

---

# Chapter 1. EASY Programming with Rascal

## *Leveraging the Extract-Analyze-SYthesize Paradigm for Meta-Programming*

Paul Klint, Jurgen Vinju, Tijs van der Storm

Mon Apr 20 13:18:47 CEST 2009

## Table of Contents

Introduction .....	2
EASY Programming .....	2
Rascal .....	4
Aim and Scope of this Book .....	5
Rascal Concepts .....	5
A Motivating Example .....	7
Preparations .....	7
Questions .....	9
The Rascal Language .....	11
Types and Values .....	11
List, Set, Map, Tuple, and Relation .....	11
User-defined Types .....	13
Typing .....	13
Expressions .....	14
Patterns .....	15
Comprehensions .....	18
Switch and Visit .....	19
Other control flow .....	19
Assertions .....	20
Declarations .....	21
Attributes .....	22
Larger Examples (TODO) .....	23
Analyzing the Component Structure of an Application .....	24
Analyzing the Structure of Java Systems .....	24
Finding Uninitialized and Unused Variables in a Program .....	26
Using Locations to Represent Program Fragments .....	27
McCabe Cyclomatic Complexity .....	29
Dataflow Analysis .....	29
Program Slicing .....	35
Built-in Operators and Functions .....	38
Benchmark .....	39
Boolean .....	39
Exception .....	40
Graph .....	40
Integer .....	41
IO .....	41
Labelled Graph .....	42
List .....	42

Location .....	44
Map .....	44
Node .....	45
Real .....	46
Relation .....	46
RSF .....	48
Resource (Eclipse only) .....	48
Set .....	48
String .....	51
Tuple .....	52
UnitTest .....	52
Value .....	53
ValueIO .....	53
View (Eclipse only) .....	53
Void .....	53
Extracting Facts from Source Code (TODO) .....	53
Workflow for Fact Extraction .....	54
Strategies for Fact Extraction .....	55
Fact Extraction using ASF+SDF .....	56
Concluding remarks .....	56
Table of Built-in Operators .....	56
Table of Built-in Functions .....	59
Bibliography .....	60

## Warning

This document is in the process of being written.

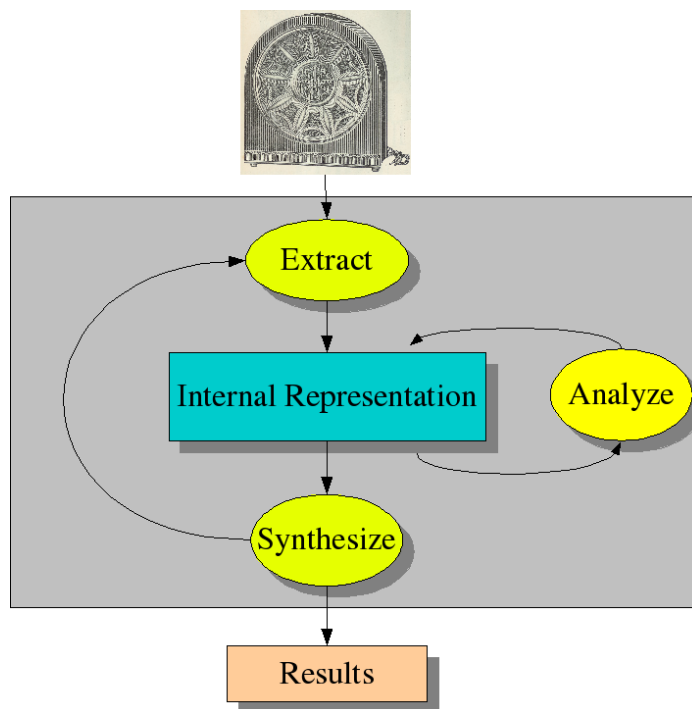
# Introduction

## EASY Programming

Many programming problems follow a fixed pattern. Starting with some *object-of-interest*, first relevant information is extracted from it and stored in an internal representation. This internal representation is then analyzed and used to synthesize results. If the synthesis indicates this, these steps can be repeated over and over again. These steps are shown in Figure 1.1, “EASY: the Extract-Analyze-Synthesize Paradigm” [3] (yes some people do like old radios!).

## Note

Find better images to illustrate this, e.g. in go a bicycle and a motor, and synthesize a motorbike.

**Figure 1.1. EASY: the Extract-Analyze-Synthesize Paradigm**

This is an abstract view on solving programming problems, but is it uncommon? No, so let's illustrate it with a few examples.

### Example 1.1. Finding security breaches

Alice is system administrator of a large online marketplace and she is looking for security breaches in her system. The object-of-interest are the system's log files. First relevant entries are extracted. This will include, for instance, messages from the SecureShell demon that reports failed login attempts. From each entry login name and originating IP address are extracted and put in a table (the internal representation in this example). These data are analysed by detecting duplicates and counting frequencies. Finally results are synthesized by listing the most frequently used login names and IP addresses.

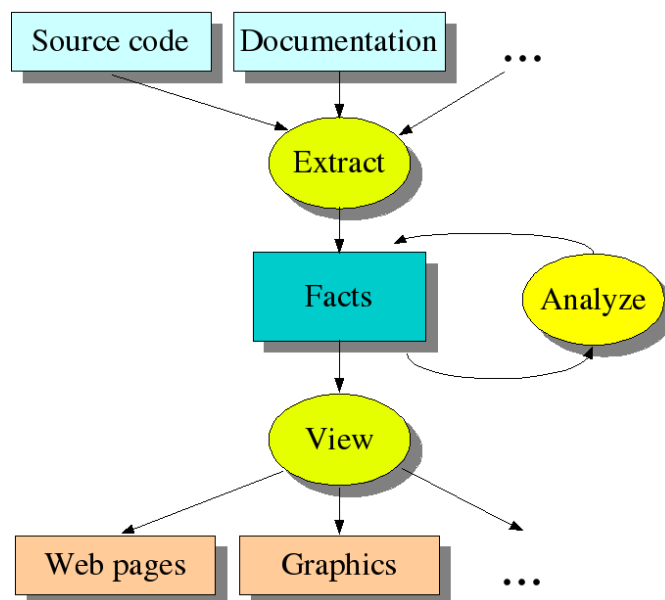
### Example 1.2. A Forensic DSL compiler

Bernd is a senior software engineer working at the Berlin headquarters of a forensic investigation lab of the German government. His daily work is to find common patterns in files stored on digital media that have been confiscated during criminal investigations. Text, audio and video files are stored in zillions of different data formats and each data format requires its own analysis technique. For each new investigation ad hoc combinations of tools are used. This makes the process very labour-intensive and error-prone. Bernd convinces his manager that designing a new domain-specific language (DSL) for forensic investigations may relieve the pressure on their lab. After designing the DSL---let's call it DERRICK---he makes an EASY implementation for it. Given a DERRICK program for a specific case under investigation, he first extracts relevant information from it and analyzes it: which media formats are relevant? Which patterns to look for? How should search results be combined? Given this new information, Java code is synthesized that uses the various existing tools and combines their results.

### Example 1.3. Renovation Financial Software

Charlotte is software engineer at a large financial institution in Paris and she is looking for options to connect an old and dusty software system to a web interface. She will need to analyze the sources of that system to understand how it can be changed to meet the new requirements. The objects-of-interest are in this case the source files, documentation, test scripts and any other available information. They have to be parsed in some way in order to extract relevant information, say the calls between various parts of the system. The call information can be represented as a binary relation between caller and callee (the internal representation in this example). This relation with 1-step calls is analyzed and further extended with 2-step calls, 3-step calls and so on. In this way call chains of arbitrary length become available. With this new information, we can synthesize results by determining the entry points of the software system, i.e. the points where calls from the outside world enter the system. Having completed this first cycle, Charlotte may be interested in which procedures can be called from the entry points and so on and so forth. Results will be typically represented as pictures that display the relationships that were found. In the case of source code analysis, a variation of our workflow scheme is quite common. It is then called the extract-analyze-view paradigm and is shown in Figure 1.2, “The extract-analyze-view paradigm” [4].

**Figure 1.2. The extract-analyze-view paradigm**



### Example 1.4. Finding Concurrency Errors

Daniel is concurrency researcher at one of the largest hardware manufacturers worldwide. He is working from an office in the Bay Area. Concurrency is the big issue for his company: it is becoming harder and harder to make CPUs faster, therefore more and of them are bundled on a single chip. Programming these multi-core chips is difficult and many programs that worked fine on a single CPU contain hard to detect concurrency errors due to subtle differences in the order of execution that results from executing the code on more than one CPU. Here is where Daniel enters the picture. He is working on tools for finding concurrency errors. First he extracts facts from the code that are relevant for concurrency problems and have to do with calls, threads, shared variables and locks. Next, he analyzes these facts and synthesizes an abstract model that captures the essentials of the concurrency behaviour of the program. Finally he runs a third-party verification tool with this model as input to do the actual verification.

