

Rscript: a Relational Approach to Program and System Understanding

Paul Klint



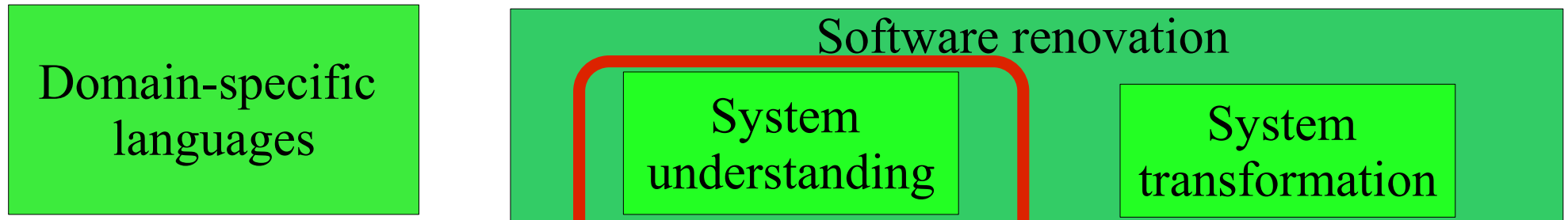
UNIVERSITEIT VAN AMSTERDAM

Structure of Presentation

- Background and context
- About program understanding
- Roadmap: Rscript

Background

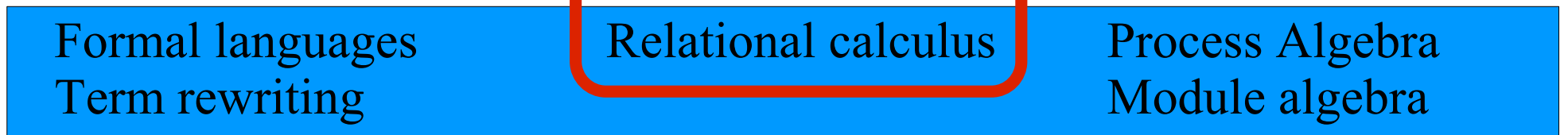
Application areas



Technology



Foundations



Compilation is a mature area

- Some new developments
 - just-in-time compilation
 - energy-aware code generation
- Many research results are not yet used widely
 - interprocedural pointer analysis
 - slicing
- Why don't we just apply all these techniques to understanding and restructuring?

Compilation is a mature area

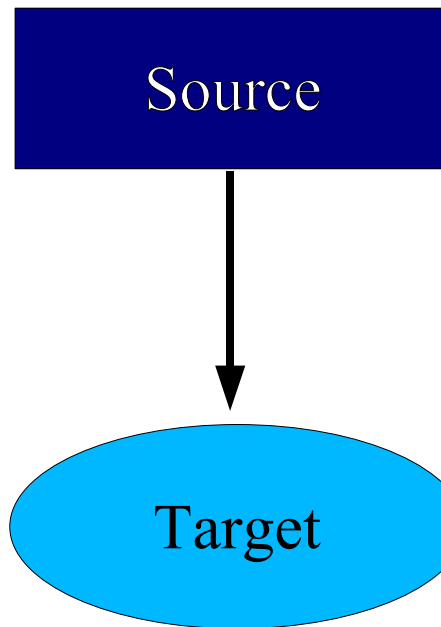
- ... of course, we do just that, but ...
- there is a **mismatch** between
 - standard compilation techniques and
 - the needs for understanding and restructuring

Compilation is ...

- A **well-defined process** with well-defined input, output and constraints
- **Input**: source program in a fixed language with well-defined syntax and semantics
- **Output**: a fixed target language with well-defined syntax and semantics
- **Constraints** are known (correctness, performance)
- A **batch-like** process

Compilation is ...

Single,
well defined,
source



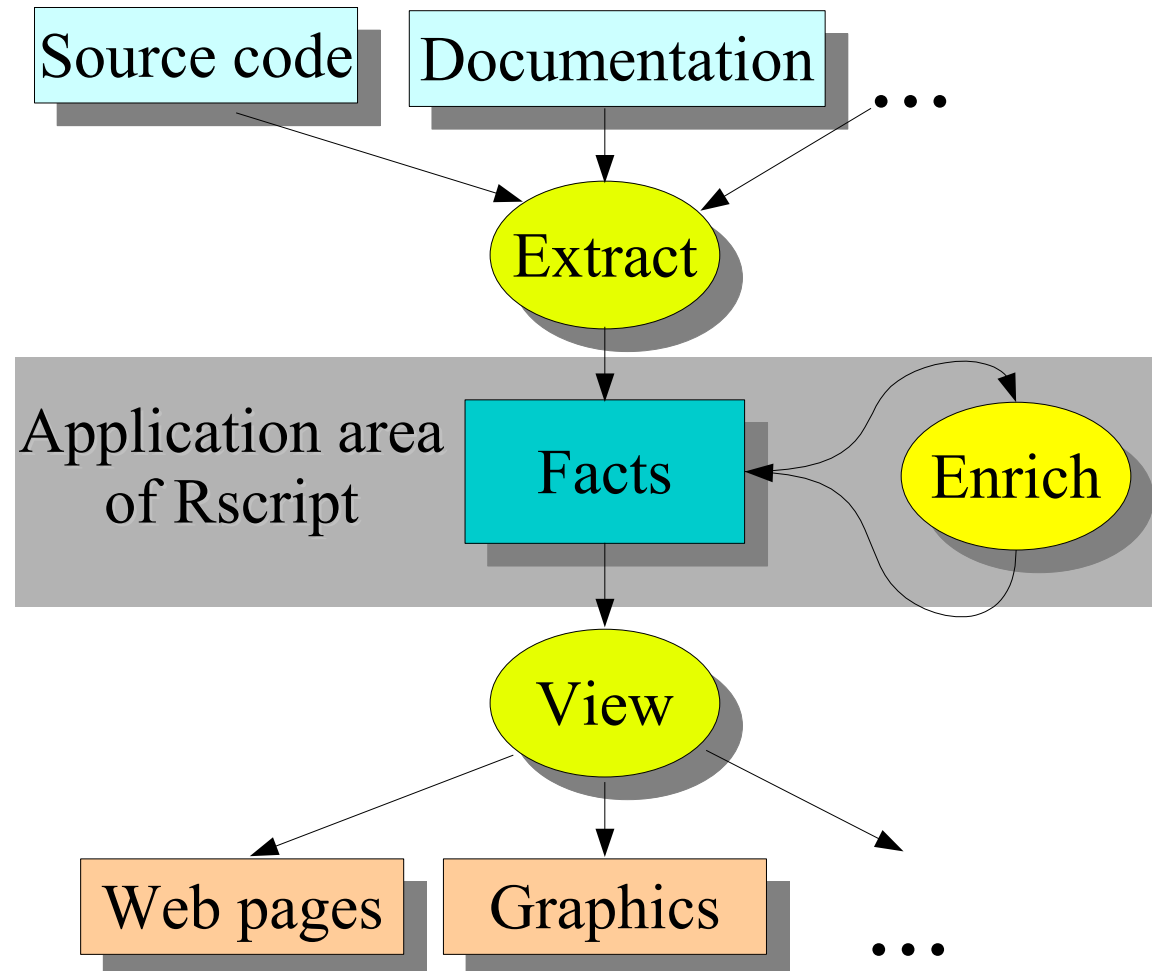
A batch-like process with
clear constraints

Single,
well
defined,
target

Understanding is ...

- An exploration process with as input
 - system artifacts (source, documentation, tests, ...)
 - implicit knowledge of its designers or maintainers
- There is no clear target language
- An interactive process:
 - **Extract** elementary facts
 - **Abstract** to get derived facts needed for analysis
 - **View** derived facts through visualization or browsing

Extract-Enrich-View Paradigm



Examples of understanding problems

- Which programs call each others?
- Which programs use which databases?
- If we change this database record, which programs are affected?
- Which programs are more complex than others?
- How much code clones exist in the code?

Examples of the results of understanding

- Textual reports indicating properties of system parts (complexity, use of certain utilities, ...)
- Same, but in hyperlinked format
- Graphs (call graphs, use def graphs for databases)
- More sophisticated visualizations

Other aspects of Understanding

- Systems consist of several source languages
- Analysis techniques over multiple language => a language-independent analysis framework is needed
- A very close link to the source text is needed

Related approaches

- Generic dataflow frameworks exist but are **not used widely**
- Relations have been used for querying of software (Rigi, GROK, RPA, ...)
 - All based on untyped, binary, relation algebra
 - Mostly used for architectural, **coarse grain**, queries

Relation-based analysis

- What happens if we use relations for fine grain software analysis (ex: find uninitialized variables)
- What happens if we use a relational calculus (as opposed to the relational algebra approaches)?
- What happens if we use term rewriting as basic computational mechanism?
 - relations can represent graphs in the rewriting world
- Could yield a unifying framework for analysis and transformation

Roadmap

- Rscript in a nutshell
- Example 1: call graph analysis
- Example 2: component structure
- Example 3: Java analysis
- Example 4: a toy language
- A vizualization experiment

Roadmap

- Rscript in a nutshell
- Example 1: call graph analysis
- Example 2: component structure
- Example 3: Java analysis
- Example 4: a toy language
- A vizualization experiment

Rscript in a Nutshell

- Basic types: **bool**, **int**, **str**, **loc** (text location in specific file with comparison operators)
- Sets, relations and associated operations (domain, range, inverse, projection, ...)
- Comprehensions
- User-defined types
- Fully typed
- **Functions** and **sets of equations** over the above

Rscript: examples

- Set: {3, 5, 3}
 - type: `set[int]`
- Set: {"y", "x", "z"}
 - type: `set[str]`
- Relation: {<"y", 3>, <"x", 3>, <"z", 5>}
 - type: `rel[str,int]`

Rscript: examples

- `rel[str,int]` $U = \{ \langle "y", 3 \rangle, \langle "x", 3 \rangle, \langle "z", 5 \rangle \}$
- `int` $Usize = \#U$
 - 3
- `rel[int,str]` $Uinv = inv(U)$
 - $\{ \langle 3, "y" \rangle, \langle 3, "x" \rangle, \langle 5, "z" \rangle \}$
- `set[str]` $Udom = domain(U)$
 - $\{ "y", "x", "z" \}$

domain:

all elements in lhs of pairs

range:

all elements in rhs of pairs

carrier:

all elements in lhs or rhs
of pairs

Comprehensions

- Comprehensions: $\{ \text{Exp} \mid \text{Gen1}, \text{Gen2}, \dots \}$
 - A generator is an enumerator or a test
 - Enumerators: $V : \text{SetExp}$ or $\langle V1, V2 \rangle : \text{RelExp}$
 - Tests: any predicate
 - consider all combinations of values in $\text{Gen1}, \text{Gen2}, \dots$
 - if some Gen_i is false, reject that combination
 - compute Exp for all legal combinations

Comprehensions

- $\{X \mid \text{int } X : \{1,2,3,4,5\}\}$
 - yields $\{1,2,3,4,5\}$
- $\{X \mid \text{int } X : \{1,2,3,4,5\}, X > 3\}$
 - yields $\{4,5\}$
- $\{\langle Y, X \rangle \mid \langle \text{int } X, \text{int } Y \rangle : \{\langle 1,10 \rangle, \langle 2,20 \rangle\}\}$
 - yields $\{\langle 10,1 \rangle, \langle 20,2 \rangle\}$

Functions

- $\text{rel}[\text{int}, \text{int}] \text{inv}(\text{rel}[\text{int}, \text{int}] R) =$
 $\{ \langle Y, X \rangle \mid \langle \text{int } X, \text{int } Y \rangle : R \}$
 - $\text{inv}(\{ \langle 1, 10 \rangle, \langle 2, 20 \rangle \})$ yields $\{ \langle 10, 1 \rangle, \langle 20, 2 \rangle \}$
- $\text{rel}[\&B, \&A] \text{inv}(\text{rel}[\&A, \&B] R) =$
 - $\{ \langle Y, X \rangle \mid \langle \&A X, \&B Y \rangle : R \}$
 - $\text{inv}(\{ \langle 1, \text{"a"} \rangle, \langle 2, \text{"b"} \rangle \})$ yields $\{ \langle \text{"a"}, 1 \rangle, \langle \text{"b"}, 2 \rangle \}$

$\&A, \&B$ indicate *any* type and are used to define polymorphic functions

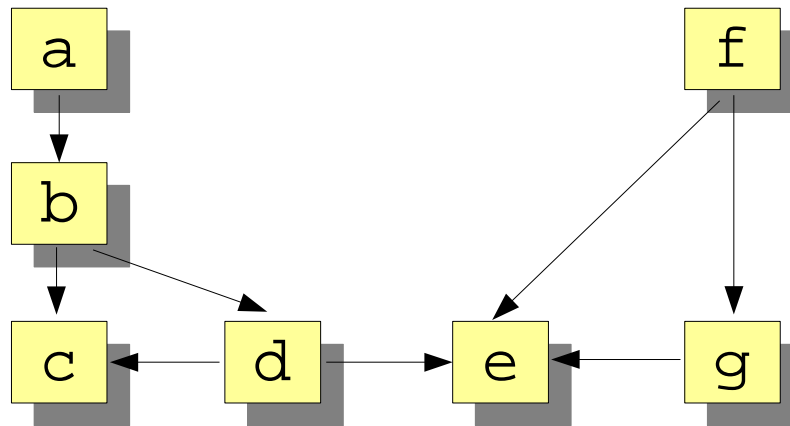
Roadmap

- Rscript in a nutshell
- Example 1: call graph analysis
- Example 2: component structure
- Example 3: Java analysis
- Example 4: a toy language
- A vizualization experiment

Roadmap

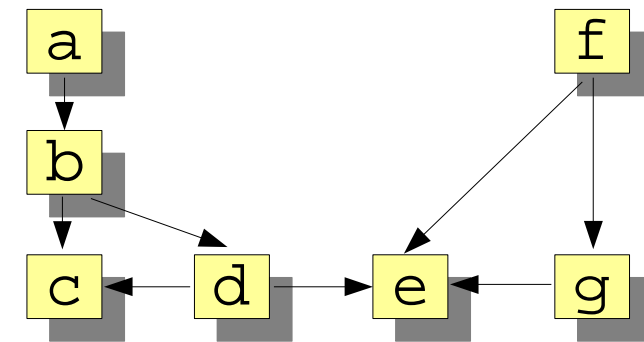
- Rscript in a nutshell
- Example 1: call graph analysis
- Example 2: component structure
- Example 3: Java analysis
- Example 4: a toy language
- A vizualization experiment

Analyzing the call structure of an application



`rel[str, str] calls = {<"a", "b">, <"b", "c">, <"b", "d">, <"d", "c">, <"d", "e">, <"f", "e">, <"f", "g">, <"g", "e">}`

Some questions



- How many calls are there?

