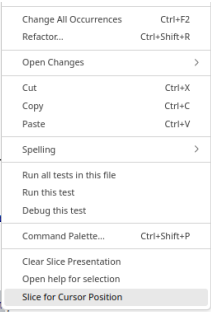


```
1 sum ← 0
2 product ← 1
3 w ← 7
4 N ← 10
5
6 for (i in 1:(N-1)) {
7   sum ← sum +
8   product ← p
9 }
10
11 cat("Sum:", sum,
12 cat("Product:", product, "\n")
13
```



A context menu is open over the code, showing various editing and development actions. The menu items include: 'Change All Occurrences' (Ctrl+F2), 'Refactor...' (Ctrl+Shift+R), 'Open Changes' (with a right arrow), 'Cut' (Ctrl+X), 'Copy' (Ctrl+C), 'Paste' (Ctrl+V), 'Spelling' (with a right arrow), 'Run all tests in this file', 'Run this test', 'Debug this test', 'Command Palette...' (Ctrl+Shift+P), 'Clear Slice Presentation', 'Open help for selection', and 'Slice for Cursor Position' (which is highlighted).



```
1 sum ← 0
2 product ← 1
3 w ← 7
4 N ← 10
5
6 for (i in 1:(N-1)) {
7   sum ← sum + i + w
8   product ← product * i
9 }
10
11 cat("Sum:", sum, "\n")
12 cat("Product:", product, "\n")
13
```

flowR: A Program Slicer for the R Programming Language

deRSE '24 | Ulm University | Florian Sihler and Prof. Matthias Tichy | March 6, 2024

The R Programming Language

- R is mainly designed for statistical computing^[1]
 - Heavily used in research (e.g., social science)^[2]
 - Ranks 6th on PYPL^[3]
 - Over 20 000 packages on CRAN^[1]
- Most users are from non-computer-science domains
- Several problems in practice^[2, 4]
 - Replication
 - Program comprehension
 - Missing tool support
- Demand for software engineering practices^[5]

[5] Thimbleby, “Improving Science That Uses Code” (2023, Oxford University Press)

[4] Wonsil et al., “Reproducibility as a Service” (2023, Software: Practice and Experience)

[3] <https://pypl.github.io/> [archived]

[2] Trisovic et al., “A Large-Scale Study on Research Code Quality and Execution” (2022, Nature Publishing Group)

[1] <https://cran.r-project.org/>

R Scripts

We analyzed 4 083 R-files^[7]



```
## set the data directory and load workspace
setwd("G:/Shared drives/Fenologija/Raksti/abeles/Zenodo/")
load("Apple_Pure_phenology_R_workspace_image.RData")

# -----
# script to recreate components included in of "Apple_Pure_phenology_R_workspace_image.RData"
# -----
# phenology data subset having at least 14 observations overlapping with meteorology data set
{
  dM <- d %>%
    filter(Year %in% (meteoH$e_obs$Year))
  dM <- dM %>%
    group_by(Variety) %>%
    summarize(N = n()) %>%
    arrange(N) %>%
    filter(N >= 14) %>% # 12 skirnes
    select(-N) %>%
    left_join(dM)
}
```

```
# set the data directory and load workspace
setwd("G:/Shared drives/Fenologija/Raksti/abeles/Zenodo/")
load("Apple_Pure_phenology_R_workspace_image.RData")

# -----
# script to recreate components included in of "Apple_Pure_phenology_R_workspace_image.RData"
# -----
# phenology data subset having at least 14 observations overlapping with meteorology data set
{
  dM <- d %>%
    filter(Year %in% (meteoH$e_obs$Year))
  dM <- dM %>%
    group_by(Variety) %>%
    summarize(N = n()) %>%
    arrange(N) %>%
    filter(N >= 14) %>% # 12 skirnes
    select(-N) %>%
    left_join(dM)
}

# -----
# FUNCTION definitions
# phenology model defined as function, see Kalvans et al, 2015, DDcos model for details
```

[7] Sihler et al., "On the Anatomy of Real-World R Code for Static Analysis" (2024, MSR)

[6] Druzde et al., *Apple phenology data set and R script, related to publication "Full flowering phenology of apple tree (Malus domestica) in Pūre orchard, Latvia from 1959 to 2019"* (2021, Zenodo)

R Scripts Fail to Replicate

74 % even fail to complete!^[2]



