# About this work

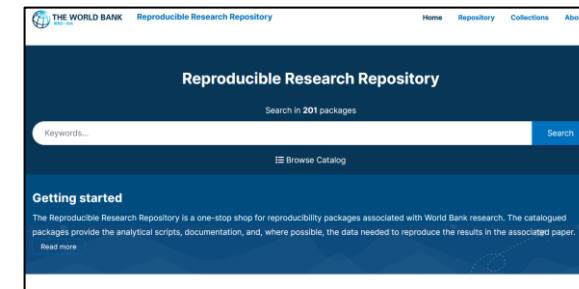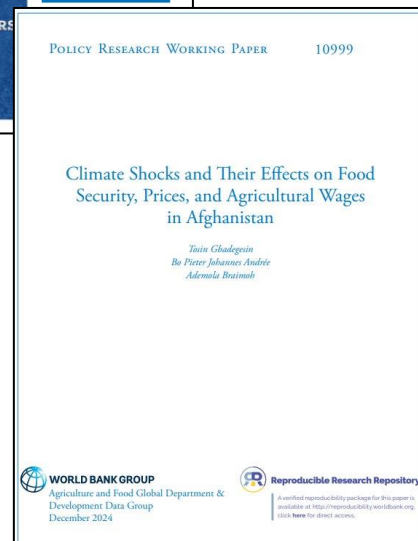- Our team has verified 200+ reproducibility packages for World Bank research
  - 77% of works reviewed used Stata
  - Only 18% are reproducible at first try

- Introducing repkit: a Stata package to address common reproducibility challenges (and make verification easier)



https://www.worldbank.org/en/research/brief/world-bank-policy-research-working-papers

https://reproducibility.worldbank.org

# 1 – Check for instabilities in your code: `reprun`

## The problem:

- Reproducibility issues arise from Stata commands that introduce unnoticed/uncontrolled randomness

- If hundreds or thousand of lines of code, costly to detect where the instability starts
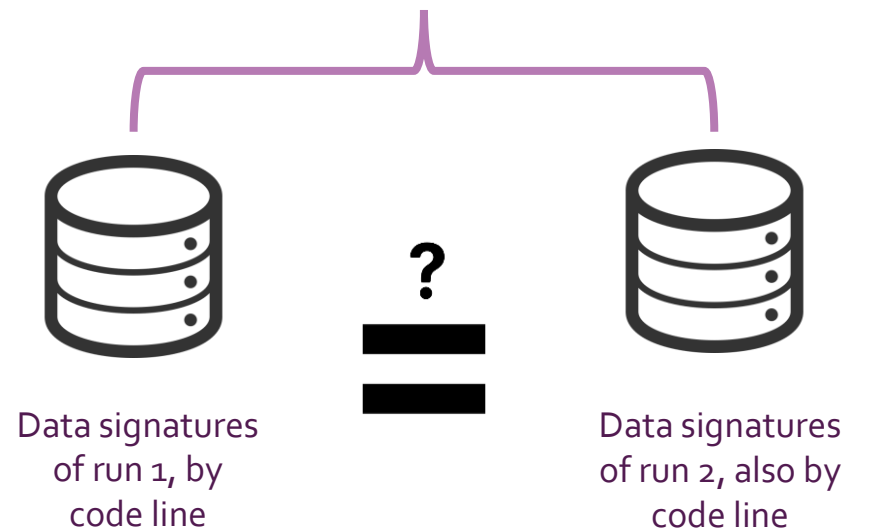
## Solution:

- `reprun` runs a do-file twice and compares intermediate results **line-by-line** across runs, flagging inconsistencies

- Works smoothly with sub do-files and loops

- Word of caution: it's only as fast as your actual code

Command
reprun "main.do"

main.do

Data signatures of run 1, by code line

?
=

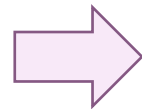Data signatures of run 2, also by code line

# 1 – Check for instabilities in your code: `reprun`

Lines 193-195: Sorting on non-unique variable and dropping observations based on sorting

```
CleanDataCoxBazar.do    ✕

192
193    sort uid1
194    by uid1:  gen dup = cond(_N==1,0,_n)
195    drop if dup==2
196    drop dup uid_r2 cbps_r2_status
197    save"RawData\instrument_cbps", replace
198
```

**Command**

```
reprun "CleanDataCoxBazaar.do", compact
```

Checking file:
└─> C:/Users/wb558768/Documents/GitHub/reprun-example/CleanDataCoxBazar.do

| Line # | Seed RNG State | | | Sort Order RNG | | | Data Checksum | | | Loop iteration: |
|--------|-------|-------|-------|-------|-------|-------|-------|-------|-------|---|
|        | Run 1 | Run 2 | Match | Run 1 | Run 2 | Match | Run 1 | Run 2 | Match |   |
| 193    |       |       |       | Change | Change | DIFF | Change | Change | DIFF |   |
| 1512   |       |       |       | Change | Change | DIFF | Change | Change | DIFF |   |
| 1515   |       |       |       | Change | Change | DIFF | Change | Change | DIFF |   |
| 1566   |       |       |       | Change | Change | DIFF | Change | Change | DIFF |   |
| 1662   |       |       |       | Change | Change | DIFF | Change | Change | DIFF |   |

`reprun` reports an inconsistency in the data checksum (data signature) of runs 1 and 2 that starts in line 193

# 2 – Avoid absolute file paths: `reproot`

## The problem:

- Code that relies on the use of absolute file paths won't work in a new computer until you adjust file paths, at least once per project

```
* Kristoffer's root path
if "`c(username)'" == "wb462869" {
    global data "C:\Users\wb462869\Dropbox\Projects\ProjectA"
    global code "C:\Users\wb462869\GitHub\ProjectA"
}
* Ben's root path
if "`c(username)'" == "bbdaniels" {
    global data "/Users/bbdaniels/Dropbox/ProjectA"
    global code "/Users/bbdaniels/GitHub/ProjectA"
}
```