## Rascal

With these examples in mind, you have a pretty good picture how EASY applies in different use cases. All these cases involve a form of *meta-programming*: software programs are the objects-of-interest

that are being analyzed and transformed. The Rascal language you are about to learn is designed for meta-programming following the EASY paradigm. It can be applied in domains ranging from compiler construction to constraint solving.

Since representation of data is central to the approach, Rascal provides a rich set of built-in data types. To support extraction and analysis, parsing and advanced pattern matching are provided. High-level control structures make analysis and synthesis of complex datastructures simple.

## Aim and Scope of this Book

## Rascal Concepts

Before explaining the Rascal language in more detail, it is good to have a general understanding of the concepts on which the language is built.

### Values

Rascal is a value-oriented language. This means that values are immutable and are always freshly constructed from existing parts and that sharing and aliasing problems are completely avoided. The language does provide assignment to variables either as the result of an explicit assignment statement or as the result of a successful match.

### Data structures

Rascal provides a rich set of datatypes. From booleans, infinite precision integers and reals to strings and source code locations. From lists, (optionally labelled) tuples and sets to (optionally labelled) maps and relations. From untyped tree structures to fully typed abstract datatypes. Syntax trees that are the result of parsing source files are represented as ADTs. There is a wealth of built-in operators and library functions available on the standard datatypes. A fragment of the abstract syntax for statements in a programming language would look as follows:

```
data STAT = asgStat(Id name, EXP exp)
           | ifStat(EXP exp, list[STAT] thenpart, list[STAT] elsepart)
           | whileStat(EXP exp, list[STAT] body)
           ;
```

### Pattern Matching

Pattern matching is *the* mechanism for case distinction in Rascal. We provide string matching based on regular expressions, list (associative) and set (associative, commutative, identity) matching, matching of abstract datatypes, and matching of concrete syntax patterns. All these forms of matching can be used in a single pattern. Patterns may contain variables that are bound when the match is successful. Anonymous (don't care) positions are indicated by an underscore (\_). The following abstract pattern matches the while statement defined above:

```
whileStat(EXP Exp, list[STAT] Stats)
```

Variables in a pattern are either explicitly declared in the pattern itself---as done in this example---or they may be declared in the context in which the pattern occurs. Patterns can also be used in an explicit match operator `:=` and can then be part of larger boolean expressions. Since a pattern match may have more than one solution, local backtracking over the alternatives of a match is provided.

### Generators

Generators enumerate the values in a given (finite) domain, be it the elements in a list, the substrings of a string, or all the nodes in a tree. Each value that is enumerated is first matched against a pattern before it can possibly contribute to the result of the generator. Examples are:

```
int x <- { 1, 3, 5, 7, 11 }  
int x <- [ 1 .. 10 ]  
asgStat(Id name, _) <- P
```

The first two generate the integer elements of a set of integers, respectively, a range of integers. The third generator traverses the complete program `P` (that is assumed to have a `PROGRAM` as value) and only yields statements that match the assignment pattern. Note the use of an anonymous variable at the `EXP` position in the pattern.

## Comprehensions and Control Structures

Rascal generalizes comprehensions in various ways. Comprehensions exist for lists, sets and maps. A comprehension consists of an expression that determines the successive elements to be included in the result and a list of generators and boolean expressions. The generators enumerate values and the boolean expressions filter them. A standard example is

```
{ x * x | int x <- [1 .. 10], x % 3 == 0 }
```

which returns the set `{9, 36, 81}`, i.e., the squares of the integers in the range `[ 1 .. 10 ]` that are divisible by 3. A more intriguing example is

```
{name | asgStat(Id name, _) <- P}
```

which returns a list of all identifiers that occur on the lefthand side of assignment statements in program `P`. Combinations of generators and boolean expressions also drive the control structures. For instance,

```
for(asgStat(Id name, _) <- P, size(Id) > 10){  
    println(Id);  
}
```

prints all identifiers in assignment statements that consist of more than 10 characters.

## Switching and Visiting

The switch statement as known from C and Java is generalized: the subject value to switch on may be an arbitrary value and the cases are arbitrary patterns. When a match fails, all its side-effects are undone and when it succeeds the statements associated with that case are executed. Visiting the elements of a datastructure is one of the most common operations in our domain and we give it first class support by way of visit expressions that resemble the switch statement. A visit expression consists of an expression that may yield an arbitrarily complex subject value and a number of cases. All the elements of the subject are visited and when one of the cases matches the statements associated with that case are executed. These cases may:

- cause some side effect;
- execute an `insert` statement that replaces the current element;
- execute a `fail` statement that causes the match for the current case to fail (and undoing all side-effects due to the successful match itself and the execution of the statements sofar).

The value of a visit expression is the original subject value with all replacements made as dictated by matching cases. The traversal order in a visit expressions can be explicitly defined by the programmer.

## Functions and Rewrite Rules

Functions are explicitly declared and are fully typed. Here is an example of a function that counts the number of assignment statements in a program:

```
int countAssignments(PROGRAM P){
```

```
int n = 0;
visit (P){
  case asgStat(Id name, _):
    n += 1;
}
return n;
}
```

Rewrite rules are the only implicit control mechanism in the language and are used to maintain invariants during computations.

## Typechecking

Rascal has a statically checked type system that prevents type errors and uninitialized variables at runtime. There are no runtime type casts as in Java and there are therefore less opportunities for run-time errors. The language provides higher-order, parametric, polymorphism. A type aliasing mechanism allows documenting specific uses of a type. Builtin operators are heavily overloaded: `1 + [2, 3]`, `[1, 2] + 3`, and `[1, 2] + [3]` all give the same result `[1, 2, 3]`, but `1 + [ "1" ]` leads to a type error.

## Execution

Following the what-you-see-is-what-you-get paradigm, control flow is completely explicit. Boolean expressions determine choices that drive the control structures. Rewrite rules form the only exception to the explicit control flow principle. Only local backtracking is provided (no surprise) in the context of boolean expressions and pattern matching; side effects are undone in case of backtracking.

# A Motivating Example

Suppose a mystery box ends up on your desk. When you open it, it contains a huge software system with several questions attached to it:

- How many procedure calls occur in this system?
- How many procedures contains it?
- What are the entry points for this system, i.e., procedures that call others but are not called themselves?
- What are the leaves of this application, i.e., procedures that are called but do not make any calls themselves?
- Which procedures call each other indirectly?
- Which procedures are called directly or indirectly from each entry point?
- Which procedures are called from all entry points?

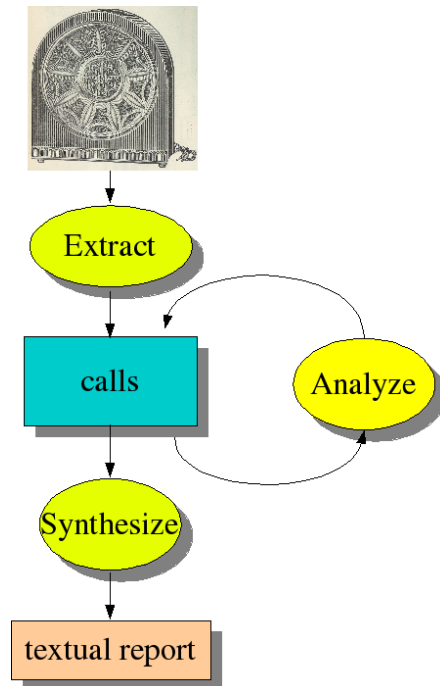
There are now two possibilities. Either you have this superb programming environment or tool suite that can immediately answer all these questions for you or you can use Rascal.

## Preparations

To illustrate this process consider the workflow in Figure 1.3, “Workflow for analyzing mystery box” [8]. First we have to extract the calls from the source code. Rascal is very at this, but to simplify this examples we assume that this call graph has already been extracted. Also keep in mind that a real call graph of a real application will contain thousands and thousands of calls. Drawing it in the way we do later on in Figure 1.4, “Graphical representation of the `calls` relation” [8] makes no sense

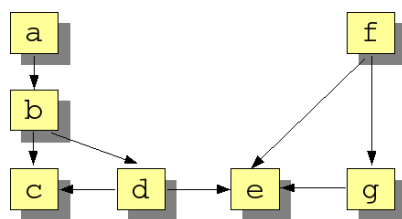
since we get a uniformly black picture due to all the call dependencies. After the extraction phase, we try to understand the extracted facts by writing queries to explore their properties. For instance, we may want to know *how many calls* there are, or *how many procedures*. We may also want to enrich these facts, for instance, by computing who calls who in more than one step. Finally, we produce a simple textual report giving answers to the questions we are interested in.