[2] Trisovic et al., "A Large-Scale Study on Research Code Quality and Execution" (2022, Nature Publishing Group)

[6] Drudze et al., *Apple phenology data set and R script, related to publication "Full flowering phenology of apple tree (Malus domestica) in Püre orchard, Latvia from 1959 to 2019"* (2021, Zenodo)

R Scripts Do Too Much



- Several analyses in one script
- Hard to comprehend
- Hard to extract/re-use parts

[6] Drudze et al., *Apple phenology data set and R script, related to publication "Full flowering phenology of apple tree (Malus domestica) in Pūre orchard, Latvia from 1959 to 2019"* (2021, Zenodo)

R Scripts Are Hard to Analyze

```
[8] mod.proj_wc <- ecospat.ESM.Projection(ESM.modeling.output=mod,  
    new.env=eval(parse(text = paste("env_",period,"_wc",sep=""))))
```

► String-based code evaluation

```
[9] pull.cat <- function(x) {  
  bins <- up_bins # (e.g., 6)  
  increments <- (range(x)[2] - range(x)[1])/(bins - 1)  
  to_return <- seq(range(x)[1], range(x)[2], increments)  
  return(to_return)  
}  
  
up.cat <- function(new_bins) {  
  up_bins = new_bins  
  body(pull.cat)[[2]] <- substitute(bins <- up_bins)  
}
```

► Self-modifying code

[9] Robertson, *Social hierarchy reveals thermoregulatory trade-offs in response to repeated stressors* (2020, Zenodo) [L. 68ff]

[8] Ma et al., *Predicting range shifts of pikas (Mammalia, Ochotonidae) in China under scenarios incorporating land-use change, climate change, and dispersal limitations* (2021, Zenodo) [L. 135f]

R Misses Sophisticated Analysis Tools

- RStudio IDE posit.co
 - Syntax-highlighting and auto-completion
 - Refactorings (rename, extract functions and variables) ← Often wrong (simple heuristics)
- R language server github.com/REditorSupport
 - Syntax-highlighting and auto-completion
 - Reference tracing & Refactorings (rename) ← Often wrong (XPath-Expressions)
- {lintr} github.com/r-lib/lintr
 - Style & syntax errors
 - Potential semantic errors
 - ↑
XPath-Expressions, packages
- {CodeDepends} github.com/duncantl/CodeDepends
 - Dependency analysis ← Only top scope
 - Creation of call-graphs

R Scripts ...

1. fail to replicate
2. **do too much**
3. are hard to analyze
4. are not well supported by tools

Better Software, Better Research

The Goal of flowR



```
## FlowR: A framework for the analysis of flow cytometry data
## Copyright (c) 2019-2021 flowR developers
## License: GPL (>= 3)
##
## This program is free software: you can redistribute it and/or modify
## it under the terms of the GNU General Public License as published by
## the Free Software Foundation, either version 3 of the License, or
## (at your option) any later version.
##
## This program is distributed in the hope that it will be useful,
## but WITHOUT ANY WARRANTY; without even the implied warranty of
## MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the
## GNU General Public License for more details.
##
## You should have received a copy of the GNU General Public License
## along with this program. If not, see <http://www.gnu.org/licenses/>.
##
## The flowR package is a framework for the analysis of flow cytometry
## data. It provides a set of functions to load, preprocess, and analyze
## flow cytometry data. The package is designed to be flexible and
## extensible, allowing users to integrate their own analysis functions
## and workflows.
##
## The flowR package is available on CRAN and can be installed using
## the following command:
##
## install.packages("flowR")
##
## For more information, please visit the flowR website:
## <http://flowR.r-lib.org>
##
## The flowR package is a part of the R ecosystem and is licensed
## under the GNU General Public License (GPL). It is a free software
## and you can redistribute it and/or modify it under the terms of
## the GPL.
##
## The flowR package is a part of the R ecosystem and is licensed
## under the GNU General Public License (GPL). It is a free software
## and you can redistribute it and/or modify it under the terms of
## the GPL.
```