- `int ncalls = # calls`

- 8

Number of elements

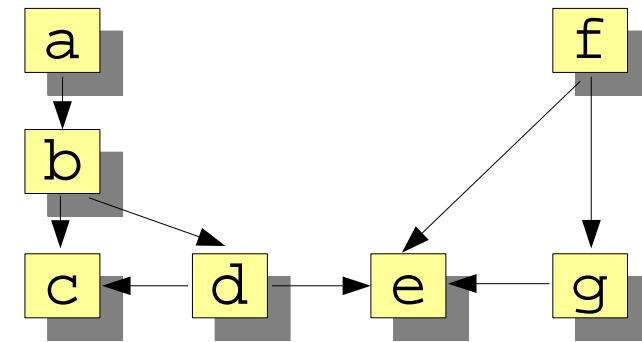
- How many procedures are there?

- `int nprocs = # carrier(calls)`

- 7

All elements in domain or range of a relations

Some questions



- What are the entry points?

- `set[str] entryPoints = top(calls)`
- `{"a", "f"}`

The *roots* of a relation
(viewed as a graph)

- What are the leaves?

- `set[str] bottomCalls = bottom(calls)`
- `{"c", "e"}`

The *leaves* of a relation
(viewed as a graph)

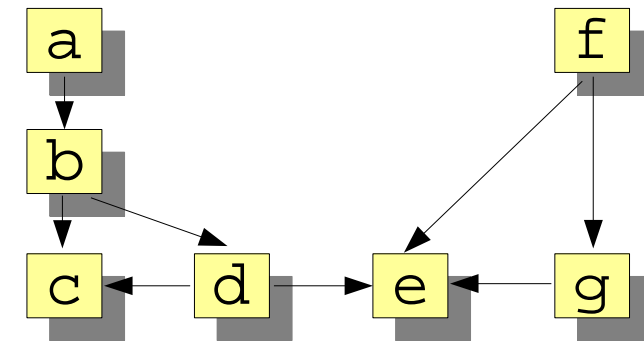
Intermezzo: Top

- The **roots** of a relation viewed as a graph
- $\text{top}(\{\langle 1,2 \rangle, \langle 1,3 \rangle, \langle 2,4 \rangle, \langle 3,4 \rangle\})$ yields $\{1\}$
- Consists of all elements that occur on the **lhs** but **not on the rhs** of a tuple
- $\text{set}[\&T] \text{ top}(\text{rel}[\&T, \&T] R) = \text{domain}(R) \setminus \text{range}(R)$

Intermezzo: Bottom

- The **leaves** of a relation viewed as a graph
- $\text{bottom}(\{ \langle 1, 2 \rangle, \langle 1, 3 \rangle, \langle 2, 4 \rangle, \langle 3, 4 \rangle \})$ yields $\{4\}$
- Consists of all elements that occur on the **rhs** but **not on the lhs** of a tuple
- $\text{set}[\&T] \text{ bottom}(\text{rel}[\&T, \&T] R) = \text{range}(R) \setminus \text{domain}(R)$

Some questions



- What are the indirect calls between procedures?

- `rel[str,str] closureCalls = calls+`

- {<"a", "b">, <"b", "c">, <"b", "d">, <"d", "c">, <"d", "e">, <"f", "e">, <"f", "g">, <"g", "e">, <"a", "c">, <"a", "d">, <"b", "e">, <"a", "e">}

The image of
domain value "a"

- What are the calls from entry point **a**?

- `set[str] calledFromA = closureCalls["a"]`

- {"b", "c", "d", "e"}

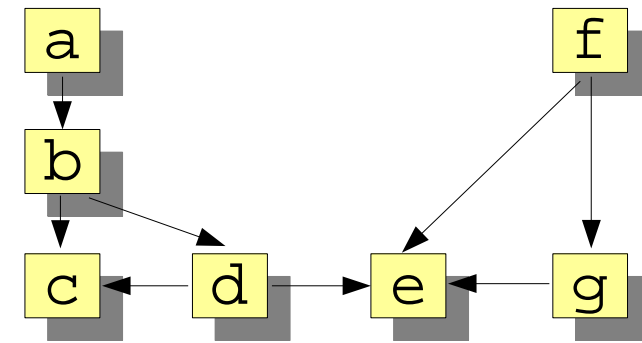
Intermezzo: right image

- Right-image of a relation: all elements that have a given value as left element (**resembles array access**)
- Notation: relation followed by [Value]
- Ex. $Rel = \{ \langle 1, 10 \rangle, \langle 2, 20 \rangle, \langle 1, 11 \rangle, \langle 3, 30 \rangle, \langle 2, 21 \rangle \}$
- $Rel[1]$ yields $\{10, 11\}$
- $Rel[\{1, 2\}]$ yields $\{10, 11, 20, 21\}$

Intermezzo: left image

- Left-image of a relation: all elements that have a given value as right element
- Notation: relation followed by $[-, \text{Value}]$
- Ex. $\text{Rel} = \{ \langle 1, 10 \rangle, \langle 2, 20 \rangle, \langle 1, 11 \rangle, \langle 3, 30 \rangle, \langle 2, 21 \rangle \}$
- $\text{Rel}[-, 10]$ yields $\{1\}$
- $\text{Rel}[-, \{10, 20\}]$ yields $\{1, 2\}$

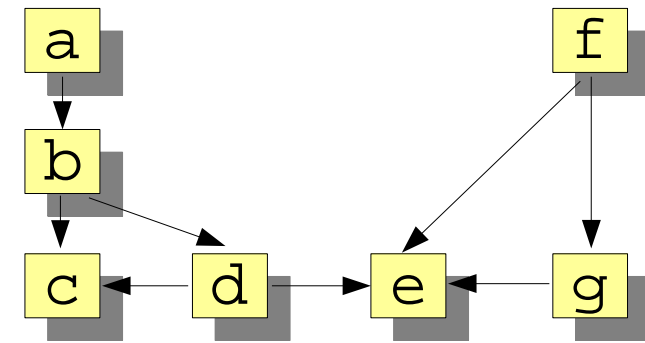
Some questions



- What are the calls to procedure **e**?
 - `set[str] callsToE = closureCalls[-,"e"]`
 - `{"a", "b", "d", "f", "g"}`

The domain of
image value "e"

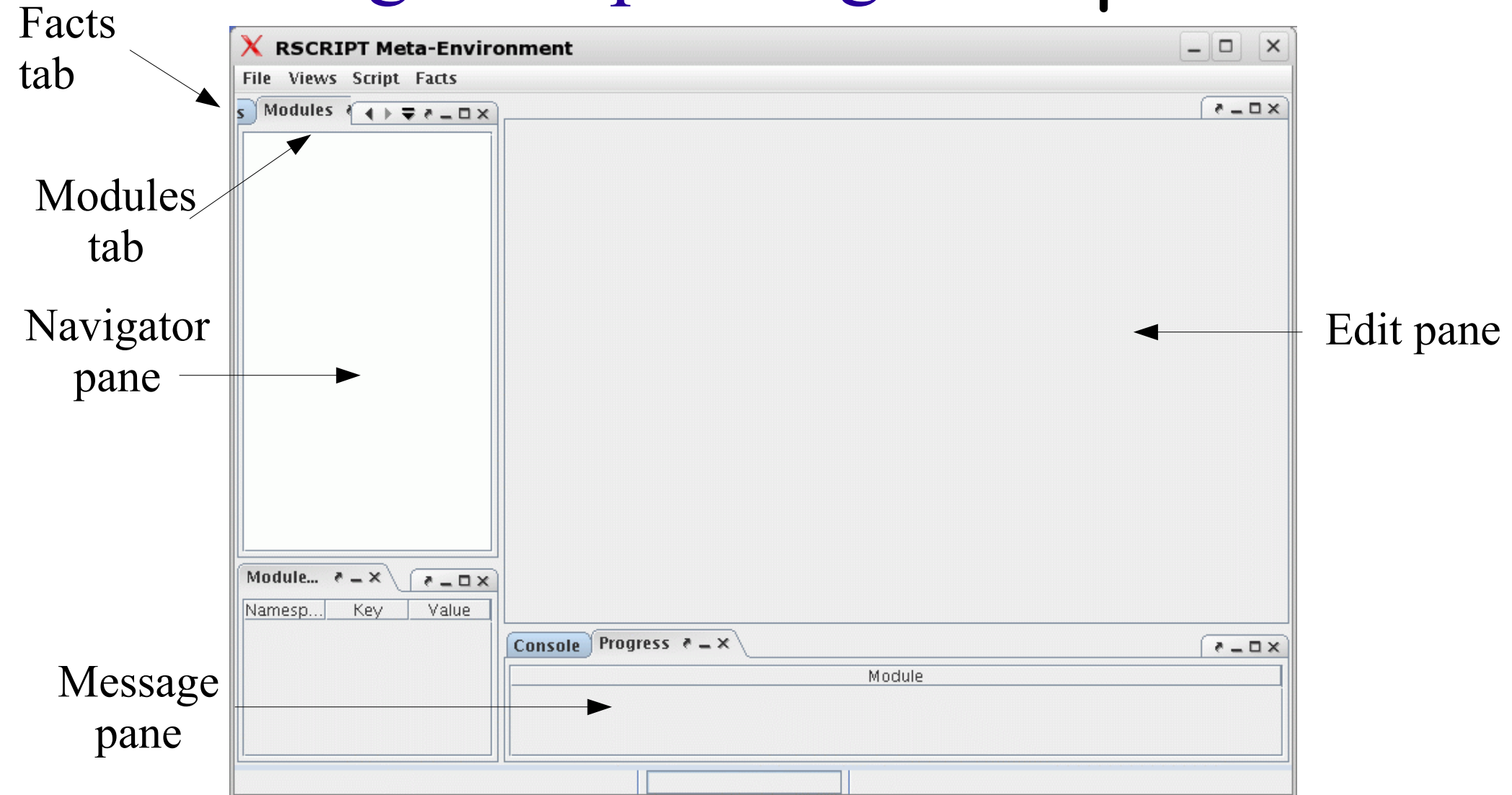
Some questions



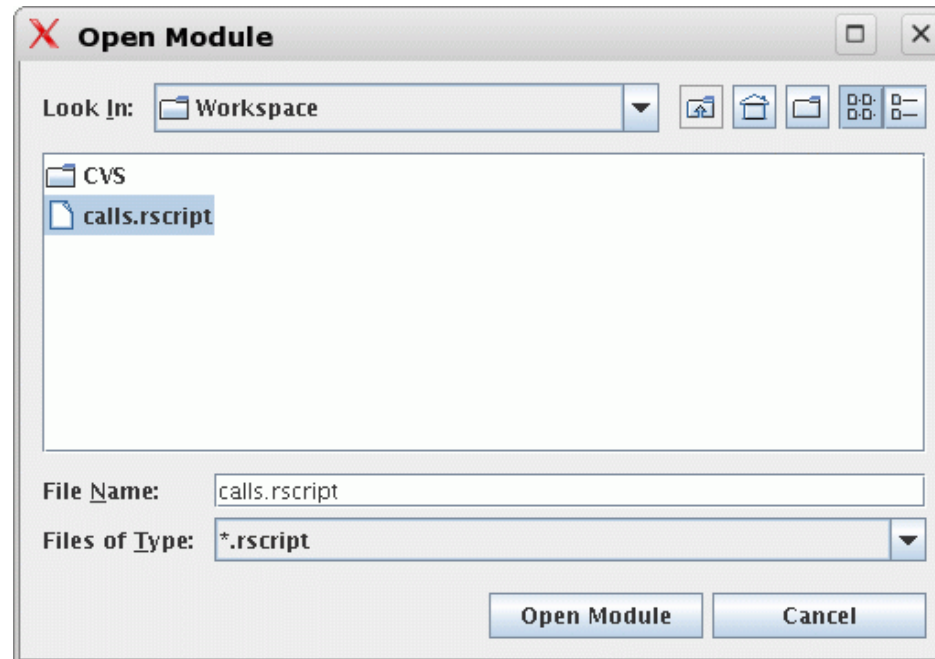
- What are the calls from entry point **f**?
 - `set[str] calledFromF = closureCalls["f"]`
 - `{"e", "g"}`
- What are the common procedures?
 - `set[str] commonProcs =`
`calledFromA inter calledFromF`
 - `{"e"}`