**Figure 1.3. Workflow for analyzing mystery box**



Now consider the call graph shown in Figure 1.4, “Graphical representation of the `calls` relation” [8]. This section is intended to give you a first impression what can be done with Rascal. Please return to this example when you have digested the detailed description of Rascal in the section called “*The Rascal Language*” [11] and the section called “*Built-in Operators and Functions*” [38].

**Figure 1.4. Graphical representation of the `calls` relation**



Rascal supports basic data types like integers and strings which are sufficient to formulate and answer the questions at hand. However, we can gain readability by introducing separately named types for the items we are describing. First, we introduce therefore a new type `proc` (an alias for strings) to denote procedures:

```
rascal> alias proc = str
done
```

Suppose that the following facts have been extracted from the source code and are represented by the relation `Calls`:

```
rascal> rel[proc , proc] Calls =
{ <"a", "b">, <"b", "c">, <"b", "d">, <"d", "c">, <"d", "e">,
```



```
<"f", "e">, <"f", "g">, <"g", "e">
}
```

This concludes the preparatory steps and now we move on to answer the questions.

## Questions

### How many procedure calls occur in this system?

To determine the numbers of calls, we simply determine the number of tuples in the `Calls` relation, as follows. First, we need the `Relation` library (described in the section called “*Relation*” [46]) so we import it:

```
rascal> import Relation;
done.
```

next we describe a new variable and calculate the number of tuples:

```
rascal> int nCalls = size(Calls)
```

The library function `size` determines the number of elements in a set or relation and is explained in the section called “*Relation*” [46]. In this example, `nCalls` will get the value 8.

### How many procedures contains it?

We get the number of procedures by determining which names occur in the tuples in the relation `Calls` and then determining the number of names:

```
rascal> set[proc] procs = carrier(Calls);
{"a", "b", "c", "d", "e", "f", "g"}
> int nprocs = size(procs);
7
```

The built-in function `carrier` determines all the values that occur in the tuples of a relation. In this case, `procs` will get the value `{"a", "b", "c", "d", "e", "f", "g"}` and `nprocs` will thus get value 7. A more concise way of expressing this would be to combine both steps:

```
rascal> int nprocs = size(carrier(Calls));
7
```

### What are the entry points for this system?

The next step in the analysis is to determine which *entry points* this application has, i.e., procedures which call others but are not called themselves. Entry points are useful since they define the external interface of a system and may also be used as guidance to split a system in parts. The `top` of a relation contains those left-hand sides of tuples in a relation that do not occur in any right-hand side. When a relation is viewed as a graph, its `top` corresponds to the root nodes of that graph. Similarly, the `bottom` of a relation corresponds to the leaf nodes of the graph. See the section called “*Graph*” [40] for more details. Using this knowledge, the entry points can be computed by determining the `top` of the `Calls` relation:

```
rascal> import Graph;
done.
rascal> set[proc] entryPoints = top(Calls);
{"a", "f"}
```

In this case, `entryPoints` is equal to `{"a", "f"}`. In other words, procedures `"a"` and `"f"` are the entry points of this application.

## What are the leaves of this application?

In a similar spirit, we can determine the *leaves* of this application, i.e., procedures that are being called but do not make any calls themselves:

```
rascal> set[proc] bottomCalls = bottom(Calls);  
{ "c", "e" }
```

In this case, `bottomCalls` is equal to `{ "c", "e" }`.

## Which procedures call each other indirectly?

We can also determine the *indirect calls* between procedures, by taking the transitive closure of the `Calls` relation:

```
rascal> rel[proc, proc] 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"> }
```

## Which procedures are called directly or indirectly from each entry point?

We know now the entry points for this application ("a" and "f") and the indirect call relations. Combining this information, we can determine which procedures are called from each entry point. This is done by indexing `closureCalls` with appropriate procedure name. The index operator yields all right-hand sides of tuples that have a given value as left-hand side. This gives the following:

```
rascal> set[proc] calledFromA = closureCalls["a"];  
{ "b", "c", "d", "e" }
```

and

```
rascal> set[proc] calledFromF = closureCalls["f"];  
{ "e", "g" }
```

## Which procedures are called from all entry points?

Finally, we can determine which procedures are called from both entry points by taking the intersection of the two sets `calledFromA` and `calledFromF`:

```
rascal> set[proc] commonProcs = calledFromA & calledFromF;  
{ "e" }
```

In other words, the procedures called from both entry points are mostly disjoint except for the common procedure "e".

## Wrap-up

These findings can be verified by inspecting a graph view of the calls relation as shown in Figure 1.4, “Graphical representation of the `calls` relation”[8]. Such a visual inspection does *not* scale very well to large graphs and this makes the above form of analysis particularly suited for studying large systems.

### Note

Add a screen dump here of the Eclipse browser.

# The Rascal Language

An Rascal program consists of one or more modules. Each module may import other modules and declare data types, variables and/or functions. The visibility of each declared entity may be `private` (only visible inside the module) or `public` (visible outside the module).

## Types and Values

### Elementary Types and Values

**Boolean.** The Booleans are represented by the type `bool` and have two values: `true` and `false`.

**Integer.** The integer values are represented by the type `int` and are written as usual, e.g., 0, 1, or 123. They can be arbitrarily large.

**Real.** The real values are represented by the type `real` and are written as usual, e.g., 1.5, or 3.14e-123. They can have arbitrary size and precision.

**String.** The string values are represented by the type `str` and consist of character sequences surrounded by double quotes. e.g., "a" or "a\ long\ string".

String literals permit *interpolation* of variable values: when `<X>` occurs inside a string literal, the value of the variable `X` is converted to string and replaces `<X>`. As a consequence, the character `<` has to be escaped as `\<` in string literals.

**Location.** Location values are represented by the type `loc` and serve as text coordinates in a specific source file. It is very handy to associate a source code location which extracted facts. Locations should *always* be generated automatically but for the curious here is an example how they look like: `loc(file:/home/paulk/pico.trm?offset=0&length=1&begin=2,3&end=4,5)`

Source locations have the following syntax:

```
loc(Url?offset=O&length=L&begin=BL,BC&end=EL,EC)
```

where:

- `Url` is an arbitrary URL.
- `O` and `L` are integer expressions giving the offset of this location to the begin of file, respectively, its length.
- `BL` and `BC` are integers expressions giving the begin line and begin column.
- `EL` and `EC` are integers expressions giving the end line and end column.

This element of a location value can be accessed and modified using the standard mechanism of field selection and field assignment. The corresponding field names are:

- `url`
- `offset`
- `length`
- `beginLine, beginColumn`
- `endLine, endColumn`.

## List, Set, Map, Tuple, and Relation

**List.** Lists are represented by the type `list[T]`, where `T` is an arbitrary type. Examples are `list[int]`, `list[tuple[int,int]]` and `list[list[str]]`. Lists are denoted by a list of

elements, separated by comma's and enclosed in bracket as in  $[E_1, E_2, \dots, E_n]$ , where the  $E_n$  ( $1 \leq i \leq n$ ) are expressions that yield the desired element type. For example,

- $[1, 2, 3]$  is of type `list[int]`,
- $\{<1,10>, <2,20>, <3,30>\}$  is of type `set[tuple[int,int]]`,
- $[<"a",10>, <"b",20>, <"c",30>]$  is of type `list[tuple[str,int]]`, and
- $[[ "a", "b"], [ "c", "d", "e"]]$  is of type `list[list[str]]`.

When variables of type list occur inside a list, their elements are automatically spliced into the surrounding list. This can be prevented by surrounding them with extra `[` and `]` brackets.

**Range.** For lists of integers, a special shorthand exists to describe ranges of integers:

- $[F \dots L]$  ranges from first element  $F$  to (and including) last element  $L$  with increments of 1.
- $[F, S, \dots E]$ , ranges from first element  $F$ , second element  $S$  to (and including) last element  $L$  with increments of  $S - F$ .

**Set.** Sets are represented by the type `set[T]`, where  $T$  is an arbitrary type. Examples are `set[int]`, `set[tuple[int,int]]` and `set[set[str]]`. Sets are denoted by a list of elements, separated by comma's and enclosed in braces as in  $\{E_1, E_2, \dots, E_n\}$ , where the  $E_n$  ( $1 \leq i \leq n$ ) are expressions that yield the desired element type. For example,

- $\{1, 2, 3\}$  is of type `set[int]`,
- $\{<1,10>, <2,20>, <3,30>\}$  is of type `set[tuple[int,int]]`,
- $\{<"a",10>, <"b",20>, <"c",30>\}$  is of type `set[tuple[str,int]]`, and
- $\{\{ "a", "b" \}, \{ "c", "d", "e" \}\}$  is of type `set[set[str]]`.

In a similar fashion as with lists, sets variables are automatically spliced into a surrounding set. This can be prevented by surrounding them with extra `{` and `}` brackets.

