

BEAD: Reproducible Computational Research Made Simple

Miklós Koren (Central European University) Krisztián Fekete (Central European University)

RSEcon, September 10, 2025¹

¹This project has received funding from the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation programme (grant agreement No 313164). The views expressed are those of the authors and do not necessarily reflect those of the ERC or the European Commission.

The Editor Says You Have One Week

- Journal editor: “substantial revision invited”
- Reviewers liked Figure 1 (life expectancy vs GDP per capita)
- Concern about health data source
- You need to:
 - Address reviewer concerns
 - Redo analysis with new data
 - Recreate Figure 1
 - Submit within one week

But Your Submission is Months Old

- Research submitted months ago
- Team has been improving data cleaning since then
- Different statistical methods now
- **First question:** How did I actually produce Figure 1?

Research Results are Functions

$$\text{Figure 1} = f(\text{code}, \text{data})$$

- Results depend on both algorithms and data
- Code under version control (Git) -> Yes
- Tagged commit at submission -> Yes
- **But what about the data?**

Data is Also a Function

$$\text{data}_1 = f(\text{code}_2, \text{data}_2)$$

- Data produced by wrangling/cleaning steps
- Which countries dropped?
- What transformations applied?
- Feature engineering details?
- **Chain of data provenance**

Real-World Data Pipelines

- Multiple datasets merged
- Many cleaning steps
- Different versions coexisting
- Green = using latest version
- Red/yellow = outdated dependencies
- Complex dependency graph

The Data Provenance Problem

Why it's complex:

- 1 **Frequent changes:** Code and data both evolve
- 2 **Complex pipelines:** Many steps, multiple datasets
- 3 **Tool heterogeneity:** Python, R, SQL, DuckDB all in one project

Team Dynamics Make it Worse

- Master/PhD students graduate and leave
- Different team members use different tools
- **Every meeting starts with:**
 - “Who knows how to reproduce this?”
 - “Who has the data?”
 - “That person already left. . .”

Existing Solutions

Version Control (Git)

- Great for code
- Not suitable for large binary data

Data Version Control (DVC)

- Similar spirit to BEAD
- More complex than needed
- dvc.org

Orchestration Tools

- Apache Airflow (Python) - airflow.apache.org
- dbt (SQL) - getdbt.com
- KNIME (no-code) - knime.com
- Too complex for heterogeneous teams

Enter BEAD

A command-line tool that ensures your output is a function of your input

- Much simpler than alternatives
- Language agnostic
- Works with heterogeneous teams
- Different experience levels
- Different operating systems

What BEAD Does NOT Do

Not a code runner

- You run your own code
- Python, R, Stata, SQL - doesn't matter

Not a file delivery system

- File system stores your files
- You copy/move files yourself

Only requirement:

- Works with flat files on file system
- Files not too big (20GB works fine)

What BEAD Enforces

Input data is immutable

- Cannot modify raw data
- Forces good practices
- Preserves data lineage

Core BEAD Concepts

The BEAD

- Self-contained computational unit
- Contains code, data, results
- Packaged as ZIP file
- Remembers exact provenance

Simple Commands

```
bead new my-analysis
```

```
bead input add source-data
```

```
bead save results
```

Demo Time

Christian will now demonstrate BEAD in action...

How BEAD Solves Our Problems

Problem	BEAD Solution
"What data did we use?"	Every bead remembers exact version
"It worked on my machine"	Exact same setup for everyone
"That person left"	Work stays reproducible
Team uses different tools	Language agnostic
Complex pipelines	Chain beads together

Real Research Example

- Multiple datasets connected
- Many cleaning steps
- Green = using latest data version
- Some steps outdated
- BEAD tracks entire dependency graph

BEAD in Practice

Step 1: Create workspace

```
bead new health-analysis
```

Step 2: Load inputs

```
bead input add wdi-data
```

```
bead input add health-metrics
```

Step 3: Run analysis

```
python clean_data.py
```

```
R --file=analyze.R
```

Step 4: Save snapshot

```
bead save figure1-v2
```

Why BEAD is Different

- **Simple:** 4 commands to learn
- **Universal:** Any language, any tool
- **Portable:** Just ZIP files
- **Secure:** Data stays on your servers
- **Transparent:** Open source, no vendor lock-in

For Research Software Engineers

- Minimal learning curve for researchers
- No infrastructure requirements
- Works with existing workflows
- Complements version control
- Enables true reproducibility

Get Started

Installation

```
pip install bead
```

Documentation

codedthinking.github.io/bead.zip

Source Code

github.com/codedthinking/bead.zip

Key Takeaways

- 1 **Data provenance is hard** - especially with changing teams
- 2 **Existing tools too complex** - for heterogeneous research teams
- 3 **BEAD keeps it simple** - focuses on one thing well
- 4 **Reproducibility becomes automatic** - not an afterthought

Thank You!

Questions?

Contact

- Web: bead.zip
- GitHub: github.com/codedthinking/bead.zip



Funded by
the European Union



European Research Council
Established by the European Commission

References

- **World Development Indicators:** data.worldbank.org/indicator
- **DVC (Data Version Control):** dvc.org
- **Apache Airflow:** airflow.apache.org
- **dbt:** getdbt.com
- **KNIME:** knime.com