Intersection

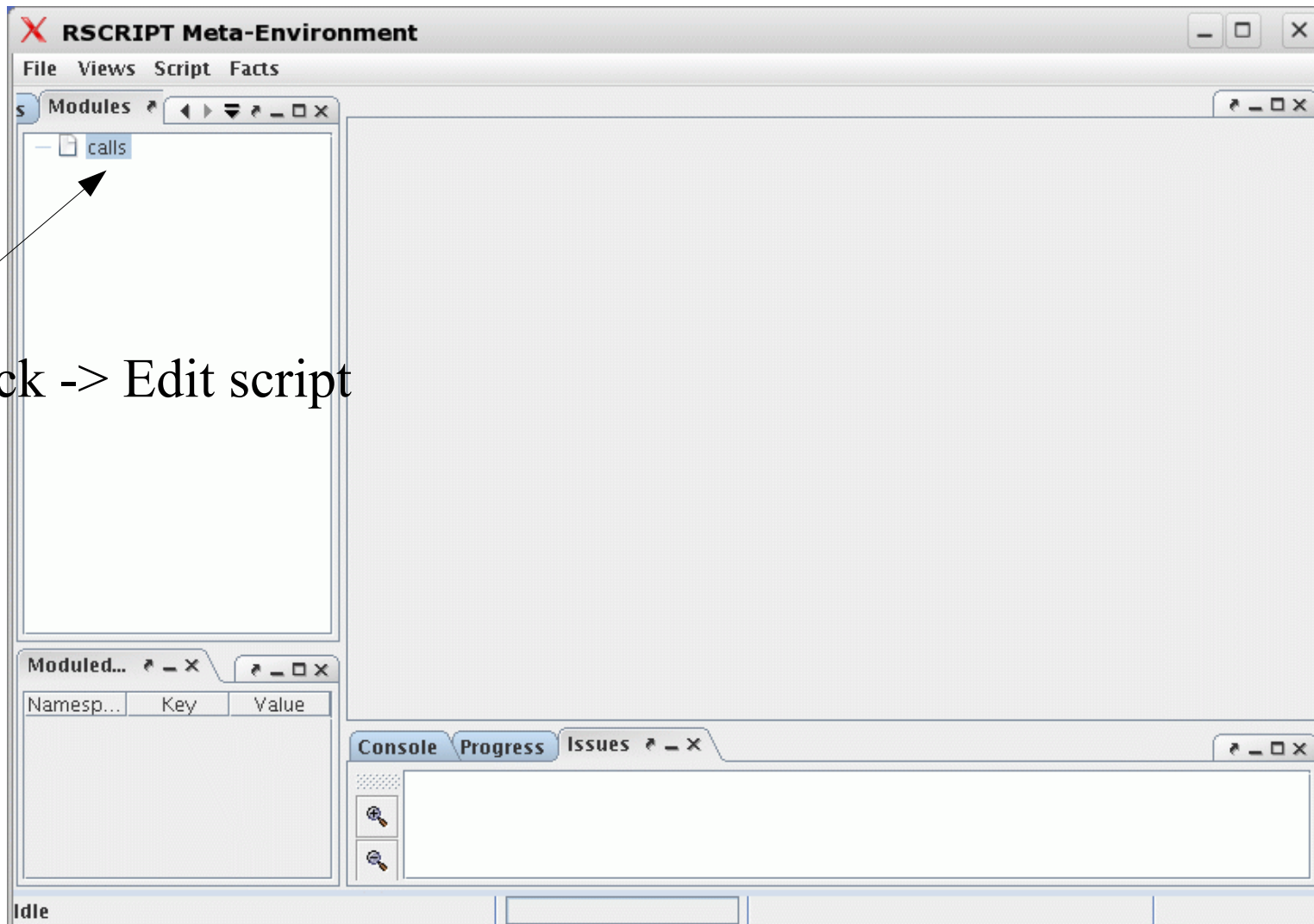
Running Rscript using rscript-meta



Script -> Open...

