# Analyzing R

**On the Anatomy of Real-World R Code for Static Analysis**

At SE '25     Ulm University | MSR '24 | **Florian Sihler**, Pietzschmann, Straub, Tichy, Diera, Dahou | February, 2025

SP | Software Engineering
Programming Languages

universität uulm

# The R Programming Language

**Is Used**
in Research[1]

**Is Used**
for Statistical Computing[2]

**Is Used (mostly)**
by Non-Programmers

- ≈70 % of scripts are not reproducible[1]

- Lacks sophisticated static-analysis tools

- *Many* powerful reflective capabilities[3]

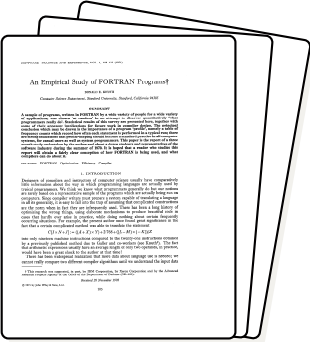- Incomplete language specification[4]

Which features are actually used?

[1] Trisovic et al., "A Large-Scale Study on Research Code Quality and Execution" (*Sci Data '22*)
[2] https://cran.r-project.org/
[3] Flückiger et al., "R melts brains: an IR for first-class environments and lazy effectful arguments" (*DLS '19*)
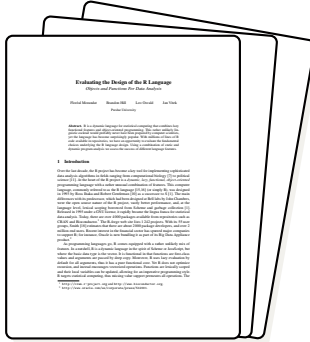[4] R Core Team, *R Language Definition*
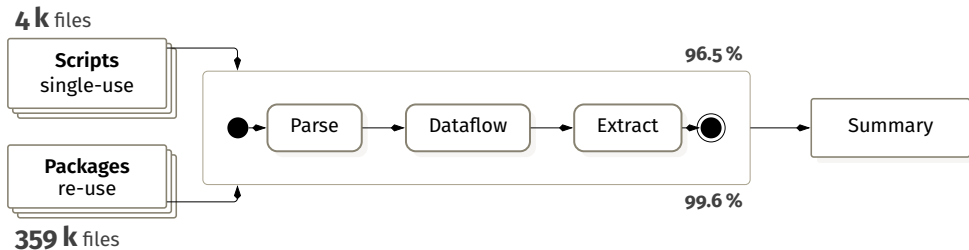
# Related Research



[5] Optimization
Fortran

[6] Syntactic Rules
Java

[7] Dynamic Usage
R

[7] Morandat et al., "Evaluating the design of the R language: Objects and functions for data analysis" (*ECOOP '12*)
[6] Qiu et al., "Understanding the Syntactic Rule Usage in Java" (*JSS '17*)
[5] Knuth, "An Empirical Study of FORTRAN Programs" (*Software: Practice and Experience '71*)

# Extraction Workflow



RQ 1: Frequent Features

RQ 2: Differences in Research Scripts and Packages

RQ 3: Insights for Static Analysis

# Overview

4.1
**Processing Errors**

4.5
**Loops**

4.2
**Metadata**

4.6
**Function Definitions**

4.3
**Assignments and Access**

4.7
**Function Calls**

4.4
**Conditionals**
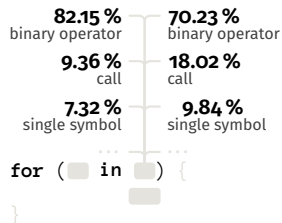
4.8
**Packages**

# Assignments   [4.3]

- ←, ←, =, →, →, `assign`, `delayedAssign`, ...

- 41 % of scripts mix ← and =

- Assignment functions are rare, but more common in packages

# Loops  [4.5]

```r
for (i in 1:10) {
    print(i)
}
```

R-Scripts



| 82.15 %<br>binary operator | 70.23 %<br>binary operator |
| 9.36 %<br>call | 18.02 %<br>call |
| 7.32 %<br>single symbol | 9.84 %<br>single symbol |

```r
for (  in  ) {
}
```

R-Packages

- Most loops have a for-i form

- Scripts contain on average 3 times as many loops

# Meta-Programming   [4.7.1]

### Evaluate Strings

```
eval(parse(text=
  paste0("v",1," ← 42")
)) # v1 ← 42
```

### Modify Functions

```
f ← function(a, b) a
body(f) ← quote(b)
f(1, 2) # 2
```

### Redefine "Keywords"

```
'for' ← \(...) "hi"
for(i in 1:10) x ← i
# "hi"
```

### Store/Load Environment

```
save.image(file="env")
# ...
load("env")
```

# Meta-Programming   [4.7.1]

### Evaluate Strings

```
eval(parse(text=
  paste0("v",1, "← 42")
)) # v1
```
1 % of scripts

3 % of packages

### Modify Functions

```
f ← function(a, b) a
body(f)
f(1, 2) # 2
```
Effectively unused

### Redefine "Keywords"

```
`for` ← \(...) "hi"
for(i i      i 2)      i
# "hi"
```
Effectively unused

### Store/Load Environment

```
save.image(file="env")
# ...
load("e
```
12 % of scripts

0.8 % of packages

# Study Results

**RQ1**
**Frequent Features**

+ Only 2 of all assignment operators
− Reflective functions
− No tests/checks in scripts

**RQ2**
**Differences**

• Scripts are longer
• Scripts prefer (for-)loops

**RQ3**
**Insights**

• Extensions for {lintr}
• No focus on reflective functions required
• Error-tolerant parsing

R has many features

only a few are used frequently

# Appendix

# References I

[1]  Ana Trisovic et al. "A Large-Scale Study on Research Code Quality and Execution". 2022

[2]  *The Comprehensive R Archive Network — cran.r-project.org.* 2024

[3]  Olivier Flückiger et al. "R melts brains: an IR for first-class environments and lazy effectful arguments". 2019

[4]  R Core Team. *R Language Definition.* 2024

[5]  Donald E. Knuth. "An Empirical Study of FORTRAN Programs". 1971

[6]  Dong Qiu et al. "Understanding the Syntactic Rule Usage in Java". Jan. 1, 2017

[7]  Floréal Morandat et al. "Evaluating the design of the R language: Objects and functions for data analysis". 2012