



# Reproducible Research in R

A Talk on How to Do the Same Thing More Than Once

Aaron Peikert

Max Planck Institute for Human Development

# Reproducibility

# Reproducibility

same thing in, same thing out

As a student:

Conceptually

SEXY

Technologically

SEXY

At the begin of my PhD:

Conceptually

SEXY

Technologically

SEXY

At the end of my PhD:

Conceptually

SEXY

Technologically

SEXY

Now:

Conceptually

SEXY

Technologically

SEXY

# Conceptual



Conceptual

Conceptual

# What is the purpose of reproducibility?

Same data → Same Black box → Same results

Same data → Same Black box → Same results  
Reproducible?

Statistical Models  
fit  
data

Statistical Models

overfit

data

By how much?



By how much?  
or  
Reproducibility and Overfit

## By how much?

- $R^2_{\text{adj.}}$
- $C_p$
- $AIC$
- Cross Validation

$$R^2_{\text{adj.}} = R^2 - (1 - R^2) \frac{p}{n - p - 1}$$

$$C_p = \frac{\sum_{i=1}^n (\hat{y}_i - y_i)^2}{\sigma_e^2} - n - 2p$$

# Reproducibility?

# Reproducibility?

$$\text{df}(\hat{y}) = \frac{\sum_{i=1}^n \text{Cov}(\hat{y}_i, y_i)}{\sigma_e^2}$$

# Reproducibility?

$$\text{df}(\hat{y}) = \frac{\sum_{i=1}^n \text{Cov}(\hat{y}_i, y_i)}{\sigma_e^2}$$

This covariance requires:

- formal derivation
- repeated computation on same or pertubated data

This covariance requires:

- formal derivation
- repeated computation on same or pertubated data

**Reproducibility.**

$$R_{\text{adj.}}^2 = R^2 - (1 - R^2) \frac{n - df}{df - 1}$$

$$C_p = \frac{\sum_{i=1}^n (\hat{y}_i - y_i)^2}{\sigma_e^2} - 3n + 2df$$



# Information Criteria

Require the Hessian around the solution for their “overfit” correction.

Computed via:

- byproduct of optimization
- via finite differences

Require the Hessian around the solution for their “overfit” correction.

Computed via:

**Reproducibility**

Cross Validation

is

Reproduction

on subsamples.

# Extended Definition

First, computational reproducibility must ensure that the same data lead to the same results.

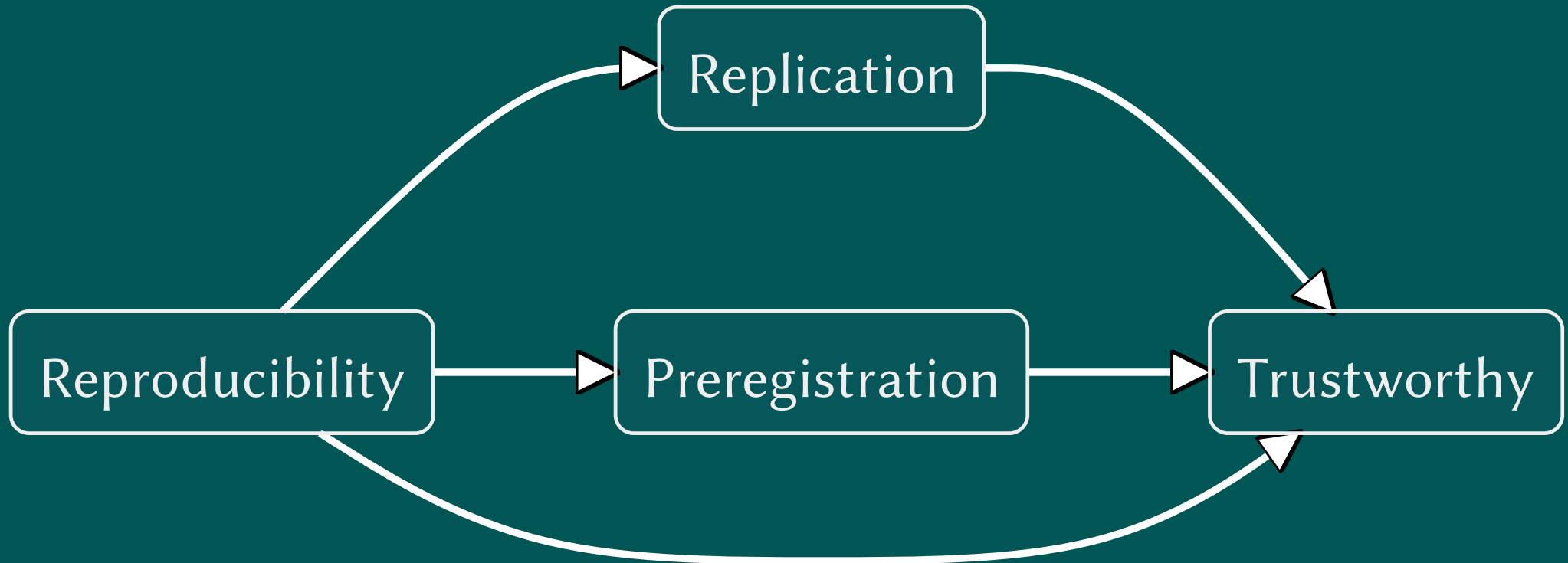
# Extended Definition

First, computational reproducibility must ensure that the same data lead to the same results.

Second, computational reproducibility must make the inductive process repeatable on similar data.



# Summary





# Technological

# Technological

For reproducibility, it really needs to be reproducible and checkable by a stranger with little time or energy to spare, because even the author will soon enough be that stranger.