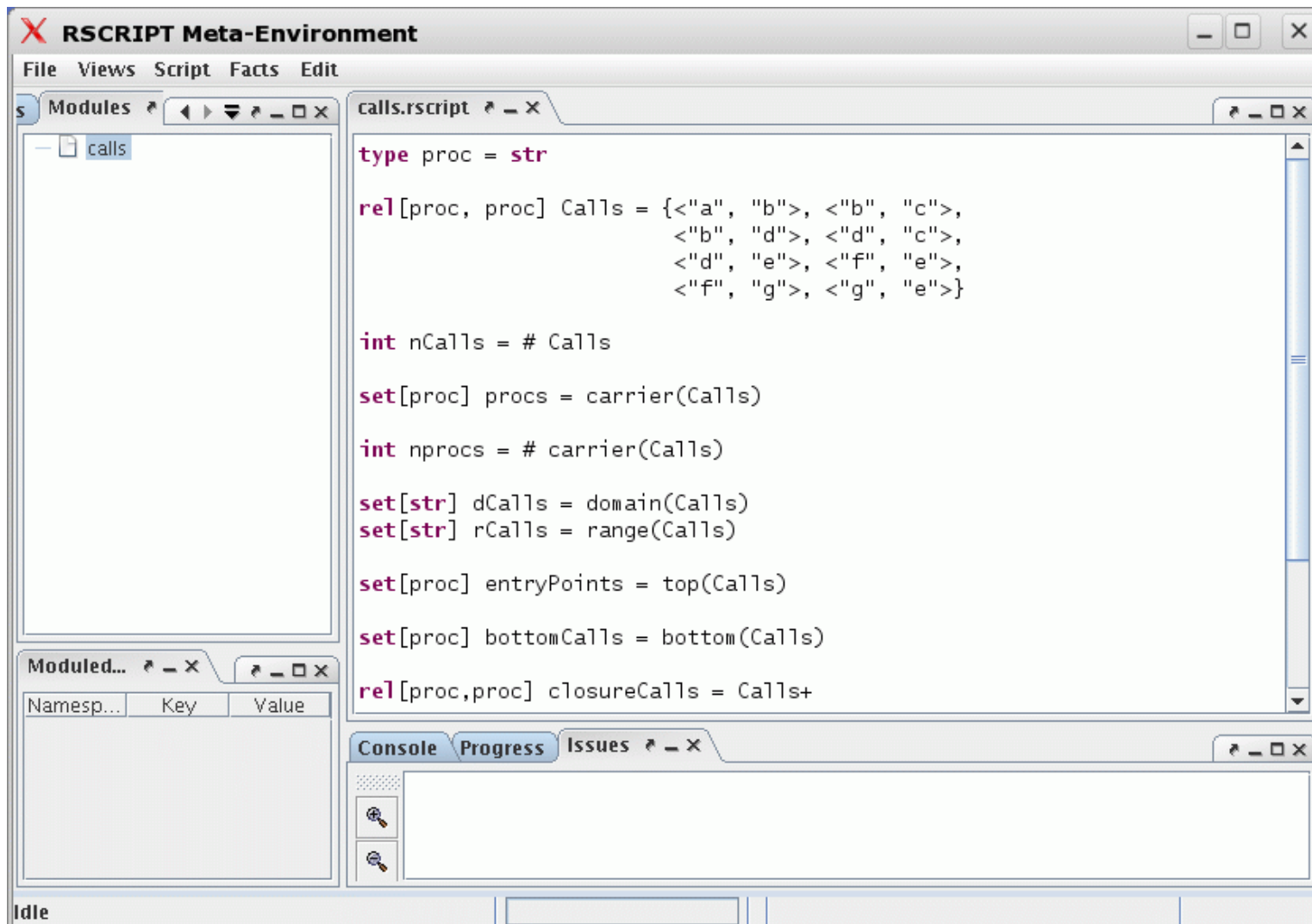


File calls has been opened

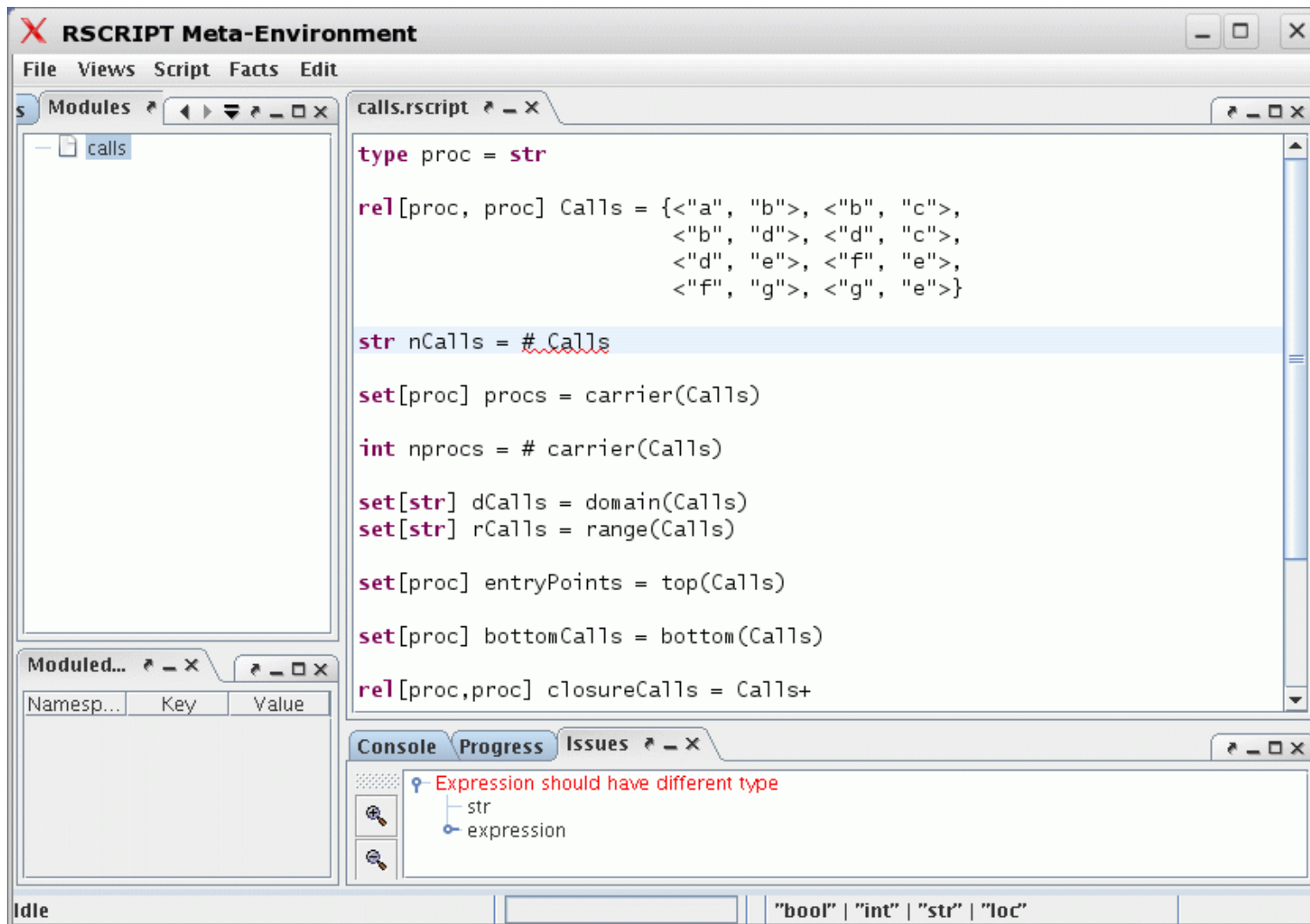


Right click -> Edit script

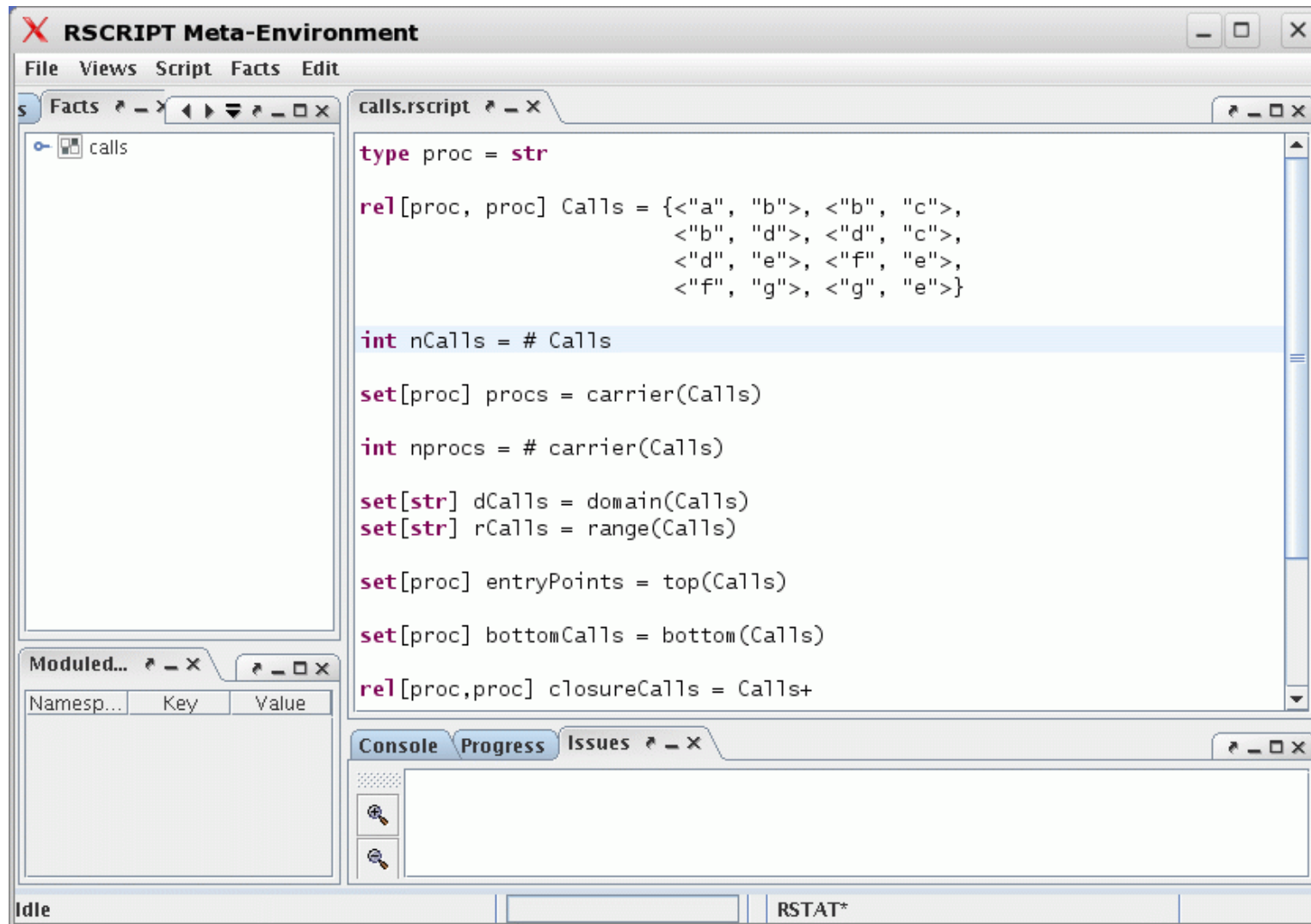
Editing calls.rscript



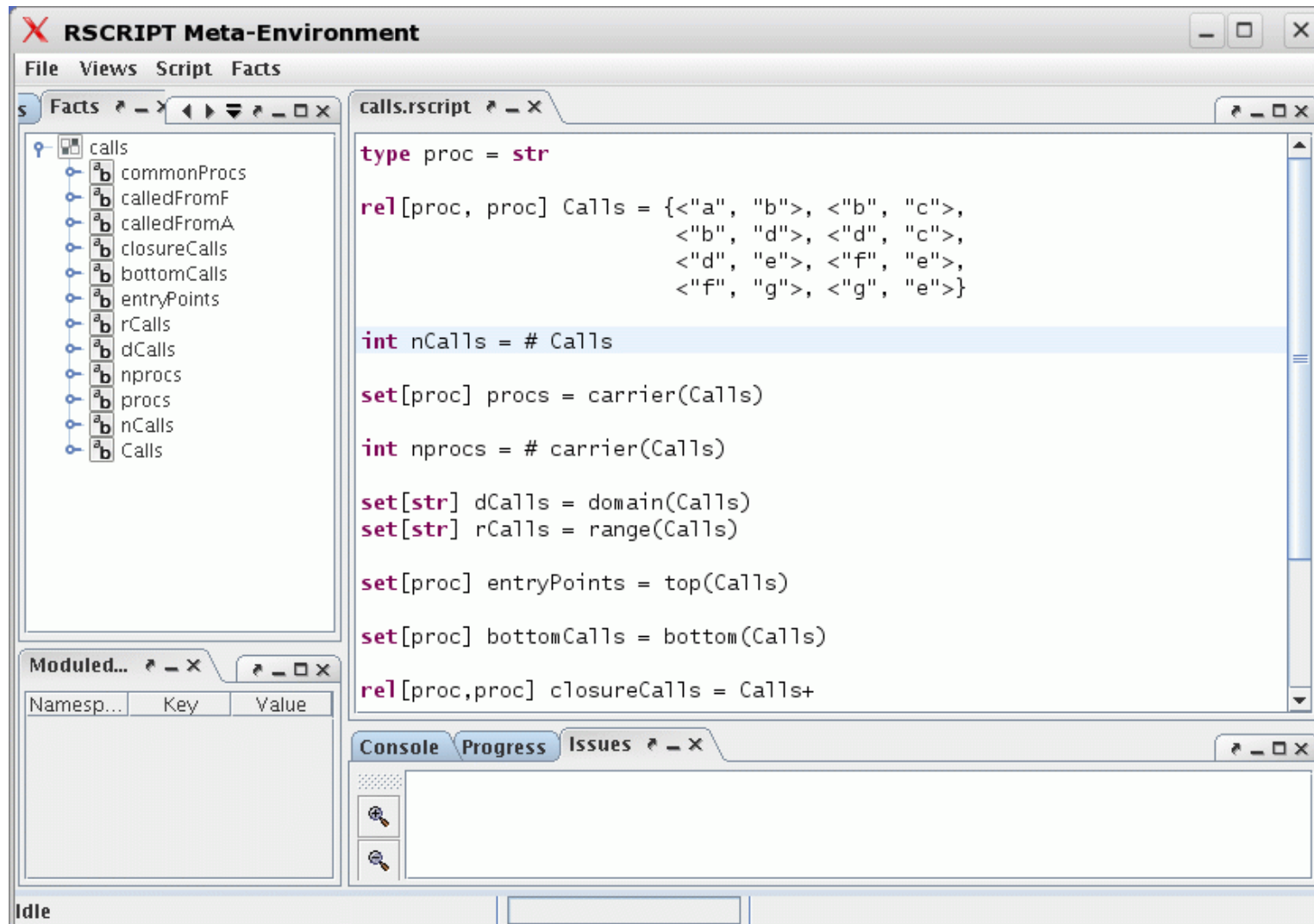
Making errors ...



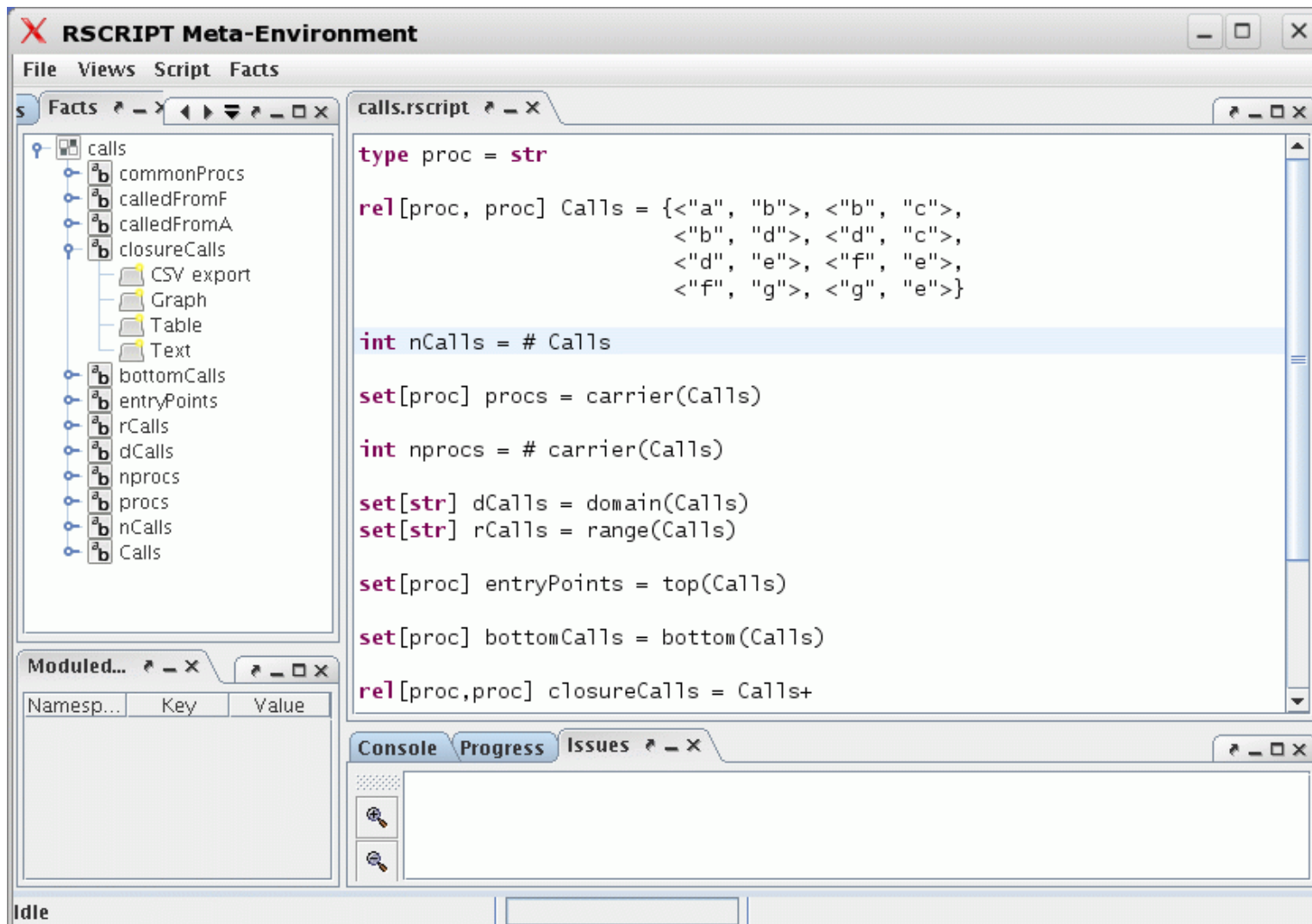
Script -> Run



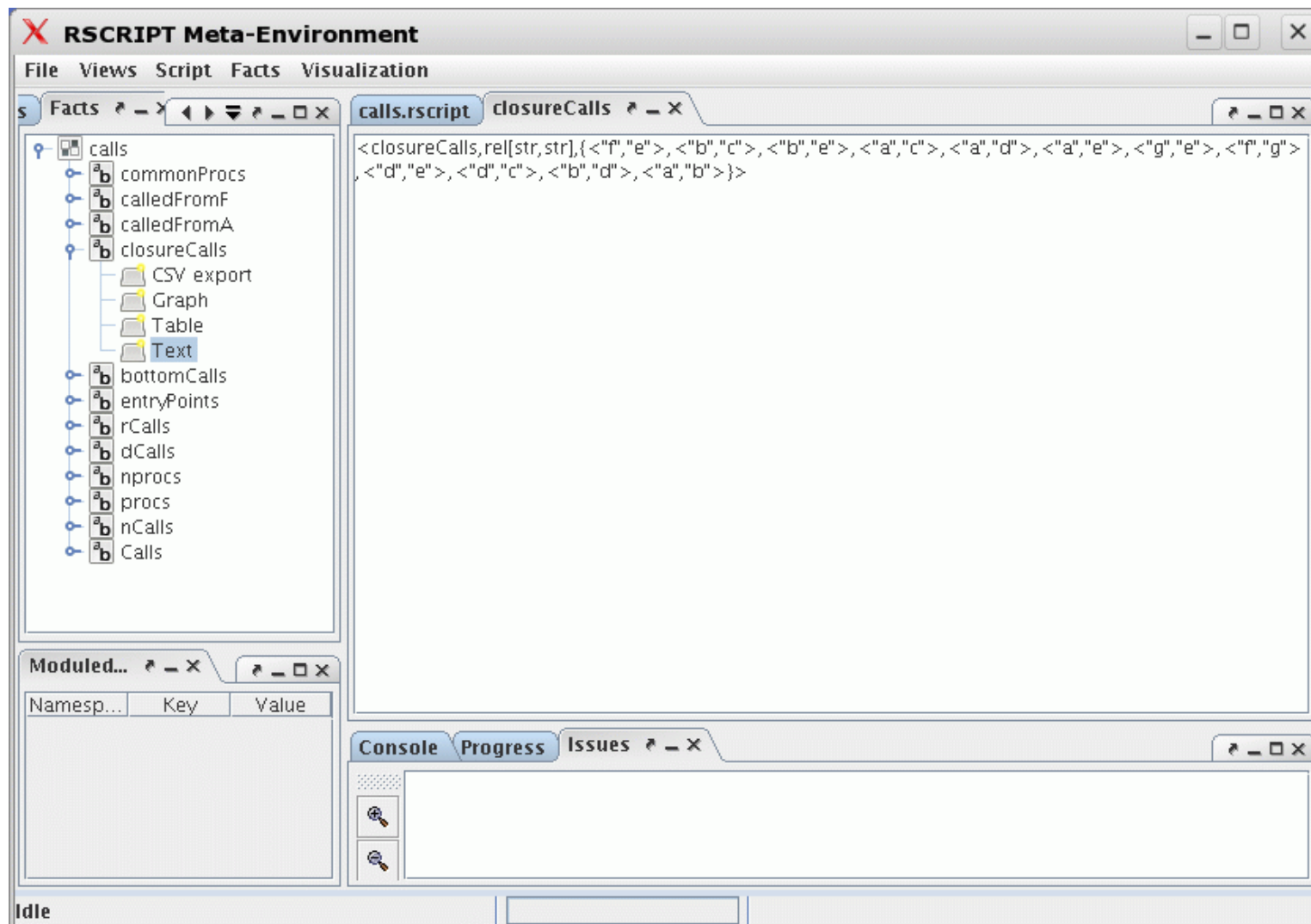
Unfolding the rstore ...