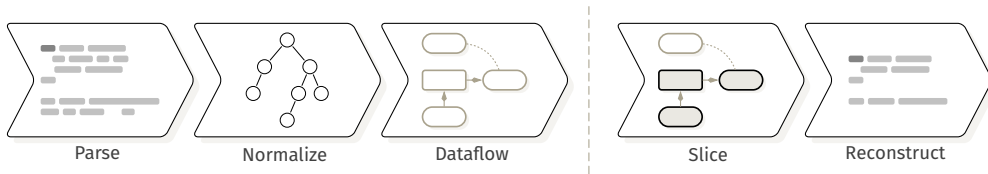


```
## FlowR: A framework for the analysis of flow cytometry data
## Copyright (c) 2019-2021 flowR developers
## License: GPL (>= 3)
##
## This program is free software: you can redistribute it and/or modify
## it under the terms of the GNU General Public License as published by
## the Free Software Foundation, either version 3 of the License, or
## (at your option) any later version.
##
## This program is distributed in the hope that it will be useful,
## but WITHOUT ANY WARRANTY; without even the implied warranty of
## MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the
## GNU General Public License for more details.
##
## You should have received a copy of the GNU General Public License
## along with this program. If not, see <http://www.gnu.org/licenses/>.
##
## The flowR package is a framework for the analysis of flow cytometry
## data. It provides a set of functions to load, preprocess, and analyze
## flow cytometry data. The package is designed to be flexible and
## extensible, allowing users to integrate their own analysis functions
## and workflows.
##
## The flowR package is available on CRAN and can be installed using
## the following command:
##
## install.packages("flowR")
##
## For more information, please visit the flowR website:
## <http://flowR.r-lib.org>
##
## The flowR package is a part of the R ecosystem and is licensed
## under the GNU General Public License (GPL). It is a free software
## and you can redistribute it and/or modify it under the terms of
## the GPL.
##
## The flowR package is a part of the R ecosystem and is licensed
## under the GNU General Public License (GPL). It is a free software
## and you can redistribute it and/or modify it under the terms of
## the GPL.
```

- Interested in a single figure
- $\approx 70\%$ reduction

[6] Drudze et al., *Apple phenology data set and R script, related to publication "Full flowering phenology of apple tree (Malus domestica) in Pūre orchard, Latvia from 1959 to 2019"* (2021, Zenodo)

The Architecture




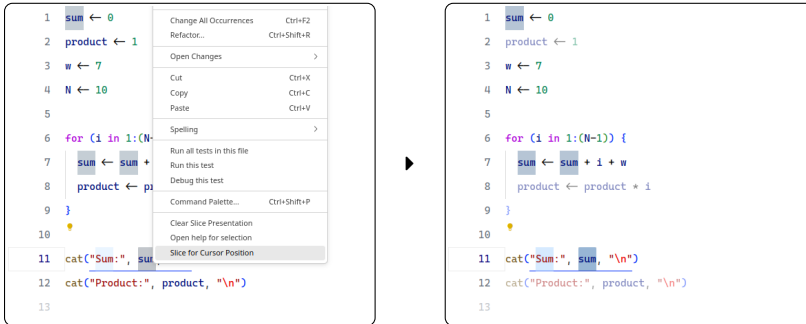
[10] Sihler, "Constructing a static program slicer for R programs" (2023, Ulm University)

[11] Weiser, "Program Slicing" (1984, IEEE Transactions on Software Engineering)





Visual Studio Code Integration

Rudimentary VSCode integration at:

 github.com/Code-Inspect/vscode-flowr



Using flowR

-  github.com/Code-Inspect/flowr
-  github.com/Code-Inspect/vscode-flowr
-  hub.docker.com/r/eagleoutice/flowr
-  npmjs.com/package/@eagleoutice/flowr

Server

```
{
  "type": "request-file-analysis",
  "id": "1",
  "filetoken": "123",
  "content": "x ← 1; x * y"
}

{
  "type": "request-slice",
  "id": "2",
  "filetoken": "123",
  "criterion": ["1@x"]
}
```

Library

```
const s = new SteppingSlicer({
  shell, tokenMap,
  request:
    requestFromInput("x ← 1; x * y"),
  criterion: ['1@x'],
})
const slice =
  await s.allRemainingSteps()
```

REPL

```
R> :parse "x ← 1; x * y"
exprlist
└─ expr
  │ └─ expr
  │ │ └─ SYMBOL "x" (1:1)
  [...]

R> :dataflow* "x ← 1; x * y"
https://mermaid.live/edit#base64:eyJj...
```

```
docker run -it --rm eagleoutice/flowr
```