## Solution:

- `reproot` manages file paths in the backend, looking in pre-defined folder locations for projects
- Needs to be set up only **once per computer**
- Works with multi-rooted projects. For example:
  - Data is on OneDrive or Dropbox
  - Code is in a GitHub repository

# 2 – No more absolute file paths: reproot

```
* Use reproot to get root paths
reproot, project("proj-a") ///
    roots("code data") prefix("prja_")
* Load data
use "${prja_data}/data/raw.dta"
* Run analysis
do  "${prja_code}/code/tab1.do"
do  "${prja_code}/code/tab2.do"
```

**reproot - settings file setup**

Description:

This sets the paths reproot will search for root files.

Select the search paths to be included:

- ☑ 2:C:/Users/WB462869/github
- ☑ 4:C:\Users\WB462869\github\repkit\src\tests\reproot
- ☑ 4:C:\Users\WB462869\Dropbox
- ☐
- ☐
- ☐
- ☐
- ☐

Add another search path:

[                    ] Browse . . .

Search depth specific to this search

4 ▲▼

Add path to possible options

# 3 – Manage dependencies simply: `repado`

## The problem:

- Missing Stata dependencies stop code execution

- Using different versions of Stata dependencies can produce different results for the same code, causing problems in reproducibility

## Solution:

- Use `repado` to temporarily change the dependencies folder (ado) in a main do-file

- Use one ado folder by project and share it with your team

- When creating a reproducibility package, include the ado folder in it

| Name | Date modified | Type |
|------|---------------|------|
| ado | 10/30/2024 9:49 AM | File folder |
| Analysis | 10/30/2024 9:49 AM | File folder |
| Cleaning | 10/30/2024 9:49 AM | File folder |
| Construct | 10/30/2024 9:49 AM | File folder |
| Preliminary | 10/30/2024 9:49 AM | File folder |
| main.do | 10/28/2024 12:57 PM | DO File |

| Name | Date modified | Type |
|------|---------------|------|
| _ | 10/30/2024 9:49 AM | File folder |
| e | 10/30/2024 9:49 AM | File folder |
| f | 10/30/2024 9:49 AM | File folder |
| i | 10/30/2024 9:49 AM | File folder |
| o | 10/30/2024 9:49 AM | File folder |
| r | 10/30/2024 9:49 AM | File folder |
| s | 10/30/2024 9:49 AM | File folder |
| backup.trk | 10/25/2024 1:14 PM | TRK File |
| stata.trk | 10/28/2024 12:57 PM | TRK File |

# 3 – Manage dependencies simply: repado

# 4 – Beautify your Stata code: `lint`

## The problem:

- Messy code is hard to understand and error-prone

- More importantly: it's not transparent at all!

## Solution:

- `lint`: a code linter for Stata

- It scans a do-file and flags coding styles issues, according to the DIME Analytics Stata Style Guide

# 4 – Beautify your Stata code: lint


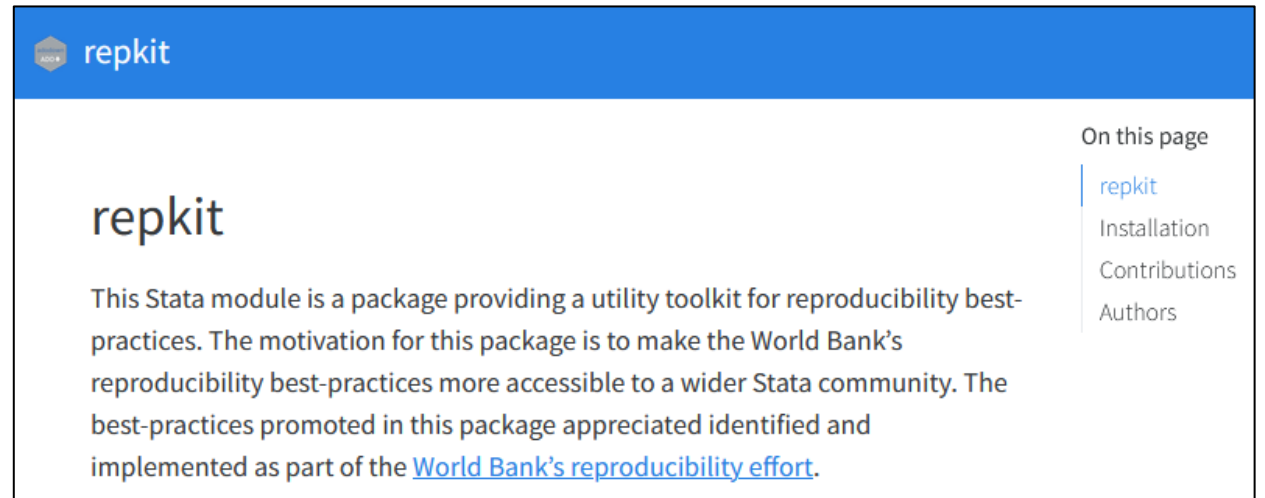
Command

```
lint "test/bad.do", verbose
```

# 4 – Beautify your Stata code: `lint`

**"Programs must be written for people to read, and only incidentally for machines to execute."**

*—Abelson, Susman, and Susman, [Structure and Interpretation of Computer Programs](#) (1985)*

# **Remember:** `repkit`

- [Package site](#)

- [GitHub repository](#)

- Comments or bugs? [Here](#)

- Currently available in SSC:

    `ssc install repkit`

# Thank you!

lsanmartin@worldbank.org