Unfolding closureCalls



closureCalls as Text

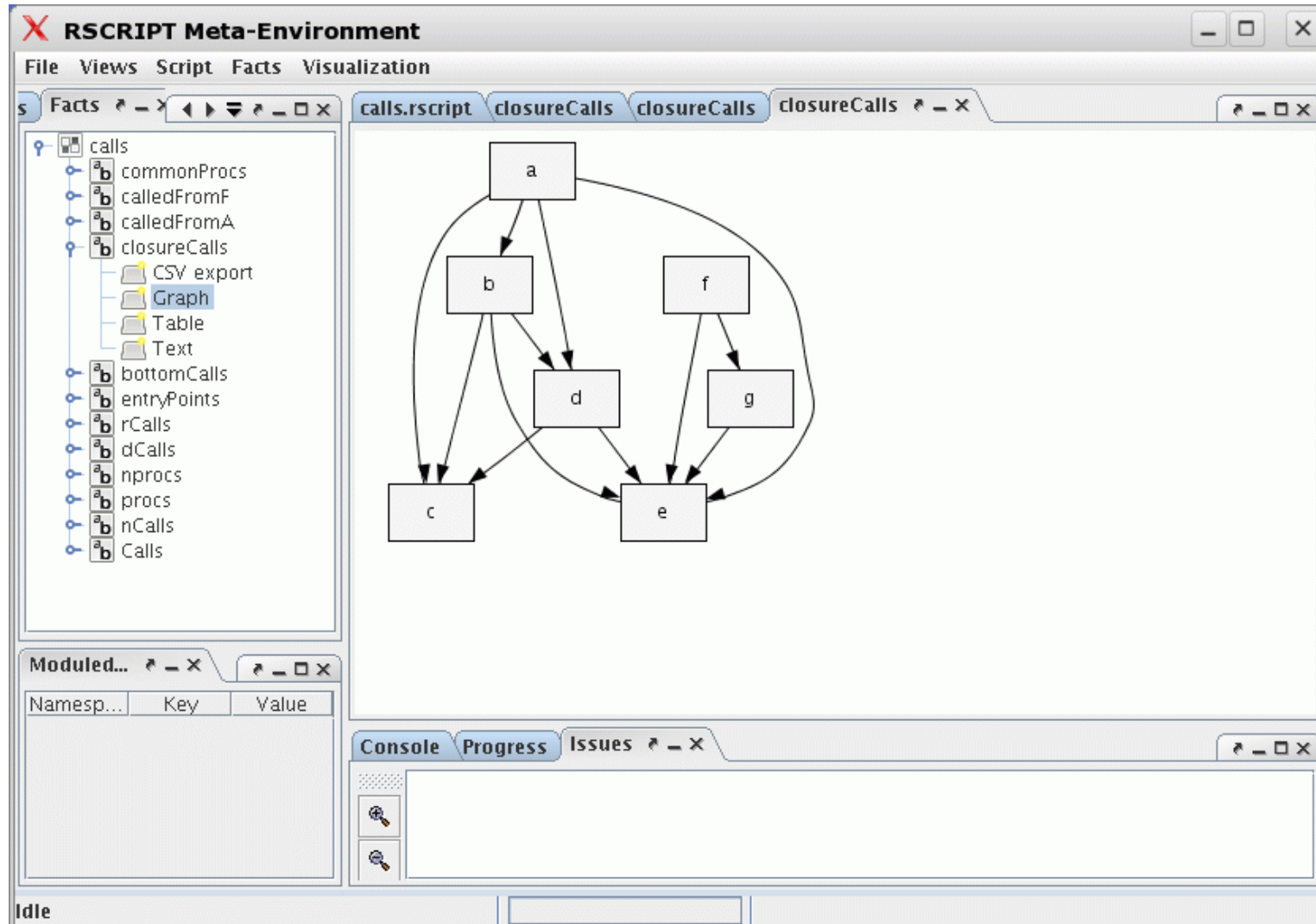


closureCalls as Table

The screenshot shows the RSCRIPT Meta-Environment interface. The main window displays a table titled 'closureCalls' with two columns: 'str [0]' and 'str [1]'. The table contains 16 rows of data. The left sidebar shows a tree view of the environment, with 'closureCalls' selected under the 'calls' folder. The bottom status bar indicates 'Idle'.

str [0]	str [1]
a	b
b	d
d	c
d	e
f	g
g	e
a	e
a	d
a	c
b	e
b	c
f	e

closureCalls as Graph



Roadmap

- Rscript in a nutshell
- Example 1: call graph analysis
- Example 2: component structure
- Example 3: Java analysis
- Example 4: a toy language
- A vizualization experiment

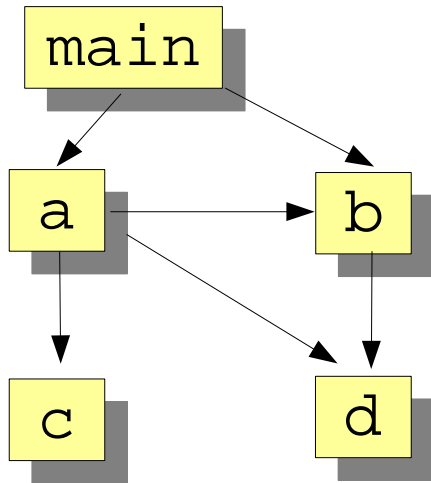
Roadmap

- Rscript in a nutshell
- Example 1: call graph analysis
- **Example 2: component structure**
- Example 3: Java analysis
- Example 4: a toy language
- A vizualization experiment

Component Structure of Application

- Suppose, we know:
 - the call relation between procedures (**Calls**)
 - the component of each procedure (**PartOf**)
- Question:
 - Can we lift the relation between procedures to a relation between components (**ComponentCalls**)?
- This is usefull for checking that real code conforms to architectural constraints

Calls

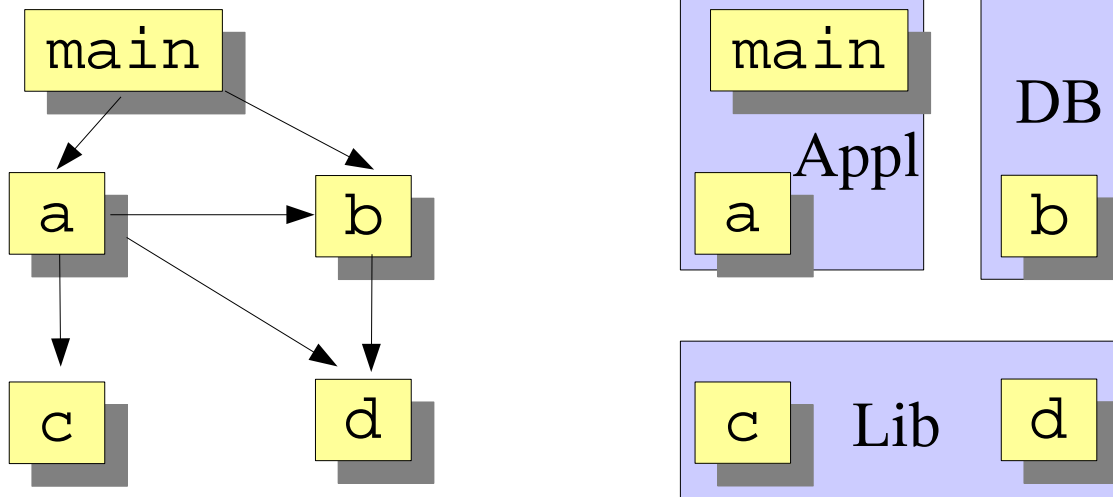


type proc = str

type comp = str

rel[proc,proc] Calls = {<"main", "a">, <"main", "b">, <"a", "b">,
<"a", "c">, <"a", "d">, <"b", "d">}

PartOf

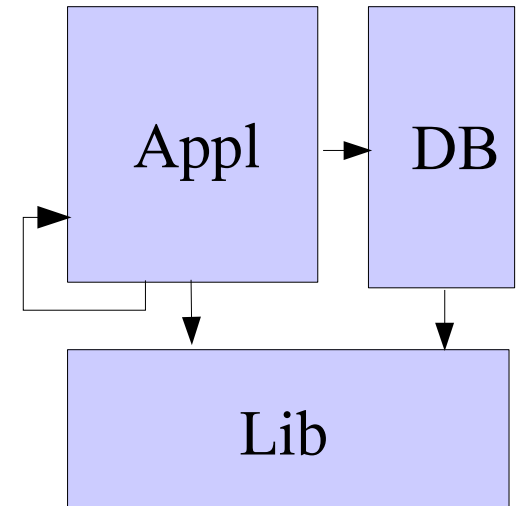
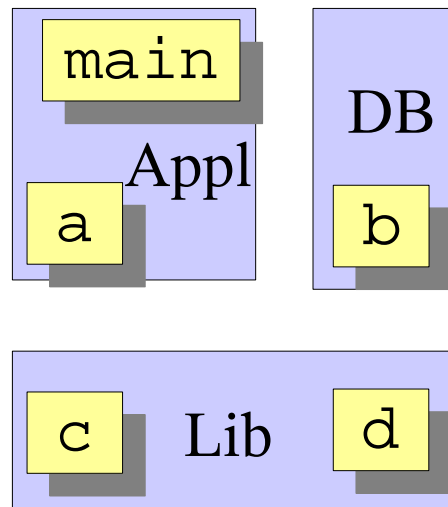
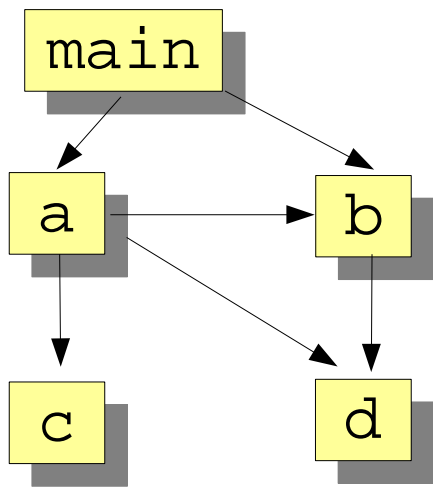


```
set[comp] Components = {"Appl", "DB", "Lib"}
```

```
rel[proc, comp] PartOf =
```

```
  {<"main", "Appl">, <"a", "Appl">, <"b", "DB">,
    <"c", "Lib">, <"d", "Lib">}
```

lift



$\text{rel}[\text{comp}, \text{comp}] \text{ lift}(\text{rel}[\text{proc}, \text{proc}] \text{ aCalls}, \text{rel}[\text{proc}, \text{comp}] \text{ aPartOf}) =$
 $\{ \langle C1, C2 \rangle \mid \langle \text{proc } P1, \text{proc } P2 \rangle : \text{aCalls},$
 $\quad \langle \text{comp } C1, \text{comp } C2 \rangle : \text{aPartOf}[P1] \times \text{aPartOf}[P2] \}$

$\text{rel}[\text{comp}, \text{comp}] \text{ ComponentCalls} = \text{lift}(\text{Calls2}, \text{PartOf})$

Resultaat: $\{ \langle \text{"DB"}, \text{"Lib"} \rangle, \langle \text{"Appl"}, \text{"Lib"} \rangle, \langle \text{"Appl"}, \text{"DB"} \rangle, \langle \text{"Appl"}, \text{"Appl"} \rangle \}$

Roadmap

- Rscript in a nutshell
- Example 1: call graph analysis
- **Example 2: component structure**
- Example 3: Java analysis
- Example 4: a toy language
- A vizualization experiment

Roadmap

- Rscript in a nutshell
- Example 1: call graph analysis
- Example 2: component structure
- **Example 3: Java analysis**
- Example 4: a toy language
- A vizualization experiment

Cyclic Dependencies

A class uses (directly or indirectly) itself

Use = methods calls, inheritance, containment

```
class ContainedClass { }  
class SuperClass {  
  class SubClass extends SuperClass {  
    ContainedClass C;  
  }  
}
```

Example of
a contained class

Motivation: cyclic class dependencies are difficult to understand/maintain

Cyclic Dependencies: Examples

```
class A { B B1; ... }  
class B extends A { ... }
```

```
class A { C C1; ... }  
class B extends A { ... }  
class C { B B1; ... }
```

Java analysis: classes in cycles

- Assume the following extracted information:
 - $\text{rel}[\text{str}, \text{str}]$ CALL
 - method call from first class to the second
 - $\text{rel}[\text{str}, \text{str}]$ INHERITANCE
 - extends and implements
 - $\text{rel}[\text{str}, \text{str}]$ CONTAINMENT
 - attribute of first class is of the type of the second class
- Question: which classes occur in a cyclic dependency?

Java analysis: cycles in classes

- Define the USE relation between two classes:
 - $\text{rel}[\text{str}, \text{str}] \text{ USE} = \text{CALL} \text{ union } \text{CONTAINMENT} \text{ union } \text{INHERITANCE}$
 - $\text{set}[\text{str}] \text{ ClassesInCycle} = \{C1 \mid \langle \text{str } C1, \text{str } C2 \rangle : \text{USE}^+, C1 \neq C2\}$
- In this way we get a set of classes that occur in a cyclic dependency, but ...
- ... which classes are in the cycle?

Java analysis: cyclic classes

- $\text{rel}[\text{str}, \text{str}] \text{ USE} = \text{CALL} \text{ union } \text{CONTAINMENT} \text{ union } \text{INHERITANCE}$
- $\text{set}[\text{str}] \text{ CLASSES} = \text{carrier}(\text{USE})$
- $\text{rel}[\text{str}, \text{str}] \text{ USETRANS} = \text{USE}^+$
- $\text{rel}[\text{str}, \text{set}[\text{str}]] = \{ \langle C, \text{USETRANS}[C] \rangle \mid \text{str } C : \text{CLASSES}, \langle C, C \rangle \text{ in USETRANS} \}$
- Each cyclic class is associated with a set of classes that form a cycle

Applications of this approach

- Search for “similar” classes
- Search for design patterns (as characterized by specific relations between the classes in the pattern)
- ...

