# Creating reproducible studies

Martijn Schuemie

# What is reproducibility?
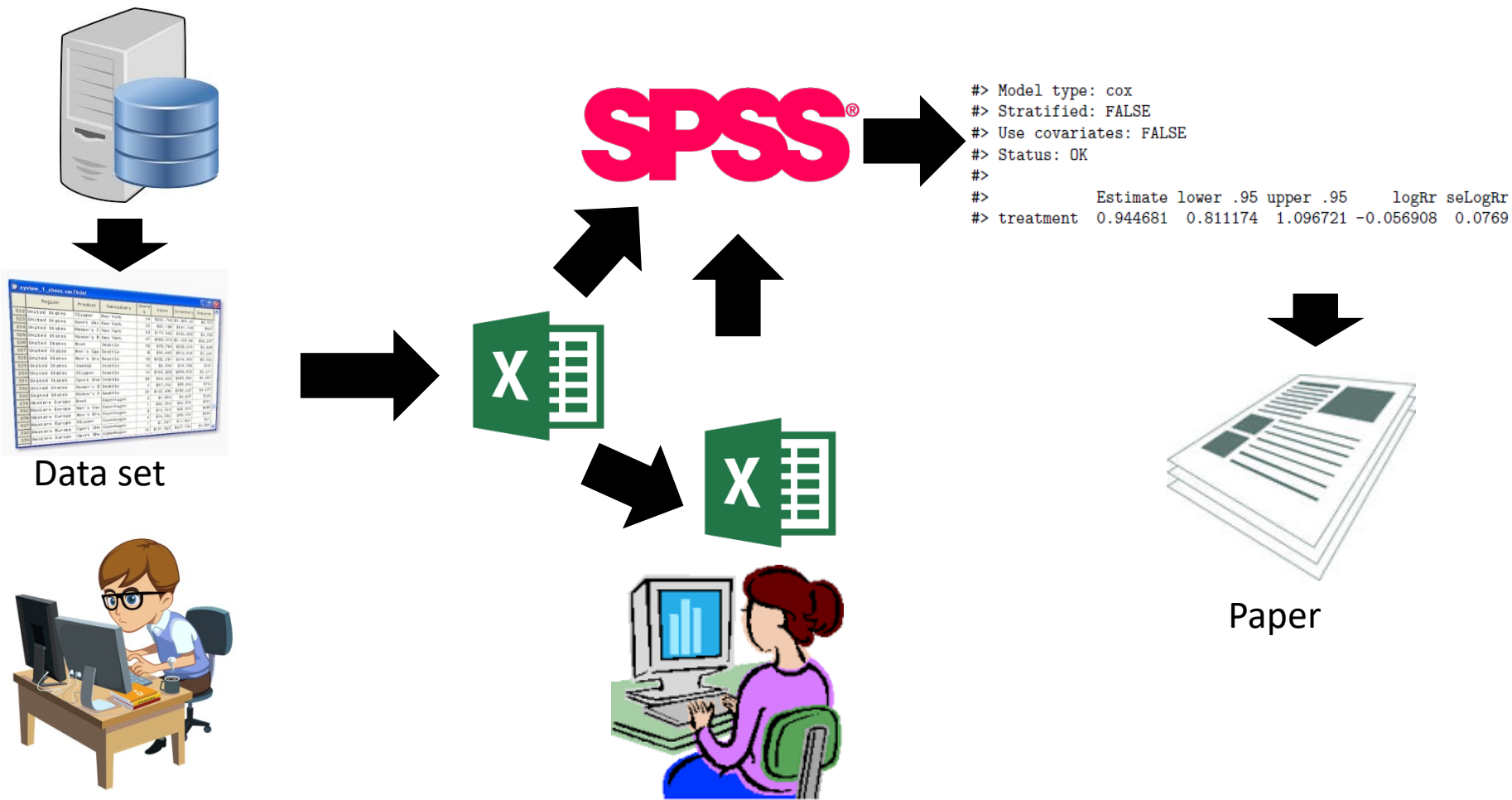
| Desired attribute | Question | Researcher | Data | Analysis | | Result |
|---|---|---|---|---|---|---|
| **Repeatable** | Identical | Identical | Identical | Identical | = | Identical |
| **Reproducible** | Identical | Different | Identical | Identical | = | Identical |
| **Replicable** | Identical | Same or different | Similar | Identical | = | Similar |
| **Generalizable** | Identical | Same or different | Different | Identical | = | Similar |
| **Robust** | Identical | Same or different | Same or different | Different | = | Similar |
| **Calibrated** | Similar (controls) | Identical | Identical | Identical | = | Statistically consistent |