**Map.** Maps are represented by the type `map[T1, T2]`, where  $T_1$  and  $T_2$  are arbitrary types. Examples are `map[int,int]`, and `map[str,int]`. Maps are denoted by a list of pairs, separated by comma's and enclosed in parentheses as in  $(K_1: V_1, \dots, K_n: V_n)$ , where the  $K_n$  ( $1 \leq i \leq n$ ) are expressions that yield the keys of the map and  $V_n$  ( $1 \leq i \leq n$ ) are expressions that yield the values for each key. Map resemble functions rather than relations in the sense that only a single value can be associated with each key. For example,

- $(\text{"pear"} : 1, \text{"apple"} : 3, \text{"banana"} : 0)$  is of type `map[str,int]`.

**Tuple.** Tuples are represented by the type `tuple[T1 L1, T2 L2, ..., Tn Ln]`, where  $T_1, T_2, \dots, T_n$  are arbitrary types and  $L_1, L_2, \dots, L_n$  are optional labels. An example of a tuple type is `tuple[str name, int freq]`. Examples are:

- $<1, 2>$  is of type `tuple[int, int]`,
- $<1, 2, 3>$  is of type `tuple[int, int, int]`,
- $<"a", 3>$  is of type `tuple[str name, int freq]`.

**Relation.** Relations are nothing more than sets of tuples, but since they are used so often we provide a shorthand notation for them. Relations are represented by the type `rel[T1L1, T2L2, ..., TnLn]`, where  $T_1, T_2, \dots, T_n$  are arbitrary types and  $L_1, L_2, \dots, L_n$  are optional labels. It is a shorthand for `set[tuple[T1L1, T2L2, ..., TnLn]]`. Examples are `rel[int,str]` and `rel[int,set[str]]`. An  $n$ -ary relations with  $m$  tuples is denoted by  $\{<E_{11}, E_{12}, \dots, E_{1n}>, \dots, <E_{m1}, E_{m2}, \dots, E_{mn}>\}$ .

$\langle E_{21}, E_{22}, \dots, E_{2n} \rangle, \dots, \langle E_{m1}, E_{m2}, \dots, E_{mn} \rangle$ , where the  $E_{ij}$  are expressions that yield the desired element type. For example,  $\{\langle 1, "a" \rangle, \langle 2, "b" \rangle, \langle 3, "c" \rangle\}$  is of type `rel[int, str]`. Examples are:

- $\{\langle 1, 10 \rangle, \langle 2, 20 \rangle, \langle 3, 30 \rangle\}$  is of type `rel[int, int]` (yes indeed, you saw this same example before and then we gave `set[tuple[int, int]]` as its type; remember that these types are interchangeable.),
- $\{\langle "a", 10 \rangle, \langle "b", 20 \rangle, \langle "c", 30 \rangle\}$  is of type `rel[str, int]`, and
- $\{\{ "a", 1, "b" \}, \{ "c", 2, "d" \}\}$  is of type `rel[str, int, str]`.

## User-defined Types

**Alias type.** Everything can be expressed using the elementary types and values that are provided by Rascal. However, for the purpose of documentation and readability it is sometimes better to use a descriptive name as type indication, rather than an elementary type. The alias declaration

```
alias T1 = T2;
```

states that the new type name  $T_1$  can be used everywhere instead of the already defined type name  $T_2$ . For instance,

```
alias ModuleId = str;  
alias Frequency = int;
```

introduces two new type names `ModuleId` and `Frequency`, both an alias for the type `str`. The use of type aliases is a good way to document your intentions. Another example is an alias definition for a graph containing integer nodes:

```
alias Graph = rel[int, int];
```

An alias definition may be parameterized. So we can generalize graphs as follows:

```
alias Graph[&Node] = rel[&Node, &Node];
```

Of course, the type variables that are used in the type in the left part should occur as parameters in the right part of the definition and vice versa.

**Abstract Data Type.** In ordinary programming languages record types or classes exist to introduce a new type name for a collection of related, named, values and to provide access to the elements of such a collection through their name. In Rascal, data declarations provide this facility. The type declaration

```
data N = Pat1 | Pat2 | ...
```

introduces a new codatatype `N` and `Pat1`, `Pat2`, are prefix patterns describing the variants of the datatype. For instance,

```
data Bool = btrue | bfalse | band(Bool L, Bool R) | bor(Bool L, Bool R);
```

defines the datatype `Bool` that contains various constants and constructor functions.

## Typing

Rascal is based on static typing, this means that as many errors and inconsistencies as possible are spotted before the program is executed. TODO

# Expressions

**Table 1.1. Non-Boolean Operators**

Operator	Description
$exp_1$ [ $name = exp_2$ ]	$exp_1$ should evaluate to a datatype with a field $name$ ; assign value $exp_2$ to that field
$exp . name$	$exp$ should evaluate to a datatype with field $name$ ; return the value of that field
$exp < field, \dots >$	
$exp_1$ [ $exp_2$ ]	The value of $exp_2$ is used as index in $exp_1$ 's value. On string, list, tuple return the element with given index; for map return the value associated with $exp_2$ 's value
$- exp$	negation on integer or real
$exp +$	transitive closure on relation
$exp *$	reflexive transitive closure on relation
$exp @ name$	value of attribute $name$ of $exp$ 's value
$exp_1$ [ @ $name = exp_2$ ]	assign value of $exp_2$ to attribute $name$ of $exp_1$ 's value
$exp_1 \circ exp_2$	Composition on relations
$exp_1 / exp_2$	Division on integer and real
$exp_1 \% exp_2$	Modulo on integer
$exp_1 * exp_2$	Multiplication on integer and real; product on list, set, relation
$exp_1 \& exp_2$	Intersection on list, set, map and relation
$exp_1 + exp_2$	Addition on integer and real; concatenation on string, list and tuple; union on set, map, and relation
$exp_1 - exp_2$	Subtraction on integer and real; difference on list, set, map, and relation
$exp_1 \text{ join } exp_2$	Join on relation

**Table 1.2. Boolean Operators**

Operator	Description
<code>! exp</code>	Negation
<code>exp ?</code>	Isdefined: true is <code>exp</code> has a well-defined value
<code>exp<sub>1</sub> in exp<sub>2</sub></code>	Element of
<code>exp<sub>1</sub> not in exp<sub>2</sub></code>	Not element of
<code>exp<sub>1</sub> &lt;= exp<sub>2</sub></code>	Less than or equal on bool, int, real or string; sublist on list; subset on set, map or relation
<code>exp<sub>1</sub> &lt; exp<sub>2</sub></code>	Less than on bool, int, real or string; strict sublist on list; strict subset on set, map or relation
<code>exp<sub>1</sub> &gt;= exp<sub>2</sub></code>	Greater than or equal on bool, int, real or string; superlist on list; superset on set, map or relation
<code>exp<sub>1</sub> &gt; exp<sub>2</sub></code>	Greater than on bool, int, real or string; strict superlist on list; strict superset on set, map or relation
<code>pat := exp</code>	Value of <code>exp</code> matches with pattern <code>pat</code>
<code>pat != exp</code>	Value of <code>exp</code> does not match with pattern <code>pat</code>
<code>exp<sub>1</sub> == exp<sub>2</sub></code>	Equality
<code>exp<sub>1</sub> != exp<sub>2</sub></code>	Inequality
<code>exp<sub>1</sub> ? exp<sub>2</sub></code>	Ifdefined Otherwise: the value of <code>exp<sub>1</sub></code> is it is well-defined, otherwise the value of <code>exp<sub>2</sub></code>
<code>exp<sub>1</sub> ? exp<sub>2</sub> : exp<sub>3</sub></code>	Conditional expression
<code>exp<sub>1</sub> ==&gt; exp<sub>2</sub></code>	Implication
<code>exp<sub>1</sub> &lt;==&gt; exp<sub>2</sub></code>	Equivalence
<code>exp<sub>1</sub> &amp;&amp; exp<sub>2</sub></code>	Boolean and
<code>exp<sub>1</sub>    exp<sub>2</sub></code>	Boolean or

## Patterns

Patterns come in three flavours:

- *Regulars patterns* to do string matching with regular expressions, see the section called “*Regular Patterns*” [15].
- *Abstract patterns* to matching on arbitrary values, see the section called “*Abstract Patterns*” [16].
- *Concrete patterns* to match syntax trees that are the result of parsing, see the section called “*Concrete patterns*” [17].

## Regular Patterns

Regular expression patterns are ordinary regular expressions that are used to match a string value and to decompose it in parts and also to compose new strings. Regular expression patterns bind variables of type `str` when the match succeeds, otherwise they do not bind anything. Their syntax and semantics parallels abstract and concrete syntax patterns as much as possible. This means that they can occur in cases of `visit` and `switch` statements, on the left-hand side of the match operator (`:=` or `!:=`) and as declarator in generators.