Roadmap

- Rscript in a nutshell
- Example 1: call graph analysis
- Example 2: component structure
- **Example 3: Java analysis**
- Example 4: a toy language
- A vizualization experiment

Roadmap

- Rscript in a nutshell
- Example 1: call graph analysis
- Example 2: component structure
- Example 3: Java analysis
- Example 4: a toy language
- A vizualization experiment

Toy program

```
begin declare x : natural, y : natural,  
            z : natural;
```

```
  x := 3;
```

```
  if    3 then
```

```
    z := y + x
```

```
  else
```

```
    x := 4
```

```
  fi
```

```
  y := z
```

```
end
```

y is undefined

z may be undefined

Toy program

$\text{rel}[\text{int}, \text{str}] \text{ DEFS} = \{ \langle 1, "x" \rangle, \langle 3, "z" \rangle, \langle 4, "x" \rangle, \langle 5, "y" \rangle \}$

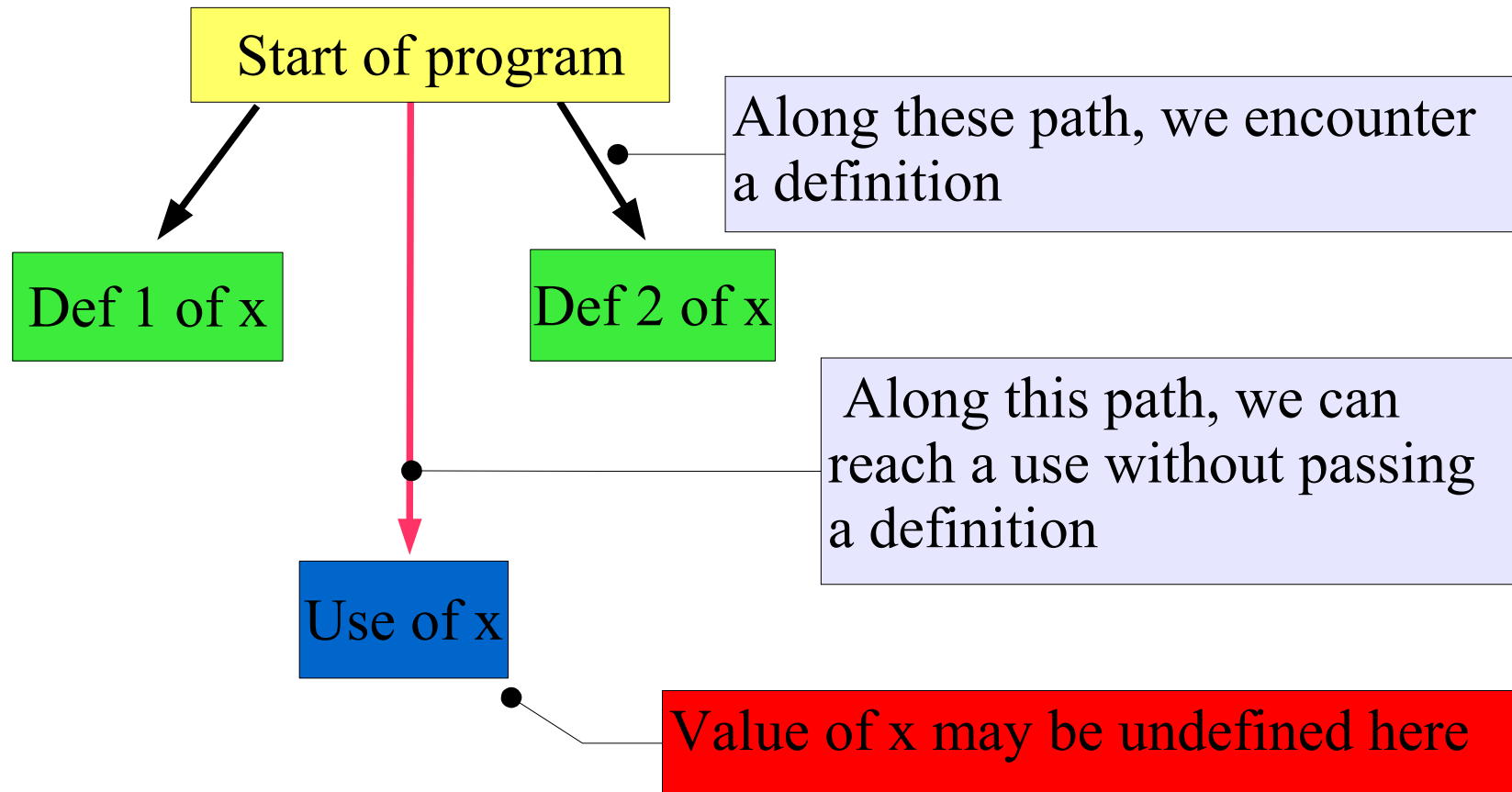
```
begin declare x : natural, y : natural,  
            z : natural;
```

```
[1] x := 3;  
    if [2] 3 then  
        [3] z := y + x  
    else  
        [4] x := 4  
    fi  
[5] y := z  
end
```

$\text{rel}[\text{int}, \text{str}] \text{ USES} = \{ \langle 3, "y" \rangle, \langle 3, "x" \rangle, \langle 5, "z" \rangle \}$

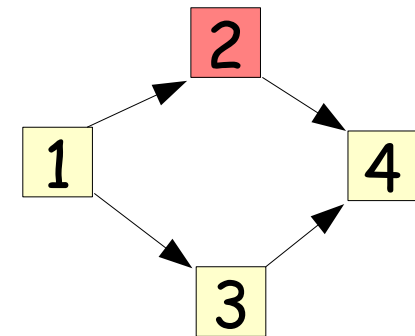
$\text{rel}[\text{int}, \text{int}] \text{ PRED} = \{ \langle 0, 1 \rangle, \langle 1, 2 \rangle, \langle 2, 3 \rangle, \langle 2, 4 \rangle, \langle 3, 5 \rangle, \langle 4, 5 \rangle \}$

Finding uninitialized variables



Intermezzo: reachX

- Reachability with exclusion of certain elements
- `set[&T] reachX(
 – set[&T] Start,
 – set[&T] Excl,
 – rel[&T,&T] Rel)`
- `reachX({1}, {2}, {<1,2>, <1,3>, <2,4>, <3,4>})`
 yields {<3,4>}



The undefined query

rel[int,str] DEFS = ...

rel[int,str] USES = ...

rel[int,int] PRED = ...

rel[int,str] UNINIT =

{ <N,V> | <int N, str V>:USES, N in reachX({0}, DEFS[-,V],PRED)}

Start from the root

Exclude all
definitions of V

There is a path from the root
to N: V is not initialized

Use the PRED relation

Reach exclude

Applying the undefined query

```
begin declare x : natural, y : natural,  
           z : natural;
```

```
[1] x := 3;
```

```
  if [2] 3 then
```

```
    [3] z := y • + x
```

```
  else
```

```
    [4] x := 4
```

```
  fi
```

```
[5] y := z •
```

```
end
```

y is undefined

z may be undefined

Result:

{<5,"z">, <3,"y">}

Some Questions

- There are several additional questions:
 - In the example so far we have worked with statement numbers but how do we make a connection with the source text? (Discussed now)
 - How do we extract relations like **PRED** and **USE** from the source text? (Discussed later)

Use locations to connect with the source text

```
rel[int,str] DEFS = ...  
rel[int,str] USES = ...  
rel[int,int] PRED = ...
```

Use *location* instead of number

```
rel[loc,str] DEFS  
rel[loc,str] USES  
rel[loc,loc] PRED  
rel[str,loc] OCCURS
```

Variable occurrence
in a statement

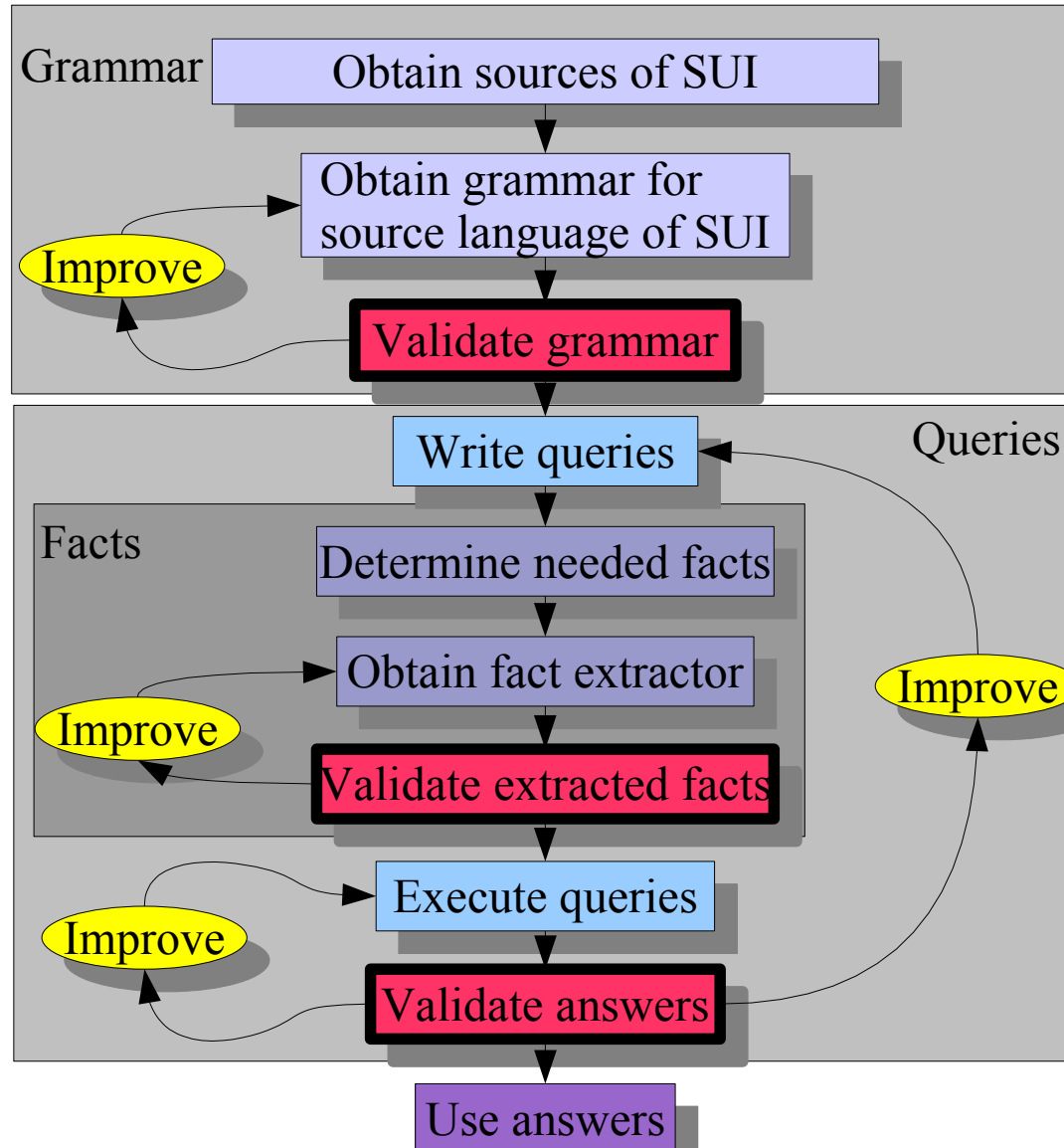
Example Rstore

```
rstore(  
  <PRED, rel[loc,loc],  
    {<area-in-file("/home/paulk/.../example.pico", area(4, 2,4, 8,84, 6)),  
      area-in-file("/home/paulk/.../example.pico", area(5, 2,5, 8,94, 6))>,  
    <area-in-file("/home/paulk/.../example.pico", area(5, 2,5, 8, 94, 6)),  
      area-in-file("/home/paulk/.../example.pico", area(6, 2,10, 4, 104, 56))>,  
    ... }>,  
  
  <DEFS, {  
    <OCCURS, rel[str,loc],  
      {<"y", area-in-file("/home/paulk/.../example.pico", area(11, 2,11, 3,164, 1))>,  
        <"z", area-in-file("/home/paulk/.../example.pico", area(11, 7,11, 8,169, 1))>,  
        ... }  
    }  
  )  
)
```

Extracting Facts

- Goal: extract facts from source code and use as input for queries
- How should fact extraction be organized?
- How to write a fact extractor?