Ensuring the **analysis** can be kept identical allows for **repeatable**, **reproducible**, **replicable**, **generalizable**, and **calibrated** science
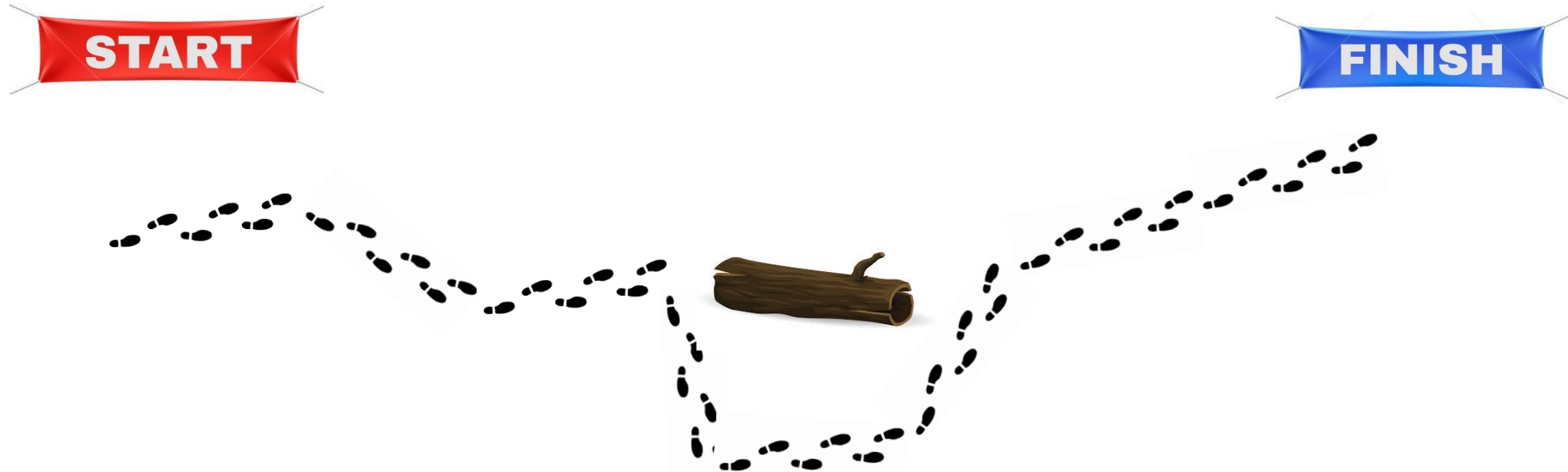
Source: The Book of OHDSI