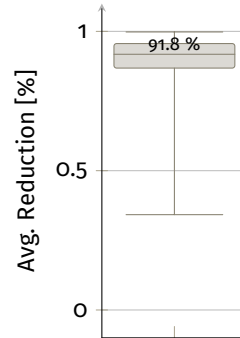
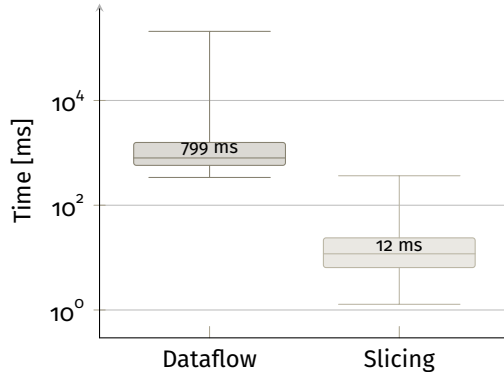
Appendix

The R Code Static Analysis Landscape

	goal	method	impl. lang.	op. assignments	func. assignments	value trace (a.i., ...)	controlflow	non-std. eval.	special operators	function calls	libraries	quotation	reflection	side effects	static scope	dynamic scope	type inference	pointer analysis	external files	pre-processors	hooks	FFI
[12] {CodeDepends}	static analysis	AST visitor	R	●	○		○	○	○	○			○	●								
[13] {codetools}	static analysis	AST visitor	R	●	○	○	○	○	○		○			●								
[14] {checkglobals}	missing libs.	AST visitor	R, C	●	○		○	○	○	○	○			○							○	
[15] {rstatic}	static analysis	AST visitor	R	●		○	●		○					●								
[16] {CodeAnalysis}	static analysis	AST visitor	R	●	○	○	●	○	○	○	○	○		●		○		○				
[17] {RTypeInference}	type inference	AST visitor	R	●		○	●		○					●		●						
[18] {pkgstats}	package insight	ctags & gtags	R, C++	○					○	○				○								
[19] {globals}	distributed env.	AST visitor	R	●		○	○	○	○	○		○		●								
[20] {Rclean}	debug/refactor	PDG traversal	R	●	○		○	○	○	○	○			●								
[21] {lintr}	linting	XPath, visitor	R	●	○		○		○	○				●								
[22] {Similar}	plagiarism	PDG, visitor	R, C++	○			○		○					●								
[23] {rco}	optimization	AST visitor	R	●		●	○							●								
[24] {cyclocomp}	code complexity	AST visitor	R				○							●								
[25] {flow}	visualize, debug	AST visitor	R, C	○			●		○		○	○							○			
[26] {PaRe}	code review	Regex	R	○			○		○													
[27] {dfgraph}	static analysis	AST visitor	R	○			○		○													
[28] {rflowgraph}	call graph	AST visitor	R						○	○												
[29] {languageserver}	editor support	XPath, visitor	R, C	●	○				○	○				●								
[30] RStudio	editor support	AST visitor	Java, C++, TS, ...	●	○				○	○				●								
[31] ROSA	optimization	visitor	C++, R	●			○				○		○	●		●	○					
[32] Random	abstract-int.	trace & visitor	R	○	○	○	○			○				○								
[4] RaaS	reproducibility	AST visitor	Python, R	●	○			○	○	○	○			●								
[33] GNU R	execute R	bytecode	C, Fortran, R	●	○	●	○	○	●	○		○		●								
[34] FastR	execute R	AST visitor	Java, C, R, ...	●		○			○	○	○			●								
[35] R̃	execute R	SSA, bytecode	C++, R, C, ...	●	●	○	●	○	●	○		○		●		○		●				
[36] renjin	execute R	SSA, CFG	R, Java, C, ...	●	○	○	●		○	○			○	●				●				
[37] pgR	execute R	bytecode	C, R, ...	●	○	●	○		○	○		○		●								
[38] MRO	execute R	bytecode	R, C, Fortran, ...	●	○	○	○															
[39] RCC	transpile C	CFG, bytecode	C/C++, Fortran, R, ...	●		○	○		○				○	●								

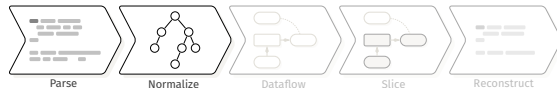
Performance Measurements

- We generated *every possible* variable of interest
- Dataflow results can be cached

