830 | 1 | Emily |
| 27514 | 1 | Barret,M. |

Could break file or
lose data

BONUS:
DuckDB's CSV Sniffer can help identify and quarantine bad rows that keep you from opening a file

ELT
TO

# LOAD: Parquet files are like a zipfile of data and metadata

## CSV

- Human readable and editable
- All data in single text file
- Row-based

```
zip_code, treated, name
06834, 1, Emily
12345, 1, Barret,M
```

## Parquet

- Compressed and machine readable
- Bundles data and metadata
- Column-based

```
Metadata
```
```
zip_code: 06834 12345
```
```
treated:1 1
```
```
name:Emily Malcolm,B
```

ELT
TO

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

**1**   **Similar scripting read/write**     Easily handled by R's {readr} or {arrow}

**2**   **Lower risk of accidental modification**     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

ELT TO

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

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

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

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

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

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

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

ELT
TO

# TRANSFORM: Data modeling provides a framework for integrating data from diverse sources and structures

Site 1

Site 1

Site 3

Site 3

Site 3

## Data Modeling

Restructure data into standard schemas representing different types of information

## Table Design

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

Unified Schema

ELT TO

# TRANSFORM: Data modeling principles can help us partition information into entities and relationships

| id | age |
|----|-----|
| 1  | 35  |
| 2  | 27  |

| id | name  |
|----|-------|
| 1  | Chris |
| 2  | Jenny |

| id | visit | test_result | name | age |
|----|-------|-------------|------|-----|
| 1  | 2025-01-05 | 13 | Chris | 35 |
| 1  | 2025-01-10 | 1  | Chris | 35 |
| 1  | 2025-01-12 | 2  | Chris | 35 |
| 2  | 2025-01-05 | 14 | Jenny | 27 |
| 2  | 2025-02-10 | 0  | Jenny | 27 |

**person**
(1 row / person)

← identified by

**persons_pii**
(1 row / person)

attends

**persons_visits**
(1 row / person x visit)

| id | visit | test_result |
|----|-------|-------------|
| 1  | 2025-01-05 | 13 |
| 1  | 2025-01-10 | 1  |
| 1  | 2025-01-12 | 2  |
| 2  | 2025-01-05 | 14 |
| 2  | 2025-02-10 | 0  |

ELT TO

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

## Activity Schema

- 1 record per entity x event
- Long and narrow
- Flexible but may require more processing

| id_person | date | event |
|-----------|------|-------|
| 1 | 2025-01-05 | treatment |
| 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 |

## One Big Table (OBT)

- 1 record per entity
- Wide
- Already imposes some analytical decisions

| id_person | treated | withdrew | death |
|-----------|---------|----------|-------|
| 1 | 2025-01-05 | | |
| 2 | 2025-01-05 | | |
| 3 | 2025-01-05 | | 2025-03-14 |
| 4 | 2025-01-05 | 2025-02-10 | |

ELT
TO

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

| id_person | date | event |
|:---:|:---:|:---:|
| 1 | 2025-01-05 | treatment |
| 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 |

Flexible to query different inclusion/exclusion criteria

Easily updated with new records
(favor appends to joins)

ELT TO

# 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")

**stg_persons**

| ID | Age | Arm |
|----|-----|-----|
| 1  | 57  | 1   |
| 2  | 45  | 0   |

LEFT JOIN

**study_results**

| ID | Age | Arm | Score |
|----|-----|-----|-------|
| 1  | 57  | 1   | 50    |
| 2  | 45  | 0   | 0     |
| 3  | NA  | NA  | 20    |
| 4  | NA  | NA  | 10    |

Seen

Silent Failure

**stg_outcomes**

| ID | Score |
|----|-------|
| 1  | 50    |
| 2  | 0     |
| 3  | 20    |
| 4  | 10    |

emilyriederer.com  ||  @emilyriederer.bsky.com

ELT TO

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

**stg_persons**

| ID | Age | Arm |
|----|-----|-----|
| 1  | 57  | 1   |
| 2  | 45  | 0   |

**Test Here:**
Same cardinality in stg_persons and stg_outcomes

**stg_outcomes**

| ID | Score |
|----|-------|
| 1  | 50    |
| 2  | 0     |
| 3  | 20    |
| 4  | 10    |

LEFT JOIN

**study_results**

| ID | Age | Arm | Score |
|----|-----|-----|-------|
| 1  | 57  | 1   | 50    |
| 2  | 45  | 0   | 0     |

Not Here

ELT TO

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

Ability to execute run starting at any stage

raw_persons_site1

raw_persons_site2

raw_persons_site3

stg_persons

Test

study_results

raw_outcomes

stg_outcomes

Define order of operations

Logging or handling of test failures, runtime errors

ELT TO

# ORCHESTRATE: Many different tools exist for orchestration

## makefile

- General purpose config file
- Runs any shell commands based on user-defined dependencies
- Higher learning curve

```
all: ./data/clean.json

./data/raw.json: ./data/raw.json ./ingest.py
        mkdir -p data
        ./ingest.py

./data/clean.json: ./data/raw.json ./clean.py
        ./clean.py
```

## dbt

- SQL framework that infers dependencies
- Compatible with any database including local compute (duckdb)
- Lower learning curve but less flexible

```
select id, name
from raw.persons
where dt_enroll = '2025-01-01'
```
persons.sql

```
select id, name, outcome
from {{ref('persons')}}
   left join
   {{ref('outcomes')}}
   using (id)
```
results.sql

Makefile example from https://datasciencesouth.com/blog/make/#a-make-data-pipeline-example
See also R package {targets} and Dagster

E L T
T O

# Big ideas from data pipelining can help improve data handling in projects of any size of complexity

**Architecture**    Start with a structured read/write process to your hard drive controlled by **R, python,** or **SQL**

**Extract**    Automate manual extraction with tools like **playwright** and **tesseract**

**Load**    Store and share **Parquet** files for type-safe resilience

**Transform**    Apply data modeling techniques like **entity-relationship mapping**, **Boyce-Codd normalization** and the **Activity schema** to standardize structure

**Test**    Test intermediate artifacts to detect failures at their source

**Orchestrate**    Adopt automated tooling to run all steps in sync with **make**, **dbt**, or **targets**

# A few more examples and resources are available online

**Architecture**    [Build NC data pipeline using local drive, duckdb, and Arrow](#)

**Extract**    [Different web scraping techniques across 6 government websites](#)

**Load**

**Transform**    [Column Names as Contracts](#)
[Column Name Contracts in dbt](#)

**Test**    [Understanding the data (error) generating process](#)
[Hypothesis-driven data testing with grouped checks](#)

**Orchestrate**

# Thank you!