Workflow Fact Extraction



SUI
=
System
Under
Investigation

Roadmap

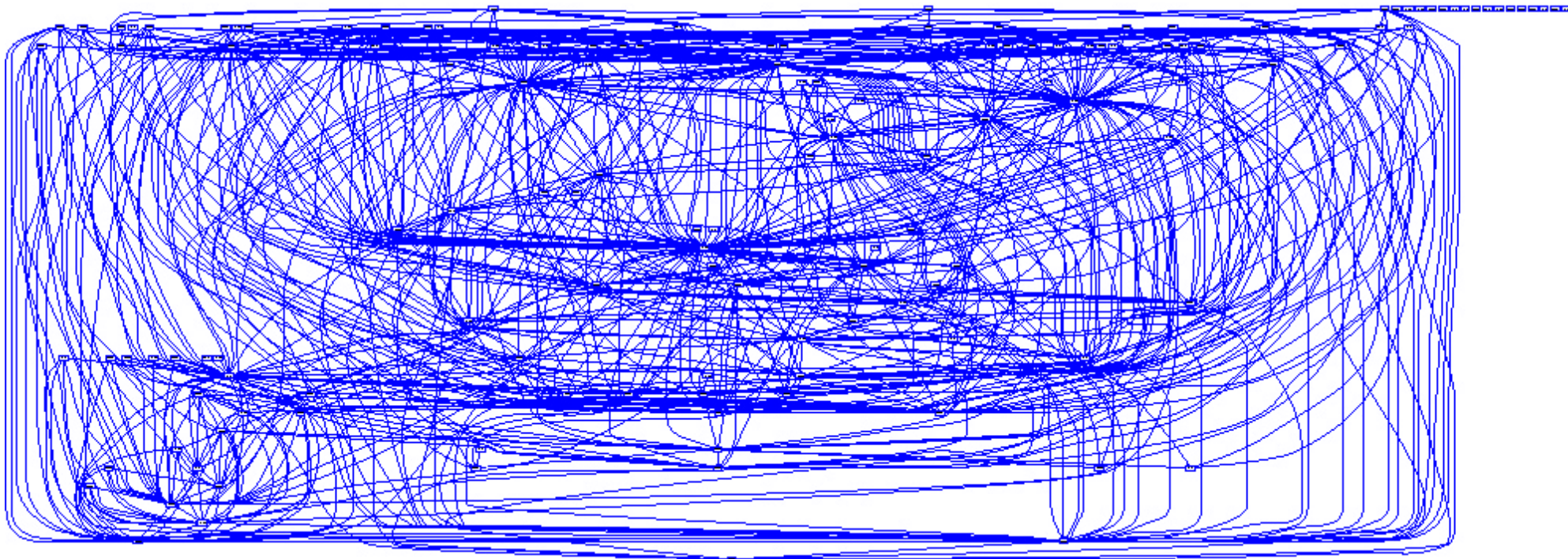
- Rscript in a nutshell
- Example 1: call graph analysis
- Example 2: component structure
- Example 3: Java analysis
- Example 4: a toy language
- A vizualization experiment

Roadmap

- Rscript in a nutshell
- Example 1: call graph analysis
- Example 2: component structure
- Example 3: Java analysis
- Example 4: a toy language
- A vizualization experiment

Issues in Program Visualization

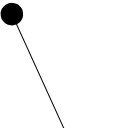
- Small graphs are nice, large graph are a disaster



(Courtesy: Arie van Deursen)

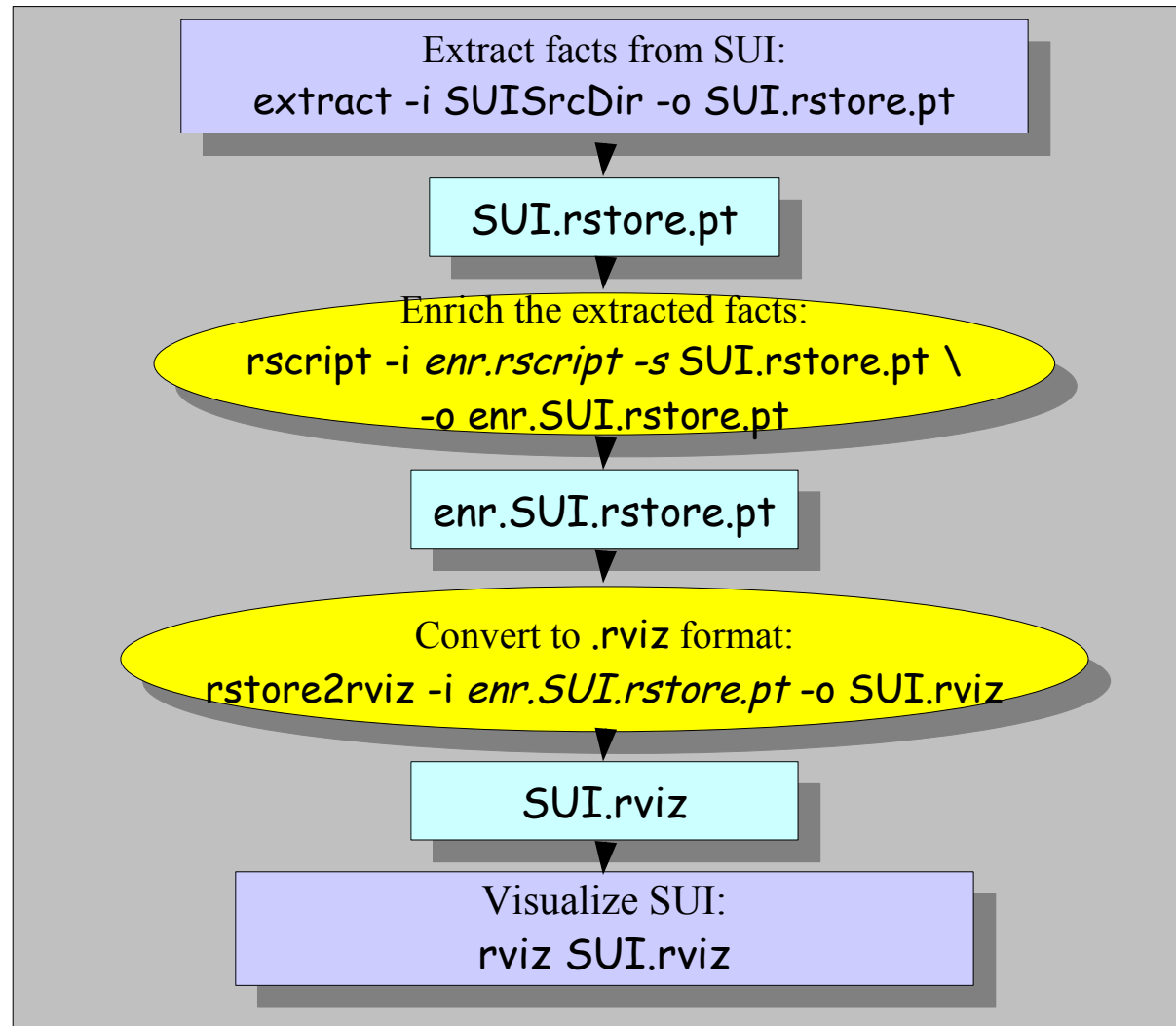
Issues in Program Visualization

- Howto display information related to source text?
- Approach (Steven Eick): use a pixel-based image of the source text
- Over 100.000 LOC on one screen!
- Experiment: visualize an Rstore for JHotDRaw (15.000 LOC)



Extraction by Hayco de Jong and
Taeke Kooiker (using ASF+SDF)