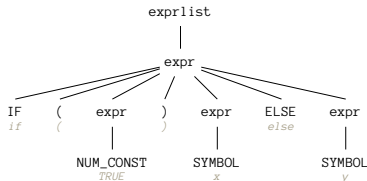


99th percentile

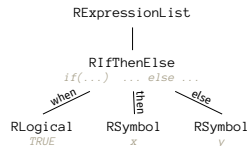
Parse & Normalize



```
parse(text="if(TRUE) x else y")
```



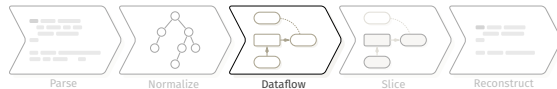
normalized →



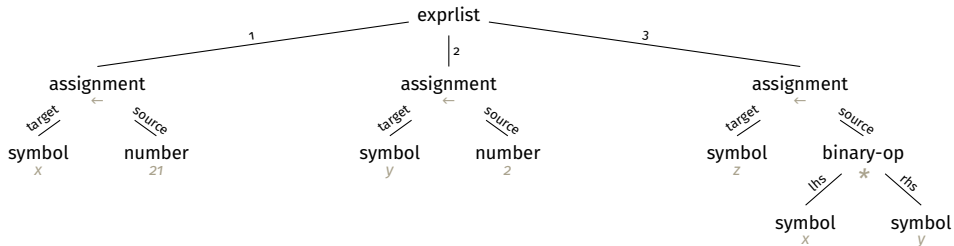
- Normalizing constants, namespacing, operators, ...
- We use the “R language definition”^[40] as a basis

[40] R Core Team, *R Language Definition* (2023)

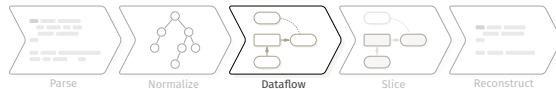
Dataflow



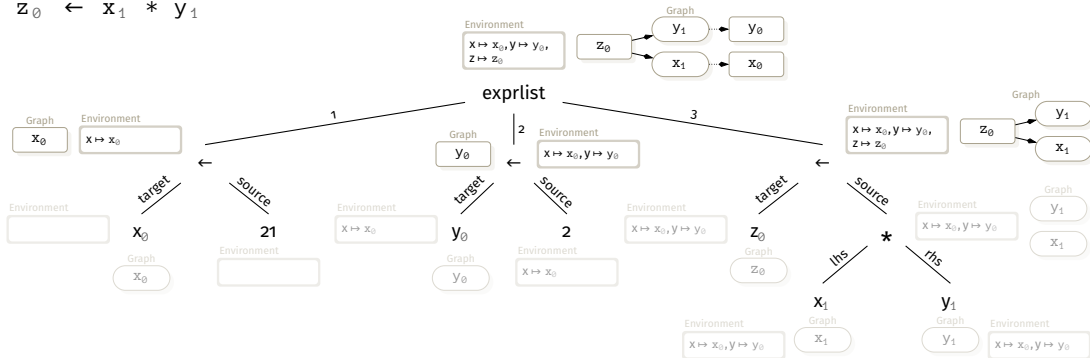
$x_0 \leftarrow 21$
 $y_0 \leftarrow 2$
 $z_0 \leftarrow x_1 * y_1$



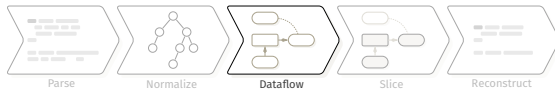
Dataflow



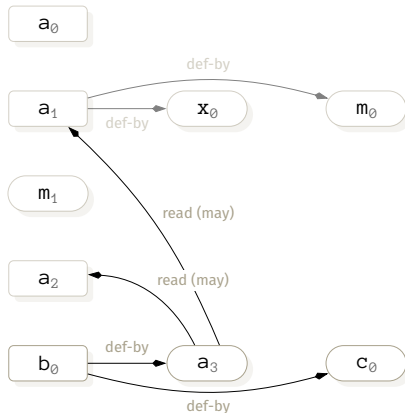
$x_0 \leftarrow 21$
 $y_0 \leftarrow 2$
 $z_0 \leftarrow x_1 * y_1$