We use a regular expression language that slightly extends the Java Regex language with the following exceptions:

- Regular expression are delimited by / and / optionally followed by a modifier (see below).
- We allow named groups, syntax `<Name : Regex>`, which introduce a variable of type `str` named `Name`. Currently, these names have to be unique in the pattern.
- Java regular expressions might have optional groups, which may introduce null bindings. Since uninitialized variables are not allowed in Rascal, we limit the kinds of expressions one can write here by not allowing nesting of named groups.
- Named groups have to be outermost, such that they can only bind in one way.
- Unlike Perl, Java uses the notation `(?Option)` inside the regular expression to set options like multi-line matching `(?m)`, case-insensitive matching `(?i)` etc. We let these options follow the regular expression.
- We allow name use in a regular expression: `<Name>` which inserts the string value of `Name` in the pattern.

Here are some examples of regular patterns.

```
/\brascal\b/i
```

does a case-insensitive match (`i`) of the word `rascal` between word boundaries (`\b`). And

```
/^.*?<word:\w+><rest:.*$>/m
```

does a multi-line match (`m`), matches the first consecutive word characters (`\w`) and assigns them to the variable `word`. The remainder of the string is assigned to the variable `rest`.

## Abstract Patterns

An abstract pattern is recursively defined and may contain the following elements:

- *Literal* of one of the basic types `bool`, `int`, `real`, `str`, or `loc`.
- A *variable declaration pattern*

```
Type Var
```

- A *variable pattern*

```
Var
```

- A *list pattern*

```
[ Pat1, Pat2, ..., Patn ]
```

- A *set pattern*

```
{ Pat1, Pat2, ..., Patn }
```

- A *tuple pattern*

```
< Pat1, Pat2, ..., Patn >
```

- A *node pattern*

```
name ( Pat1, Pat2, ..., Patn )
```

## Caution

Map patterns are currently not supported.



## Concrete patterns

### Caution

Concrete patterns are not yet implemented.

A *concrete pattern* is a (possibly quoted) concrete syntax fragment that may contain variables. We want to cover the whole spectrum from maximally quoted patterns that can unambiguously describe **any** syntax fragment to minimally quoted patterns as we are used to in ASF+SDF. A concrete pattern may have the following forms:

- A *variable declaration pattern*

```
<Type Var>
```

- A *variable pattern*

```
<Var>
```

- A *quoted pattern*

```
[ | Token1 Token2 ... Tokenn | ]
```

Inside a quoted pattern arbitrary lexical tokens may occur, but the characters <, > and | have to be escaped as \<, \>, \|. Quoted patterns may contain variable declaration patterns and variable patterns.

- A *typed quoted pattern*

```
Symbol [ | Token1 Token2 ... Tokenn | ]
```

is a quoted pattern that is preceded by an SDF symbol to define its desired syntactic type.

- An *unquoted pattern*

```
Token1 Token2 ... Tokenn
```

is a quoted pattern without the surrounding quotes.

- Inside syntax patterns, layout is ignored.

Examples:

- Quoted syntax pattern with two pattern variable declarations:

```
[ | while <EXP Exp> do <{STATEMENT ";"}* Stats> od | ]
```

- Quoted syntax pattern with two pattern variable uses:

```
[ | while <Exp> do <Stats> od | ]
```

- Identical to the previous example, but with a declaration of the desired syntactic type:

```
STATEMENT [ | while <Exp> do <Stats> od | ]
```

- Unquoted syntax pattern with two pattern variable declarations:

```
while <EXP Exp> do <{STATEMENT ";"}* Stats> od
```

- Unquoted syntax pattern with two pattern variable uses:

```
while <Exp> do <Stats> od
```

Obviously, with less quoting and type information, the probability of ambiguities increases. Our assumption is that a type checker can resolve most of them.

## Comprehensions

We will use the familiar notation for *list comprehension*

$$[E \mid G_1, \dots, G_n]$$

to denote the construction of a list consisting of the successive values of the *contributing expression*  $E$ . The values and the resulting list are determined by  $E$  and the *generators*  $G_1, \dots, G_n$ .  $E$  is computed for all possible combinations of values produced by the generators. Each generator may introduce new variables that can be used in subsequent generators as well as in the expressions  $E$ . A generator can use the variables introduced by preceding generators. Generators may enumerate all the values in a set or relation, or they may perform an arbitrary test.

In addition to list comprehensions, Rascal also supports *set comprehension*

$$\{E \mid G_1, \dots, G_n\}$$

that also serve as relation comprehension, provided that  $E$  is of a tuple type.

Finally, *map comprehensions* are written as:

$$(E_1 : E_2 \mid G_1, \dots, G_n)$$

Since the entries in a map require both a key and a value for each entry, two expressions are needed in this case.

## Enumerator

An enumerator generates all the values in a given list, set, map, tuple, relation or abstract datatype. They have the following form:

$$P \leftarrow E$$

where  $P$  is a pattern and  $E$  is an expression. An enumerator is evaluated as follows:

- Expression  $E$  is evaluated and may have an arbitrary value  $V$ .
- The elements of  $V$  are enumerated one by one.
- Each element value is matched against the pattern  $P$ . There are two cases:
  - The match succeeds, any variables in  $P$  are bound, and the next generator in the comprehension is evaluated. The variables that are introduced by an enumerator are only available to generators that appear later (i.e., to the right) in the comprehension. When this enumerator is the last generator in the comprehension its contributing expression is evaluated.
  - The match fails, no variables are bound. If  $V$  has more elements, a next element is tried. Otherwise, a previous generator (i.e., to the left) is tried. If this enumerator is the first generator in the comprehension, the evaluation of the comprehension is complete.

These are examples of enumerators:

- `int N <- {1, 2, 3, 4, 5},`
- `str K <- KEYWORDS`, where `KEYWORDS` should evaluate to a value of `set[str]`.
- `<str K, int N> <- {<"a",10>, <"b",20>, <"c",30>}`.
- `<str K, int N> <- FREQUENCIES`, where `FREQUENCIES` should evaluate to a value of `type rel[str,int]`.

- `<str K, 10> <- FREQUENCIES`, will only generate pairs with 10 as second element.

## Note

Type information will be used to check the plausibility of an enumerator and guard you against mistakes. An impossible enumerator like `int N <- {"apples", "oranges"}` will be flagged as an error since the pattern can never match.

## Note

An enumerator may be preceeded by a *strategy indication*:

- `top-down`
- `bottom-up` (this is the default)

These take only effect for enumerators that produce the elements of an abstract data type and determine the order in which the elements are enumerated.

## Test

A test is a boolean-valued expression. If the evaluation yields `true` this indicates that the current combination of generated values up to this test is still as desired and execution continues with subsequent generators. If the evaluation yields `false` this indicates that the current combination of values is undesired, and that another combination should be tried by going back to a previous generator.

Examples:

- `N >= 3` tests whether `N` has a value greater than or equal 3.
- `S == "coffee"` tests whether `S` is equal to the string `"coffee"`.

In both examples, the variable (`N`, respectively, `S`) should have been introduced by a generator that occurs earlier in the comprehension.

## Examples of Comprehensions

- `{X | int X : {1, 2, 3, 4, 5}, X >= 3}` yields the set `{3, 4, 5}`.
- `{<X, Y> | int X : {1, 2, 3}, int Y : {2, 3, 4}, X >= Y}` yields the relation `{<2, 2>, <3, 2>, <3, 3>}`.
- `{<Y, X> | <int X, int Y> : {<1, 10>, <2, 20>}}` yields the inverse of the given relation: `{<10, 1>, <20, 2>}`.
- `{X, X * X | X : {1, 2, 3, 4, 5}, X >= 3}` yields the set `{3, 4, 5, 9, 16, 25}`.

## Switch and Visit

## Other control flow

### While

### Do

### For

## One

## All

## Solve

It is also possible to define mutually dependent sets of equations:

```
equations
  initial
    T1 V1 init I1
    ...
    Tn Vn init In
  satisfy
    V1 = E1
    ...
    Vn = En
end equations
```

In the `initial` section, the variables  $V_i$  are declared and initialized. In the `satisfy` section, the actual set of equations is given. The expressions  $E_i$  may refer to any of the variables  $V_i$  (and to any variables declared earlier). This set of equations is solved by evaluating the expressions  $E_i$ , assigning their value to the corresponding variables  $V_i$ , and repeating this as long as the value of one of the variables was changed. This is typically used for solving a set of dataflow equations. Example:

- Although transitive closure is provided as a built-in operator, we can use equations to define the transitive closure of a relation. Recall that  $[R^+ = R \cup (R \circ R) \cup (R \circ R \circ R) \cup \dots]$ . This can be expressed as follows.