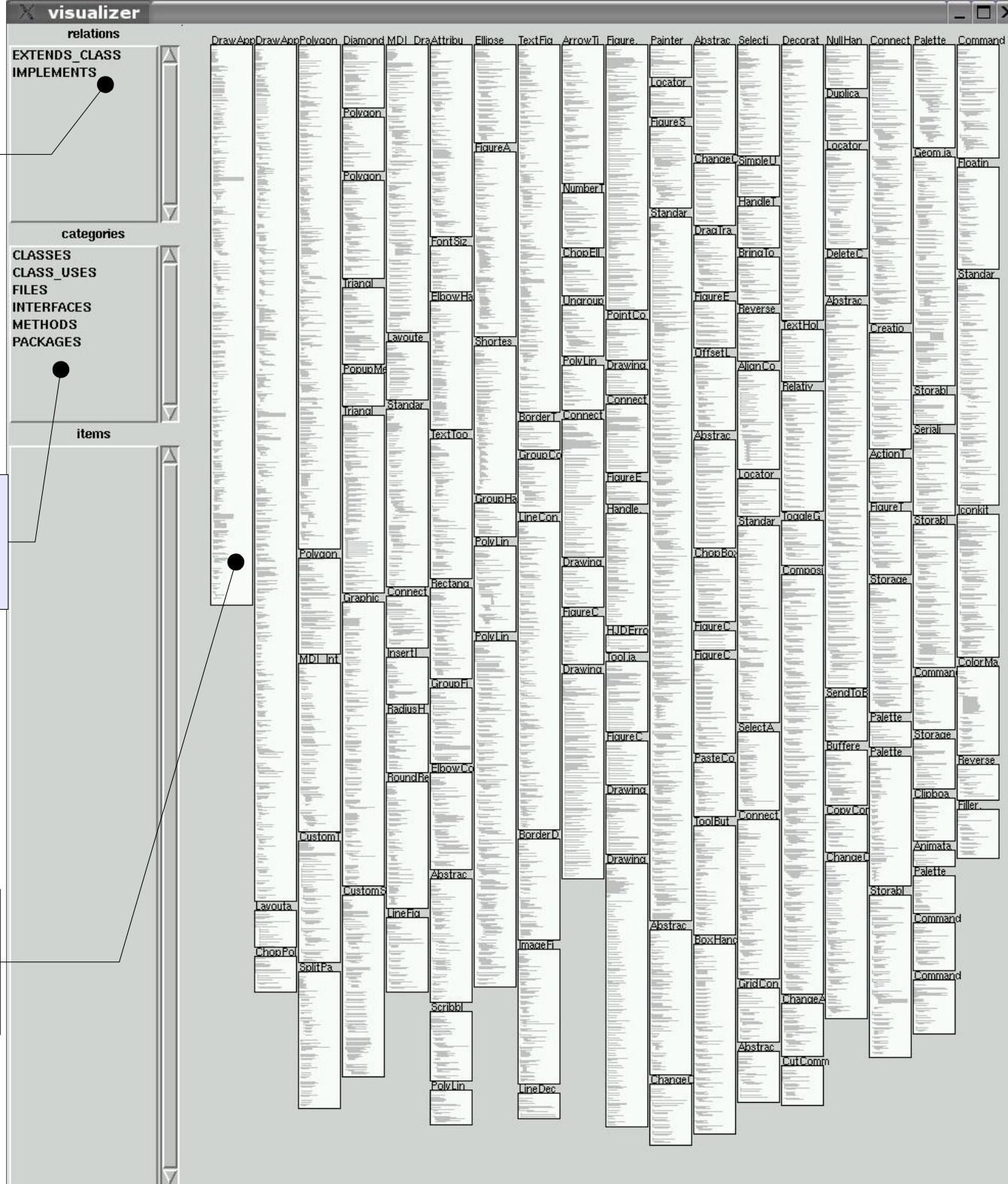
Visualization Workflow

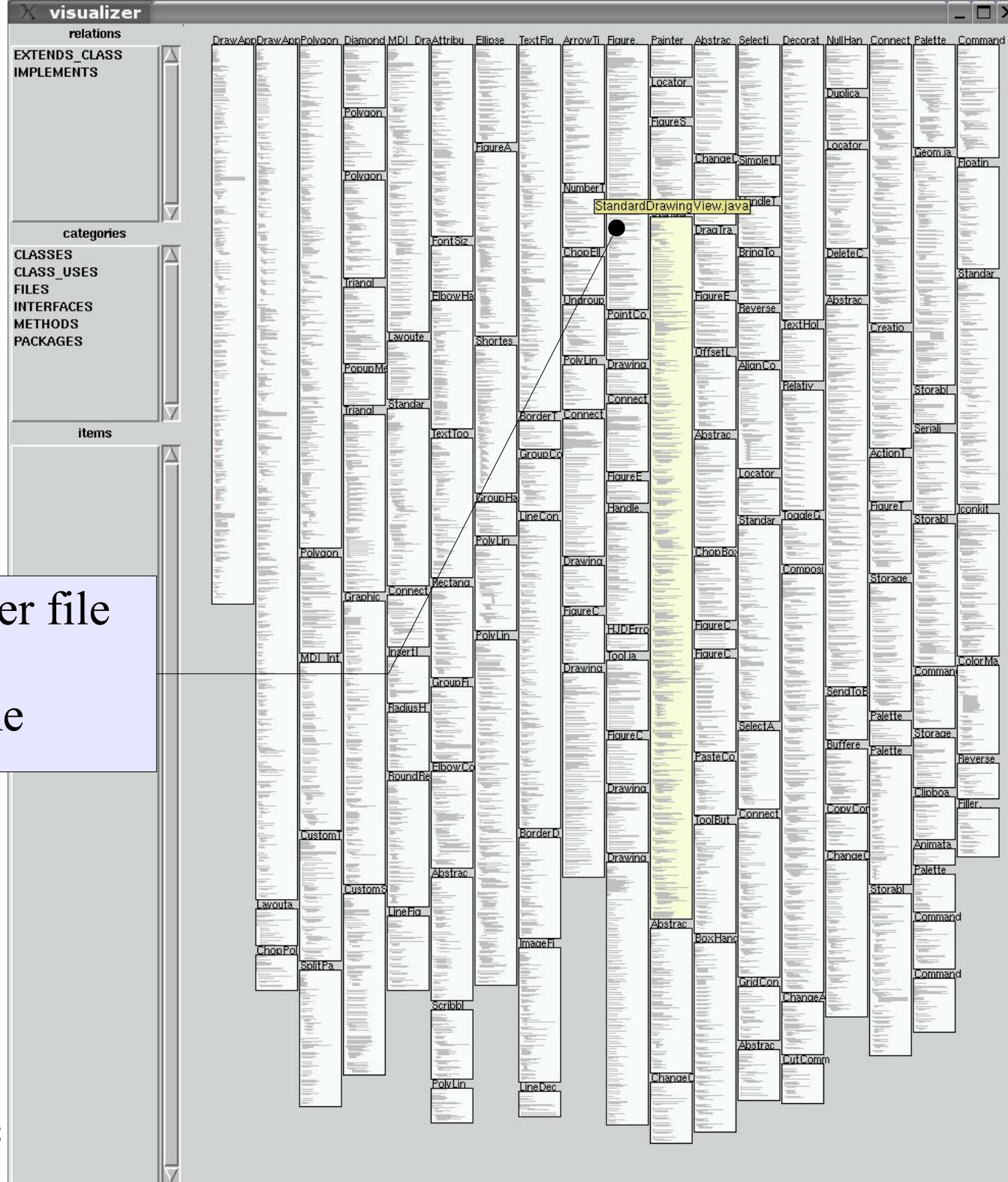


Relations

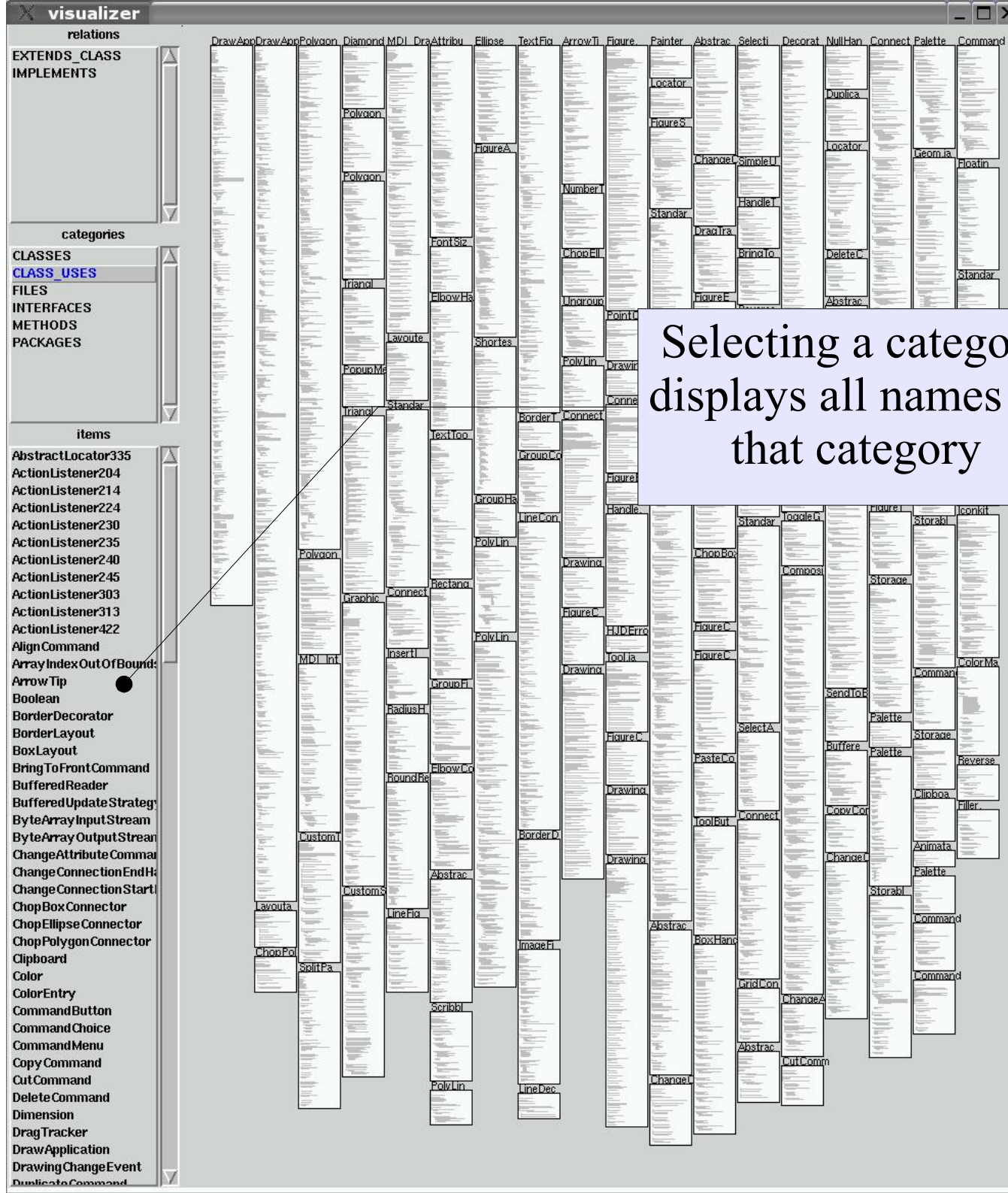
Categories
of names

Rectangle
per file

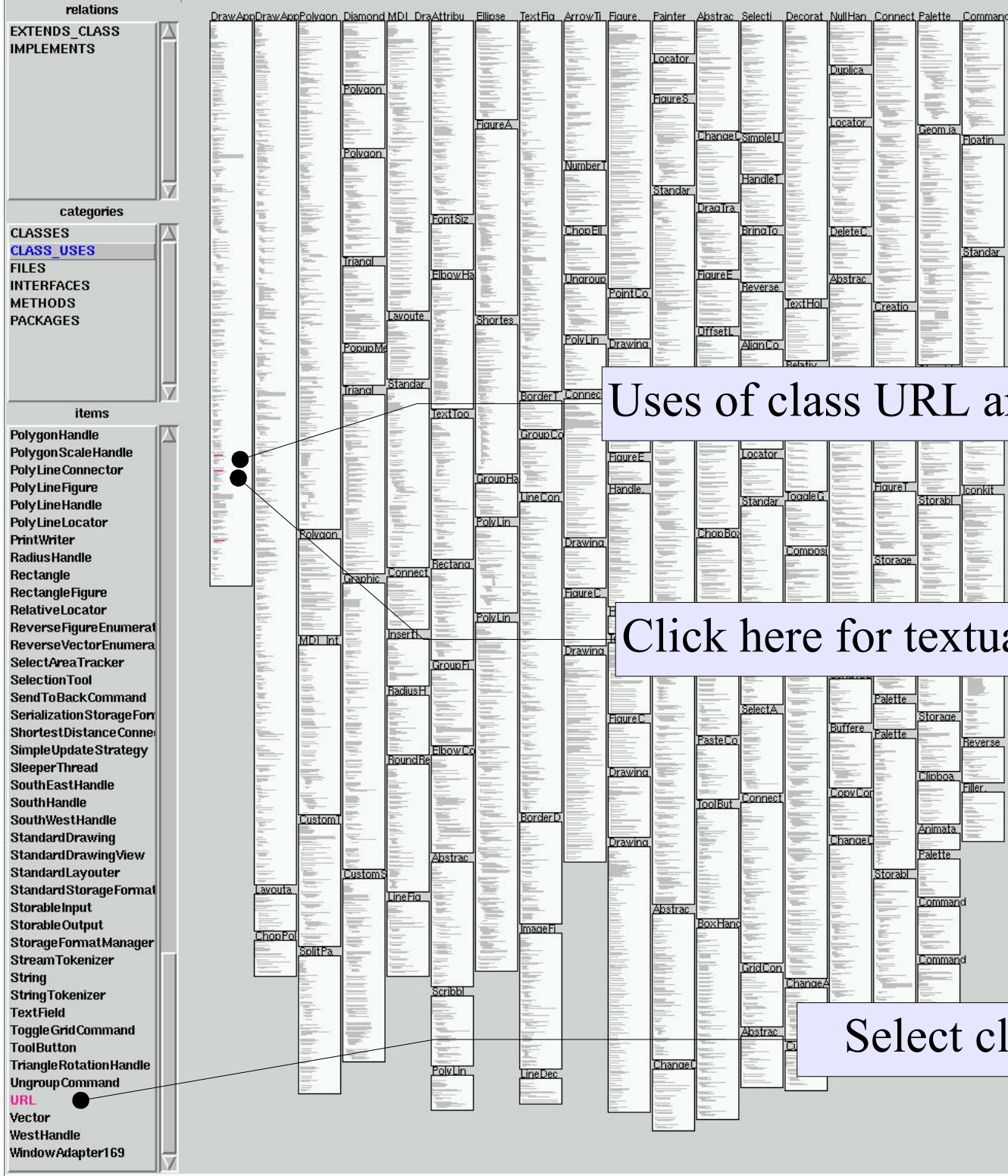


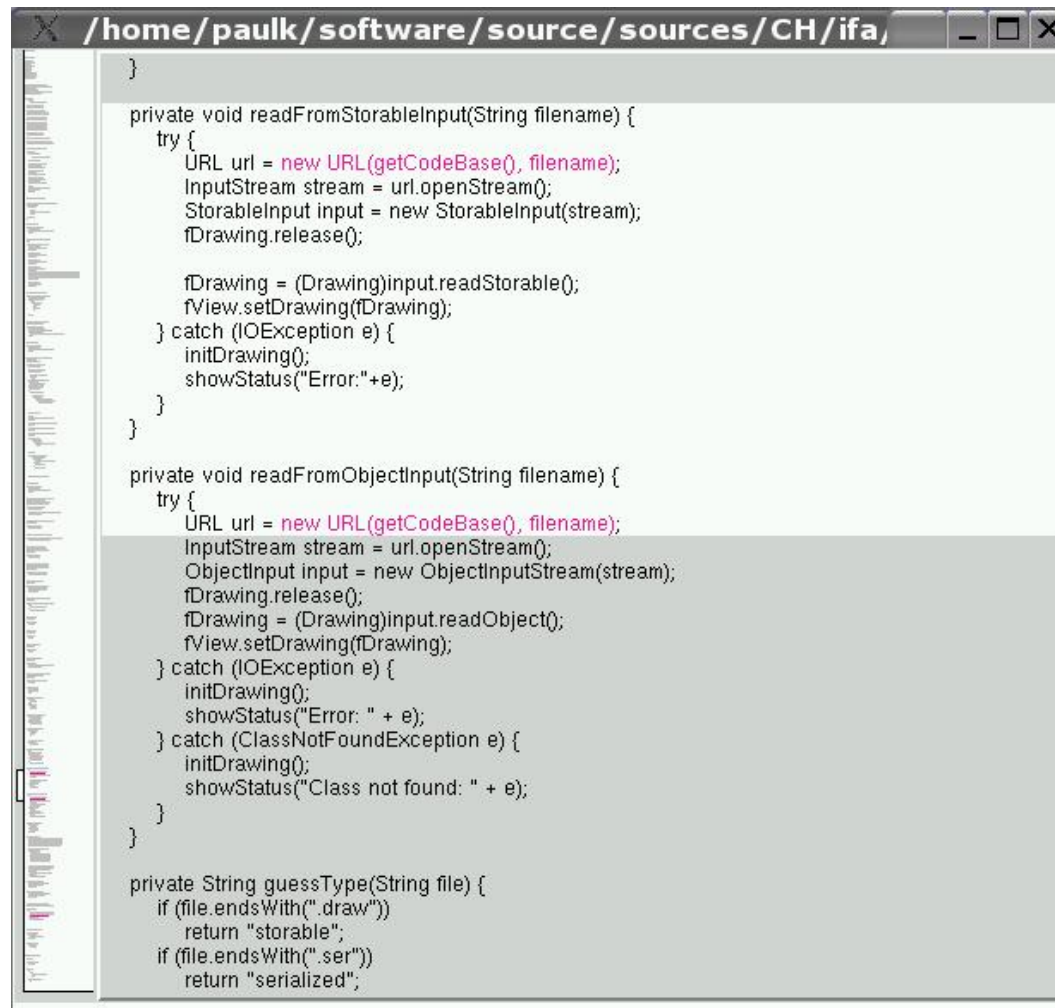


Hovering over file
shows
full name



Selecting a category displays all names in that category





The screenshot shows a code editor window with the title bar `/home/paulk/software/source/sources/CH/ifa,`. The editor contains the following Java code:

```
}  
  
private void readFromStorableInput(String filename) {  
    try {  
        URL url = new URL(getCodeBase(), filename);  
        InputStream stream = url.openStream();  
        StorableInput input = new StorableInput(stream);  
        fDrawing.release();  
  
        fDrawing = (Drawing)input.readStorable();  
        fView.setDrawing(fDrawing);  
    } catch (IOException e) {  
        initDrawing();  
        showStatus("Error:" + e);  
    }  
}  
  
private void readFromObjectInput(String filename) {  
    try {  
        URL url = new URL(getCodeBase(), filename);  
        InputStream stream = url.openStream();  
        ObjectInput input = new ObjectInputStream(stream);  
        fDrawing.release();  
        fDrawing = (Drawing)input.readObject();  
        fView.setDrawing(fDrawing);  
    } catch (IOException e) {  
        initDrawing();  
        showStatus("Error: " + e);  
    } catch (ClassNotFoundException e) {  
        initDrawing();  
        showStatus("Class not found: " + e);  
    }  
}  
  
private String guessType(String file) {  
    if (file.endsWith(".draw"))  
        return "storable";  
    if (file.endsWith(".ser"))  
        return "serialized";  
}
```

Wrap up: Rscript

- A simple, language-independent, relational calculus
- Fully typed
- Equation solver (\Rightarrow dataflow equations)
- Areas allow close link with source text
- Implementation: ASF+SDF
- IDE: **rscript-meta**
 - an instance of The Meta-Environment

Wrap up : Rscript

- Calls analysis
- Lifting of procedure calls to component relations
- Unitialized/unused variables
- McCabe & friends
- Clones in C code
- Dataflow analysis
 - reaching definitions
 - live variables
- Program slicing
- Java & ToolBus analysis
- Feature Descriptions/
package dependencies

Wrap up: visualization

- A **lot** of work to do but promising start
- Alternative pixel representations?
- Treemaps for directory structure of files?
- Colormaps for displaying metrics?
- Implementation: Tcl/Tk but may change to Swing
- Some simple visualizations are included in `rscrip-meta`

Further reading

- P. Klint, How understanding and restructuring differ from compiling: a rewriting approach, IWPC03
- P. Klint, A tutorial introduction to Rscript on www.meta-environment.org
- www.cwi.nl/~paulk/publications/all.html