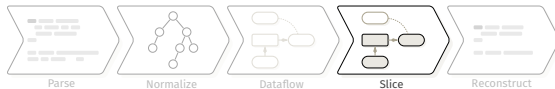
Resulting Dataflow



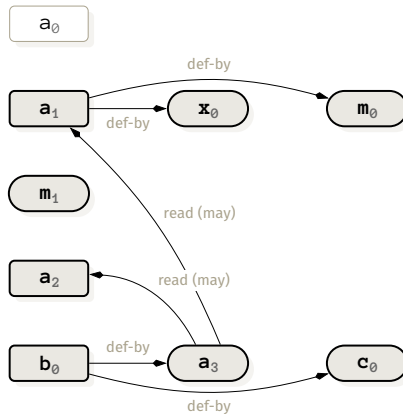
```
> a0 ← 3  
a1 ← x0 * m0  
  
if(m1 > 3) {  
  a2 ← 5  
}  
  
b0 ← a3 + c0
```



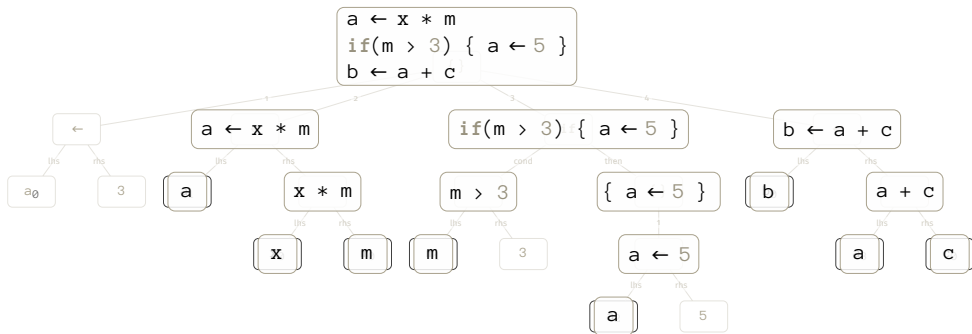
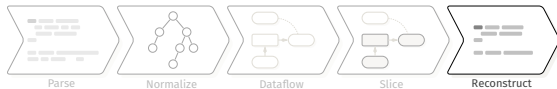
Slicing, I



```
a0 ← 3  
a1 ← x0 * m0  
  
if(m1 > 3) {  
  a2 ← 5  
}  
  
b0 ← a3 + c0
```

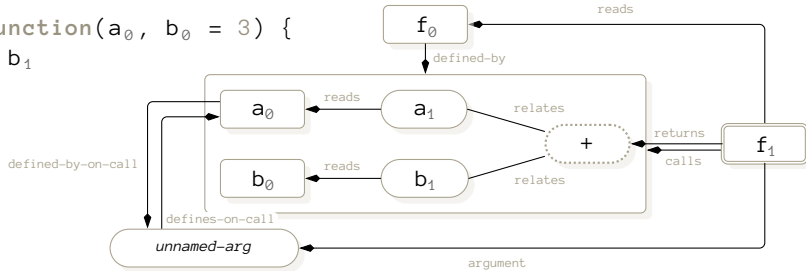


Slicing, II



There Is More...

```
f0 ← function(a0, b0 = 3) {  
  a1 + b1  
}  
f1(39)
```



Definition-Retrieval

```
paste(  
  "(*|descendant-or-self::exprlist/*)[self::FUNCTION or self::OP-LAMBDA]/  
    following-sibling::SYMBOL_FORMALS[text() = '{token_quote}' and @line1 <= {  
      row}]",  
  "(*|descendant-or-self::exprlist/*)[LEFT_ASSIGN[preceding-sibling::expr[count  
    (*)=1]/SYMBOL[text() = '{token_quote}' and @line1 <= {row}] and following-  
    sibling::expr[@start > {start} or @end < {end}]]]",  
  "(*|descendant-or-self::exprlist/*)[RIGHT_ASSIGN[following-sibling::expr[count  
    (*)=1]/SYMBOL[text() = '{token_quote}' and @line1 <= {row}] and preceding-  
    sibling::expr[@start > {start} or @end < {end}]]]",  
  "(*|descendant-or-self::exprlist/*)[EQ_ASSIGN[preceding-sibling::expr[count(*)=  
    1]/SYMBOL[text() = '{token_quote}' and @line1 <= {row}] and following-  
    sibling::expr[@start > {start} or @end < {end}]]]",  
  "forcond/SYMBOL[text() = '{token_quote}' and @line1 <= {row}]",  
  sep = "|")
```