### Warning

Fix expression.

```
rel[int,int] R = {<1,2>, <2,3>, <3,4>}

equations
  initial
    rel[int,int] T init R
  satisfy
    T = T union (T o R)
end equations
```

The resulting value of `T` is as expected:

```
{<1,2>, <2,3>, <3,4>, <1, 3>, <2, 4>, <1, 4>}
```

## Assertions

An `assert` statement may occur everywhere where a declaration is allowed. It has the form

```
assert L: E
```

where  $L$  is a string that serves as a label for this assertion, and  $E$  is a boolean-value expression. During execution, a list of true and false assertions is maintained. When the script is executed as a *test suite* a summary of this information is shown to the user. When the script is executed in the standard fashion, the `assert` statement has no affect. Example:

- `assert "Equality on Sets 1": {1, 2, 3, 1} == {3, 2, 1, 1}`

## Declarations

### Variable Declarations

A variable declaration has the form

$$T \ V = E$$

where  $T$  is a type,  $V$  is a variable name, and  $E$  is an expression that should have type  $T$ . The effect is that the value of expression  $E$  is assigned to  $V$  and can be used later on as  $V$ 's value. Double declarations are not allowed. As a convenience, also declarations without an initialization expression are permitted and have the form

$$T \ V$$

and only introduce the variable  $V$ . Examples:

- `int max = 100` declares the integer variable `max` with value 100.
- The definition

```
rel[str,int] day = {<"mon", 1>, <"tue", 2>, <"wed",3>,
                  <"thu", 4>, <"fri", 5>, <"sat",6>, <"sun",7>}
```

declares the variable `day`, a relation that maps strings to integers.

### Local Variable Declarations

Local variables can be introduced as follows:

$$E \text{ where } T_1 \ V_1 = E_1, \dots, T_n \ V_n = E_n \text{ end where}$$

First the local variables  $V_i$  are bound to their respective values  $E_i$ , and then the value of expression  $E$  is yielded.

### Function Declarations

A function declaration has the form

$$T \ F(T_1 \ V_1, \dots, T_n \ V_n) = E$$

Here  $T$  is the result type of the function and this should be equal to the type of the associated expression  $E$ . Each  $T_i \ V_i$  represents a typed formal parameter of the function. The formal parameters may occur in  $E$  and get their value when  $F$  is invoked from another expression. Example:

- The function declaration

```
rel[int, int] invert(rel[int,int] R) =
  {<Y, X> | <int X, int Y> : R }
```

yields the inverse of the argument relation  $R$ . For instance, `invert({<1,10>, <2,20>})` yields `{<10,1>, <20,2>}`.

**Parameterized types in function declarations.** The types that occur in function declarations may also contain *type variables* that are written as `&` followed by an identifier. In this way functions can be defined for arbitrary types. Examples:

- The declaration

```
rel[&T2, &T1] invert2(rel[&T1,&T2] R) =
```

```
{<Y, X> | <&T1 X, &T2 Y> : R }
```

yields an inversion function that is applicable to any binary relation. For instance,

- `invert2({<1,10>, <2,20>})` yields `{<10,1>, <20,2>}`,
- `invert2({<"mon", 1>, <"tue", 2>})` yields `{<1, "mon">, <2, "tue">}`.
- The function

```
<&T2, &T1> swap(&T1 A, &T2 B) = <B, A>
```

can be used to swap the elements of pairs of arbitrary types. For instance,

- `swap(<1, 2>)` yields `<2,1>` and
- `swap(<"wed", 3>)` yields `<3, "wed">`.

## Attributes

Attributes are adornments of data and programs and come in two flavours:

- *node annotations* that allow associating one or more named values with nodes in a tree.
- *declaration tags* that allow associating one or more named values to a declaration in a Rascal program.

The former are intended to attach application data to values, like adding position information or control flow information to source code or adding visualization information to a relation. The latter are intended to add metadata to a Rascal program and allow to influence the execution of the Rascal program, for instance, by adding memoization hints or database mappings for relations.

## Node annotations

An annotation may be associated with any node value. An annotation has a name and the type of its value is explicitly declared. Any value of any named type can be annotated and the type of these annotations can be declared precisely.

For instance, we can add to certain syntactic constructs of programs (e.g., `EXPRESSION`) an annotation with name `posinfo` that contains location information:

```
anno loc EXPRESSION @ posinfo;
```

or location information could be added for all syntax trees:

```
anno loc node @ posinfo;
```

We can add to the graph datatype introduced earlier, the annotation with name `LayoutStrategy` that defines which graph layout algorithm to apply to a particular graph, e.g.,

```
data LayoutStrategy = "dot" | "tree" | "force" |  
                    "hierarchy" | "fisheye";
```

```
anno LayoutStrategy Graph @ strategy;
```

The following constructs are provided for handling annotations:

- `Val @ Anno`: get the value of annotation *Anno* of value *Val* (may be undefined!).
- `Val1[@Anno = Val2]`: set the value of annotation *Anno* of the value *Val1* to *Val2*.
- `Var @ Anno = Val`: set the value of annotation *Anno* of the value of variable *Var* to *Val*.

## Declaration tags

### Warning

Tags are not yet implemented.

### Warning

The syntax of tags has to be aligned with the syntax of annotations. This is done in the examples below but not yet in the syntax.

All declarations in a Rascal program may contain (in fixed positions depending on the declaration type) one or more declaration tags (`tag`). A tag is defined by declaring its name, the declaration type to which it can be attached, and the name and type of the annotation. The declaration type `all`, makes the declaration tag applicable for all possible declaration types. All declaration tags have the generic format `@Name{ ... }`, with arbitrary text between the brackets that is further constrained by the declared type. Here is an example of a license tag:

```
tag str license on module;
```

This will allow to write things like:

```
module Booleans
@license{This module is distributed under the GPL}
...
```

Other examples of declaration tags are:

```
tag str todo on all           %% a todo note for all declaration types
tag void deprecated on function %% marks a deprecated function
tag int memo on function      %% bounded memoization of
                             %% function calls
tag str doc on all            %% documentation string
tag str primitive on function %% a primitive, built-in, function
```

Here is an example of a documentation string as used in the Rascal standard library:

```
public &T max(set[&T] R)
  @doc{Maximum of a set: max}
{
  &T result = arb(R);
  for(&T E : R){
    result = max(result, E);
  }
  return result;
}
```

## Larger Examples (TODO)

Now we will have a closer look at some larger applications of Rascal. We start by analyzing the global structure of a software system. You may now want to reread the example of call graph analysis given earlier in the section called “*A Motivating Example*” [7] as a reminder. The component structure of an application is analyzed in the section called “*Analyzing the Component Structure of an Application*” [24] and Java systems are analyzed in the section called “*Analyzing the Structure of Java Systems*” [24]. Next we move on to the detection of initialized variables in the section called “*Finding Uninitialized and Unused Variables in a Program*” [26] and we explain how source code locations can be included in a such an analysis (the section called “*Using Locations to Represent Program Fragments*” [27]). As an example of computing code metrics, we describe the calculation of McCabe’s cyclomatic complexity in the section called “*M McCabe Cyclomatic Complexity*”. Several examples of dataflow analysis follow in the section called “*Dataflow Analysis*” [29]. A description of program slicing concludes the chapter (the section called “*Program Slicing*” [35]).

## Analyzing the Component Structure of an Application

A frequently occurring problem is that we know the call relation of a system but that we want to understand it at the component level rather than at the procedure level. If it is known to which component each procedure belongs, it is possible to *lift* the call relation to the component level as proposed in [Kri99]. First, introduce new types to denote procedure calls as well as components of a system:

```
type proc = str
type comp = str
```

Given a calls relation `Calls2`, the next step is to define the components of the system and to define a `PartOf` relation between procedures and components.