# What do observational studies currently look like?



Data set

SPSS

```
#> Model type: cox
#> Stratified: FALSE
#> Use covariates: FALSE
#> Status: OK
#>
#>            Estimate  lower .95  upper .95      logRr  seLogRr
#> treatment  0.944681   0.811174   1.096721  -0.056908   0.0769
```
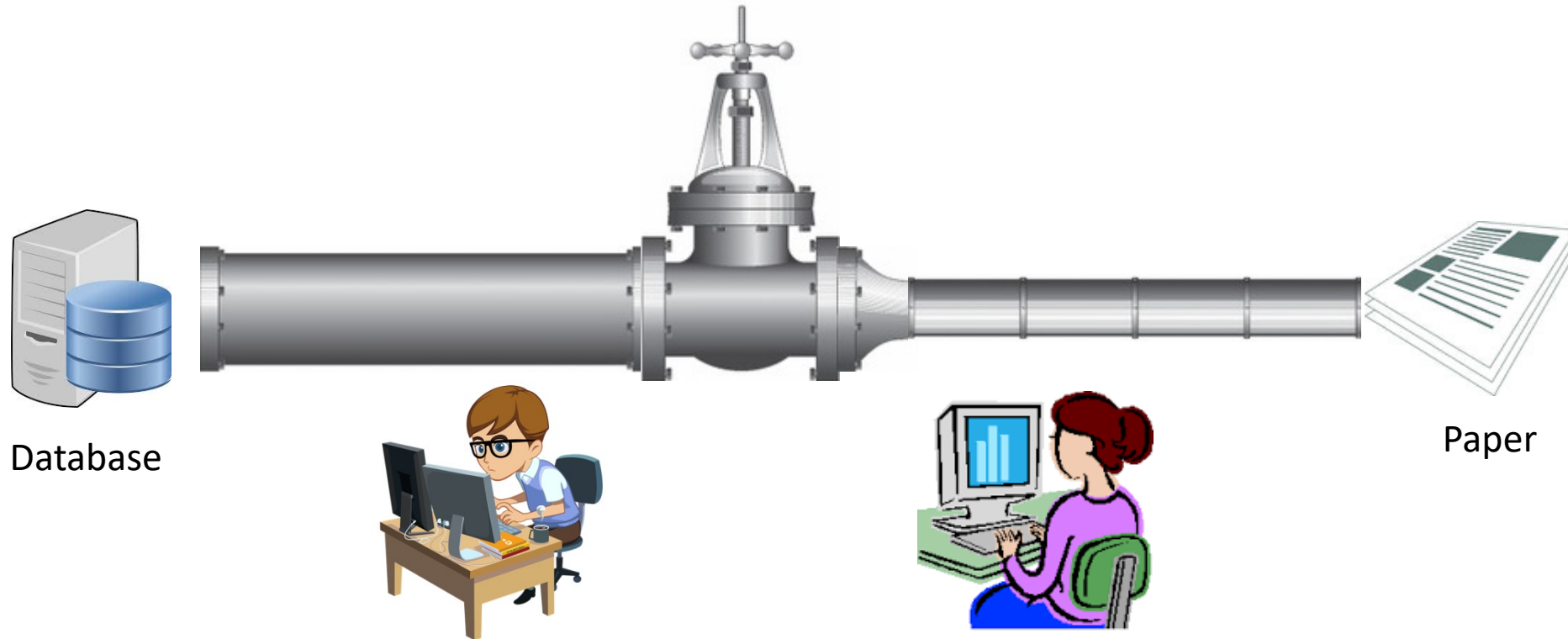
Paper

# A journey from data set to paper



Most epidemiologists view a study as a journey from data set to paper.
- The protocol might be your map
- You will come across obstacles that you will have to overcome
- Several steps will require manual intervention
- In the end, it will be impossible to retrace your exact steps

# What should OHDSI studies look like?



Database

Paper

A study should be like a pipeline
- A fully automated process from database to paper
- 'Performing a study' = building the pipeline

# No peeking while developing the pipeline!



Source: The Book of OHDSI

# OHDSI study packages available as open source

- Each OHDSI study comes with a protocol as well as a fully executable R package
- R packages and protocols are available on

https://data.ohdsi.org/OhdsiStudies/

# Anatomy of a typical study package

# Preserving the compute environment





Pro:
- Runs at most institutions
- Integrated in R and R-Studio
- Lightweight

Cons:
- Only preserves R package versions, not R itself
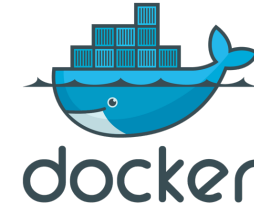- Still being developed

Pro:
- Preserves entire compute environment

Cons:
- Some OHDSI sites are not allowed to use Docker
- Docker images can get big
- Outside of R: requires additional tools + knowledge to use

# Summary

- Our analyses must be 100% repeatable (on same or different data) to allow for **repeatable**, **reproducible**, **replicable**, **generalizable**, and **calibrated** science.

- Each study must be written as a **pipeline**, automatically transforming data in the **CDM** to the **study outputs** (tables, figures, etc.).

- Current OHDSI practice is to create a **study package** (in ohdsi-studies) and capture the compute environment in *renv*.
  - Exploring additional use of Docker