References I

- [1] *The Comprehensive R Archive Network — cran.r-project.org.* 2024
- [2] Ana Trisovic et al. “A Large-Scale Study on Research Code Quality and Execution”. 2022
- [3] *PYPL — Popularity of Programming Language index.* 2024
- [4] Joseph Wonsil et al. “Reproducibility as a Service”. 2023
- [5] Harold Thimbleby. “Improving Science That Uses Code”. 2023
- [6] Inese Drudze et al. *Apple phenology data set and R script, related to publication "Full flowering phenology of apple tree (Malus domestica) in Pūre orchard, Latvia from 1959 to 2019".* June 2021
- [7] Florian Sihler et al. “On the Anatomy of Real-World R Code for Static Analysis”. 2024
- [8] Liang Ma et al. *Predicting range shifts of pikas (Mammalia, Ochotonidae) in China under scenarios incorporating land-use change, climate change, and dispersal limitations.* Aug. 2021
- [9] Joshua Robertson. *Social hierarchy reveals thermoregulatory trade-offs in response to repeated stressors.* Oct. 2020
- [10] Florian Sihler. “Constructing a static program slicer for R programs”. 2023

References II

- [11] Mark Weiser. “Program Slicing”. July 1984
- [12] Duncan Lang et al. *CodeDepends. Analysis of R Code for Reproducible Research and Code Comprehension.* 2018
- [13] Luke Tierney. *codetools: Code Analysis Tools for R.* 2023
- [14] Joris Chau. *checkglobals: Static Analysis of R-Code Dependencies.* 2023
- [15] Nick Ulle and Duncan Temple Lang. *rstatic: Low-level Static Analysis Tools for R Code.* 2019
- [16] Duncan Lang et al. *CodeAnalysis. Tools for static analysis of R code.* 2023
- [17] Nick Ulle and Duncan Temple Lang. *RTypeInference: Infer Types of Inputs and Outputs for R Expressions.* 2021
- [18] Mark Padgham. *pkgstats.* 2021
- [19] Henrik Bengtsson. *globals: Identify Global Objects in R Expressions.* 2022
- [20] Matthew Lau. *Rclean: A Tool for Writing Cleaner, More Transparent Code.* 2022
- [21] Jim Hester et al. *lintr: A ‘Linter’ for R Code.* 2023
- [22] Maciej Bartoszek and Marek Gagolewski. *SimilaR: R Source Code Similarity Evaluation.* 2020

References III

- [23] Juan Cruz Rodriguez. *rco: The R Code Optimizer*. 2021
- [24] Gabor Csardi. *cyclocomp: Cyclomatic Complexity of R Code*. 2023
- [25] Antoine Fabri. *flow: View and Browse Code Using Flow Diagrams*. 2023
- [26] Maarten van Kessel. *PaRe: A Way to Perform Code Review or QA on Other Packages*. 2023
- [27] Dan Kary. *dfgraph: Visualize R Code with Data Flow Graphs*.
- [28] Evan Patterson. *The algebra and machine representation of statistical models*. 2020
- [29] Randy Lai. *languageserver: Language Server Protocol*. 2023
- [30] Posit team. *RStudio: Integrated Development Environment for R*. 2023
- [31] Rathijit Sen et al. *ROSA: R Optimizations with Static Analysis*. 2017
- [32] Gianluca Amato and Francesca Scozzari. “Random: R-Based Analyzer for Numerical Domains”. 2012
- [33] R Core Team. *R: A Language and Environment for Statistical Computing*. 2023
- [34] Tomas Kalibera et al. “A fast abstract syntax tree interpreter for R”. 2014
- [35] Olivier Flückiger et al. “Sampling optimized code for type feedback”. 2020

References IV

- [36] Alexander Bertram. “Renjin: A new r interpreter built on the jvm”. 2013
- [37] Radford M Neal. “Speed Improvements in pqR: Current Status and Future Plans”. 2014
- [38] *Microsoft R Open Source*. 2019
- [39] John Garvin. *RCC: A compiler for the R language for statistical computing*. 2004
- [40] R Core Team. *R Language Definition*. 2023