— Gwern Branwen

# Now you!

Try to add your name to the next slide. Every time I show this presentation your name will be there.

### Awesome People:

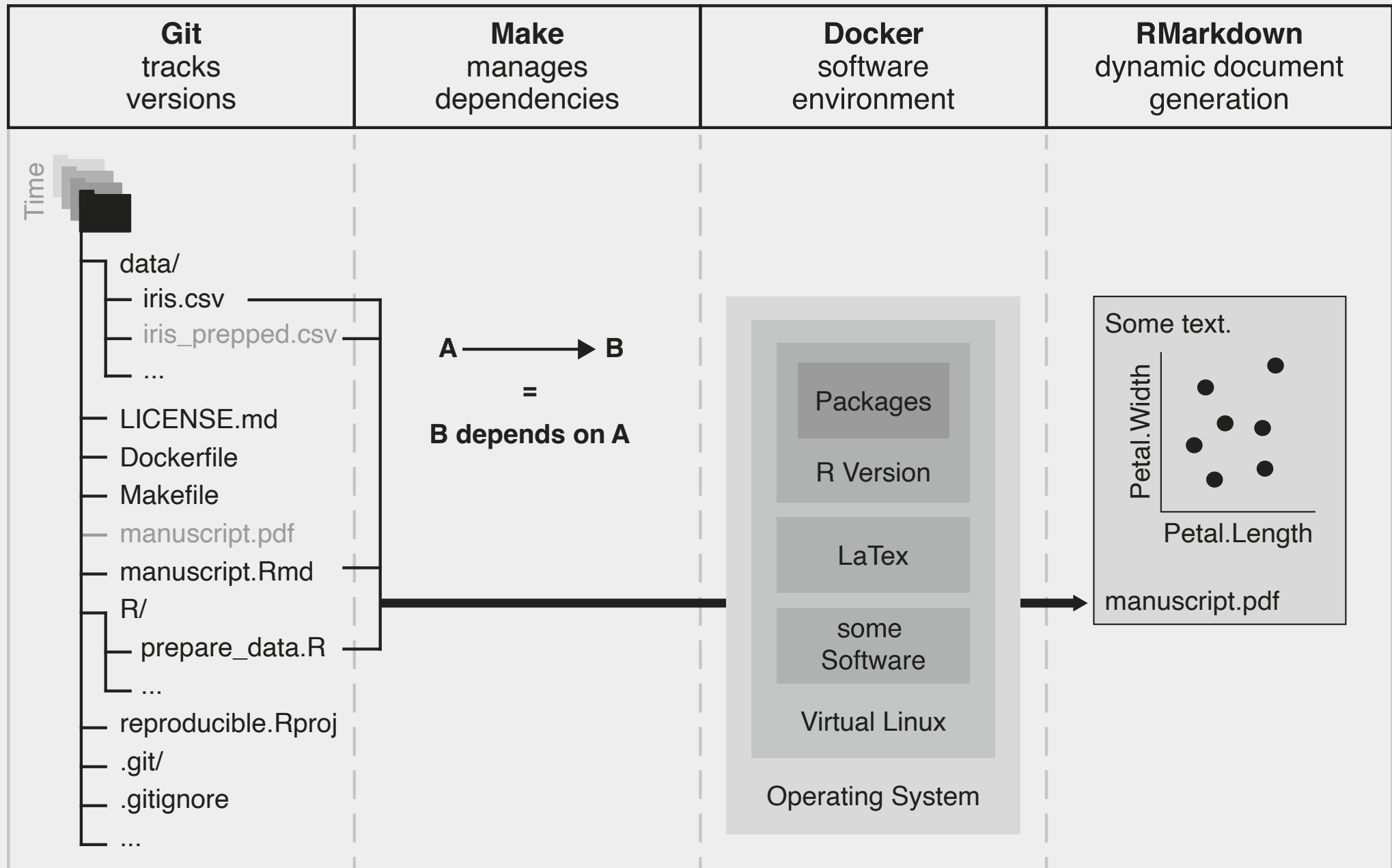
- Aaron Peikert

# Four Problems with Reproducibility

1. versioning
2. copy&paste errors
3. software dependencies
4. linking everything together

# Four Solutions for Reproducibility

1. version control
2. dynamic documents
3. software management
4. workflow orchestration



# ACRISS

Association of Car Rental Industry Systems Standards



Category	Type	Trans / Driven wheels	Fuel / air-con
E: Economy	F: SUV	A: Auto (drive unspecified)	R: Unspecified Fuel With AC
I: Intermediate	T: Convertible	B: Auto 4WD	D: Diesel With AC
S: Standard	C: 2/4 Door	D: Auto AWD	H: Hybrid With AC
F: Fullsize	D: 4-5 Door	M: Manual (drive unspecified)	E: Electric With AC
P: Premium	S: Sport		V: Petrol With AC

EDMR

# The Rental Car Model

Ride it like you stole it



# Four Solutions for Reproducibility

1. version control
2. dynamic documents
3. software management
4. workflow orchestration

R + RMarkdown + Docker + Make + Git

<https://github.com/aaronpeikert/reproducible-research>

Julia + RMarkdown + Pkg.jl + GitHub Actions + Git

<https://github.com/formal-methods-mpi/pkgmanuscript/blob/main/Dockerfile#L14>

Lua + Quarto + GitHub Actions + GitHub Actions + Git  
<https://github.com/aaronpeikert/repro-talk>

Python + Quarto + Docker + GitHub Actions + Git

<https://github.com/formal-methods-mpi/projects/pull/41/files>



R + RMarkdown + Docker + Make + Git

<https://github.com/aaronpeikert/bayes-prereg/pull/97>

# Summary