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

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

rel[proc, comp] PartOf = {<"main", "Appl">, <"a", "Appl">,
                        <"b", "DB">, <"c", "Lib">, <"d", "Lib">
                      }
```

Actual lifting, amounts to translating each call between procedures by a call between components. This is achieved by the following function `lift`:

```
rel[comp,comp] lift(rel[proc,proc] aCalls, rel[proc,comp] aPartOf)=
  { <C1, C2> | <proc P1, proc P2> : aCalls,
              <comp C1, comp C2> : aPartOf[P1] x aPartOf[P2]
  }
```

In our example, the lifted call relation between components is obtained by

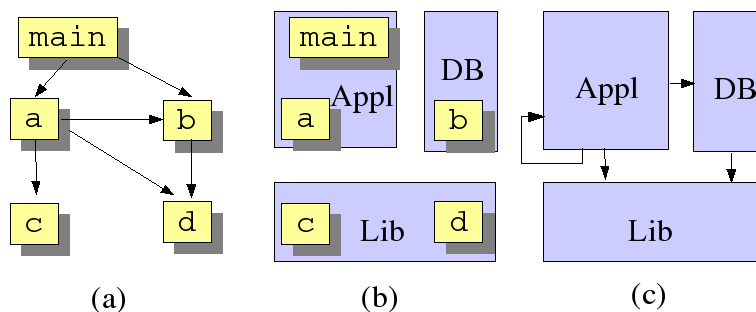
```
rel[comp,comp] ComponentCalls = lift(Calls2, PartOf)
```

and has as value:

```
{<"DB", "Lib">, <"Appl", "Lib">, <"Appl", "DB">, <"Appl", "Appl">}
```

The relevant relations for this example are shown in Figure 1.5, “(a) `Calls2`; (b) `PartOf`; (c) `ComponentCalls`.” [24].

**Figure 1.5. (a) `Calls2`; (b) `PartOf`; (c) `ComponentCalls`.**



## Analyzing the Structure of Java Systems

Now we consider the analysis of Java systems (inspired by [BNL03]). Suppose that the type `class` is defined as follows



```
type class = str
```

and that the following relations are available about a Java application:

- `rel[class, class] CALL`: If  $\langle C_1, C_2 \rangle$  is an element of `CALL`, then some method of  $C_2$  is called from  $C_1$ .
- `rel[class, class] INHERITANCE`: If  $\langle C_1, C_2 \rangle$  is an element of `INHERITANCE`, then class  $C_1$  either extends class  $C_2$  or  $C_1$  implements interface  $C_2$ .
- `rel[class, class] CONTAINMENT`: If  $\langle C_1, C_2 \rangle$  is an element of `CONTAINMENT`, then one of the attributes of class  $C_1$  is of type  $C_2$ .

To make this more explicit, consider the class `LocatorHandle` from the `JHotDraw` application (version 5.2) as shown here:

```
package CH.ifa.draw.standard;

import java.awt.Point;
import CH.ifa.draw.framework.*;
/**
 * A LocatorHandle implements a Handle by delegating the
 * location requests to a Locator object.
 */
public class LocatorHandle extends AbstractHandle {
    private Locator    fLocator;
    /**
     * Initializes the LocatorHandle with the given Locator.
     */
    public LocatorHandle(Figure owner, Locator l) {
        super(owner);
        fLocator = l;
    }
    /**
     * Locates the handle on the figure by forwarding the request
     * to its figure.
     */
    public Point locate() {
        return fLocator.locate(owner());
    }
}
```

It leads to the addition to the above relations of the following tuples:

- To `CALL` the pairs  $\langle \text{"LocatorHandle"}, \text{"AbstractHandle"} \rangle$  and  $\langle \text{"LocatorHandle"}, \text{"Locator"} \rangle$  will be added.
- To `INHERITANCE` the pair  $\langle \text{"LocatorHandle"}, \text{"AbstractHandle"} \rangle$  will be added.
- To `CONTAINMENT` the pair  $\langle \text{"LocatorHandle"}, \text{"Locator"} \rangle$  will be added.

Cyclic structures in object-oriented systems makes understanding hard. Therefore it is interesting to spot classes that occur as part of a cyclic dependency. Here we determine cyclic uses of classes that include calls, inheritance and containment. This is achieved as follows:

```
rel[class, class] USE = CALL union CONTAINMENT union INHERITANCE
set[str] ClassesInCycle =
    {C1 | <class C1, class C2> : USE+, C1 == C2}
```

First, we define the `USE` relation as the union of the three available relations `CALL`, `CONTAINMENT` and `INHERITANCE`. Next, we consider all pairs  $\langle C_1, C_2 \rangle$  in the transitive closure of the `USE` relation

such that  $C_1$  and  $C_2$  are equal. Those are precisely the cases of a class with a cyclic dependency on itself. Probably, we do not only want to know which classes occur in a cyclic dependency, but we also want to know which classes are involved in such a cycle. In other words, we want to associate with each class a set of classes that are responsible for the cyclic dependency. This can be done as follows.

```
rel[class,class] USE = CALL union CONTAINMENT union INHERITANCE
set[class] CLASSES = carrier(USE)
rel[class,class] USETRANS = USE+
rel[class,set[class]] ClassCycles =
  {<C, USETRANS[C]> | class C : CLASSES, <C, C> in USETRANS }
```

First, we introduce two new shorthands: `CLASSES` and `USETRANS`. Next, we consider all classes  $C$  with a cyclic dependency and add the pair  $\langle C, \text{USETRANS}[C] \rangle$  to the relation `ClassCycles`. Note that `USETRANS[C]` is the right image of the relation `USETRANS` for element  $C$ , i.e., all classes that can be called transitively from class  $C$ .

## Finding Uninitialized and Unused Variables in a Program

Consider the following program in the toy language Pico: (This is an extended version of the example presented earlier in [Kli03].)

```
[ 1] begin declare x : natural, y : natural,
[ 2]             z : natural, p : natural;
[ 3]   x := 3;
[ 4]   p := 4;
[ 5]   if q then
[ 6]     z := y + x
[ 7]   else
[ 8]     x := 4
[ 9]   fi;
[10]   y := z
[11] end
```

Inspection of this program learns that some of the variables are being used before they have been initialized. The variables in question are  $q$  (line 5),  $y$  (line 6), and  $z$  (line 10). It is also clear that variable  $p$  is initialized (line 4), but is never used. How can we automate these kinds of analysis? Recall from the section called “*EASY Programming*”[2] that we follow Extract-Analyze-SYnthesize paradigm to approach such a problem. The first step is to determine which elementary facts we need about the program. For this and many other kinds of program analysis, we need at least the following:

- The *control flow graph* of the program. We represent it by a relation `PRED` (for predecessor) which relates each statement with each predecessors.
- The *definitions* of each variable, i.e., the program statements where a value is assigned to the variable. It is represented by the relation `DEFS`.
- The *uses* of each variable, i.e., the program statements where the value of the variable is used. It is represented by the relation `USES`.

In this example, we will use line numbers to identify the statements in the program. (In the section called “*Using Locations to Represent Program Fragments*”[27], we will use locations to represent statements.) Assuming that there is a tool to extract the above information from a program text, we get the following for the above example:

```
type expr = int
type varname = str
expr ROOT = 1
rel[expr,expr] PRED = { <1,3>, <3,4>, <4,5>, <5,6>, <5,8>,
```

```

                                <6,10>, <8,10>
                                }
rel[expr,varname] DEFS = { <3,"x">, <4,"p">, <6,"z">,
                           <8,"x">, <10,"y">
                           }
rel[expr,varname] USES = { <5,"q">, <6,"y">, <6,"x">, <10,"z"> }

```

This concludes the extraction phase. Next, we have to enrich these basic facts to obtain the initialized variables in the program. So, when is a variable  $V$  in some statement  $S$  initialized? If we execute the program (starting in  $ROOT$ ), there may be several possible execution path that can reach statement  $S$ . All is well if *all* these execution path contain a definition of  $V$ . However, if one or more of these path do *not* contain a definition of  $V$ , then  $V$  may be uninitialized in statement  $S$ . This can be formalized as follows:

```

rel[expr,varname] UNINIT =
  { <E, V> | <expr E, varname V>: USES,
        E in reachX({ROOT}, DEFS[-,V], PRED)
  }

```

We analyze this definition in detail:

- $\langle \text{expr } E, \text{varname } V \rangle : \text{USES}$  enumerates all tuples in the  $\text{USES}$  relation. In other words, we consider the use of each variable in turn.
- $E \text{ in } \text{reachX}(\{\text{ROOT}\}, \text{DEFS}[-,V], \text{PRED})$  is a test that determines whether statement  $S$  is reachable from the  $ROOT$  without encountering a definition of variable  $V$ .
  - $\{\text{ROOT}\}$  represents the initial set of nodes from which all path should start.
  - $\text{DEFS}[-,V]$  yields the set of all statements in which a definition of variable  $V$  occurs. These nodes form the exclusion set for  $\text{reachX}$ : no path will be extended beyond an element in this set.
  - $\text{PRED}$  is the relation for which the reachability has to be determined.
  - The result of  $\text{reachX}(\{\text{ROOT}\}, \text{DEFS}[-,V], \text{PRED})$  is a set that contains all nodes that are reachable from the  $ROOT$  (as well as all intermediate nodes on each path).
  - Finally,  $E \text{ in } \text{reachX}(\{\text{ROOT}\}, \text{DEFS}[-,V], \text{PRED})$  tests whether expression  $E$  can be reached from the  $ROOT$ .
- The net effect is that  $\text{UNINIT}$  will only contain pairs that satisfy the test just described.

When we execute the resulting Rascal code (i.e., the declarations of  $ROOT$ ,  $\text{PRED}$ ,  $\text{DEFS}$ ,  $\text{USES}$  and  $\text{UNINIT}$ ), we get as value for  $\text{UNINIT}$ :

```
{<5, "q">, <6, "y">, <10, "z">}
```

and this is in concordance with the informal analysis given at the beginning of this example.

As a bonus, we can also determine the *unused* variables in a program, i.e., variables that are defined but are used nowhere. This is done as follows:

```
set[var] UNUSED = range(DEFS) \ range(USES)
```

Taking the range of the relations  $\text{DEFS}$  and  $\text{USES}$  yields the variables that are defined, respectively, used in the program. The difference of these two sets yields the unused variables, in this case  $\{\text{"p"}\}$ .

## Using Locations to Represent Program Fragments

### Warning

Fix the following

```
\begin{figure}[tb] \begin{center} \epsfig{figure=figs/meta-pico.eps,width=6cm} \hspace*{0.5cm}
\epsfig{figure=figs/pico-example.eps,width=6cm} \end{center} \hrulefill \caption{\label{FIG:meta-
pico}Checking undefined variables in Pico programs using the ASF+SDF Meta-Environment. On the
left, main window of Meta-Environment with error messages related to Pico program shown on the
right.{\bf THIS FIGURE IS OUTDATED}} \end{figure}
```

One aspect of the example we have just seen is artificial: where do these line numbers come from that we used to indicate expressions in the program? One solution is to let the extraction phase generate *locations* to precisely indicate relevant places in the program text. Recall from the section called “*Elementary Types and Values*” [11], that a location consists of a file name, a begin line, a begin position, an end line, and an end position. Also recall that locations can be compared: a location  $A_1$  is smaller than another location  $A_2$ , if  $A_1$  is textually contained in  $A_2$ . By including locations in the final answer of a relational expression, external tools will be able to highlight places of interest in the source text.

The first step, is to define the type `expr` as aliases for `loc` (instead of `int`):

```
type expr = loc
type varname = str
```

Of course, the actual relations are now represented differently. For instance, the `USES` relation may now look like

```
{ <area-in-file("/home/paulk/example.pico",
               area(5,5,5,6,106,1)), "q">,
  <area-in-file("/home/paulk/example.pico",
               area(6,13,6,14,127,1)), "y">,
  <area-in-file("/home/paulk/example.pico",
               area(6,17,6,18,131,1)), "x">,
  <area-in-file("/home/paulk/example.pico",
               area(10,7,10,8,168,1)), "z">
}
```

The definition of `UNINIT` can be nearly reused as is. The only thing that remains to be changed is to map the (expression, variable-name) tuples to (variable-name, variable-occurrence) tuples, for the benefit of the precise highlighting of the relevant variables. We define a new type `var` to represent variable occurrences and an auxiliary set that `VARNAMES` that contains all variable names:

```
type var = loc
set[varname] VARNAMES = range(DEFS) union range(USES)
```

Remains the new definition of `UNINIT`:

```
rel[var, varname] UNINIT =
  { <V, VN> | var-name VN : VARNAMES,
              var V : USES[- ,VN],
              expr E : reachX({ROOT}, DEFS[- ,VN], PRED),
              V <= E
  }
```

This definition can be understood as follows:

- `var-name VN : VARNAMES` generates all variable names.
- `var V : USES[- ,VN]` generates all variable uses `V` for variables with name `VN`.
- As before, `expr E : reachX({ROOT}, DEFS[- ,VN], PRED)` generates all expressions `E` that can be reached from the start of the program without encountering a definition for variables named `VN`.
- `V <= E` tests whether variable use `V` is enclosed in that expression (using a comparison on locations). If so, we have found an uninitialized occurrence of the variable named `VN`.

## Warning

Fix reference

In Figure~\ref{FIG:meta-pico} it is shown how checking of Pico programs in the ASF+SDF Meta-Environment looks like.

## McCabe Cyclomatic Complexity

The *cyclomatic complexity* of a program is defined as  $e - n + 2$ , where  $e$  and  $n$  are the number of edges and nodes in the control flow graph, respectively. It was proposed by McCabe [McC76] as a measure of program complexity. Experiments have shown that programs with a higher cyclomatic complexity are more difficult to understand and test and have more errors. It is generally accepted that a program, module or procedure with a cyclomatic complexity larger than 15 is *too complex*. Essentially, cyclomatic complexity measures the number of decision points in a program and can be computed by counting all if statement, case branches in switch statements and the number of conditional loops. Given a control flow in the form of a predecessor relation `rel[stat,stat]` PRED between statements, the cyclomatic complexity can be computed in an Rascal as follows:

```
int cyclomatic-complexity(rel[stat,stat] PRED) =
    #PRED - #carrier(PRED) + 2
```

The number of edges  $e$  is equal to the number of tuples in PRED. The number of nodes  $n$  is equal to the number of elements in the carrier of PRED, i.e., all elements that occur in a tuple in PRED.

## Dataflow Analysis

*Dataflow analysis* is a program analysis technique that forms the basis for many compiler optimizations. It is described in any text book on compiler construction, e.g. [ASU86]. The goal of dataflow analysis is to determine the effect of statements on their surroundings. Typical examples are:

- Dominators (the section called “*Dominators*”[29]): which nodes in the flow dominate the execution of other nodes?
- Reaching definitions (the section called “*Reaching Definitions*”[31]): which definitions of variables are still valid at each statement?
- Live variables (the section called “*Live Variables*”[34]): of which variables will the values be used by successors of a statement?
- Available expressions: an expression is available if it is computed along each path from the start of the program to the current statement.

## Dominators

A node  $d$  of a flow graph *dominates* a node  $n$ , if every path from the initial node of the flow graph to  $n$  goes through  $d$  [ASU86] (Section 10.4). Dominators play a role in the analysis of conditional statements and loops. The function `dominators` that computes the dominators for a given flow graph PRED and an entry node ROOT is defined as follows:

```
rel[stat,stat] dominators(rel[stat,stat] PRED, int ROOT) =
    DOMINATES
where
    set[int] VERTICES = carrier(PRED)

    rel[int,set[int]] DOMINATES =
        { <V, VERTICES \ {V, ROOT} \ reachX({ROOT}, {V}, PRED)> |
          int V : VERTICES }
endwhere
```

First, the auxiliary set VERTICES (all the statements) is computed. The relation DOMINATES consists of all pairs  $\langle S, \{S_1, \dots, S_n\} \rangle$  such that

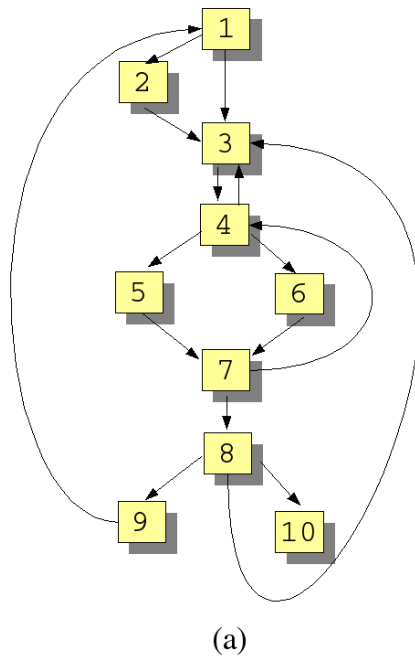
- $S_i$  is not an initial node or equal to  $S$ .
- $S_i$  cannot be reached from the initial node without going through  $S$ .

Consider the flow graph

```
rel[int,int] PRED = {
<1,2>, <1,3>,
<2,3>,
<3,4>,
<4,3>,<4,5>, <4,6>,
<5,7>,
<6,7>,
<7,4>,<7,8>,
<8,9>,<8,10>,<8,3>,
<9,1>,
<10,7>
}
```

It is illustrated in Figure 1.6, “Flow graph” [30]

**Figure 1.6. Flow graph**

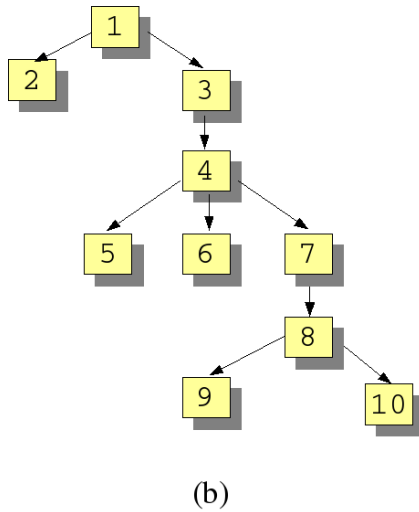


The result of applying dominators to it is as follows:

```
{<1, {2, 3, 4, 5, 6, 7, 8, 9, 10}>,
<2, {}>,
<3, {4, 5, 6, 7, 8, 9, 10}>,
<4, {5, 6, 7, 8, 9, 10}>,
<5, {}>,
<6, {}>,
<7, {8, 9, 10}>,
<8, {9, 10}>,
<9, {}>,
<10, {}>}
```

The resulting *dominator tree* is shown in Figure 1.7, “Dominator tree”[31]. The dominator tree has the initial node as root and each node  $d$  in the tree only dominates its descendants in the tree.

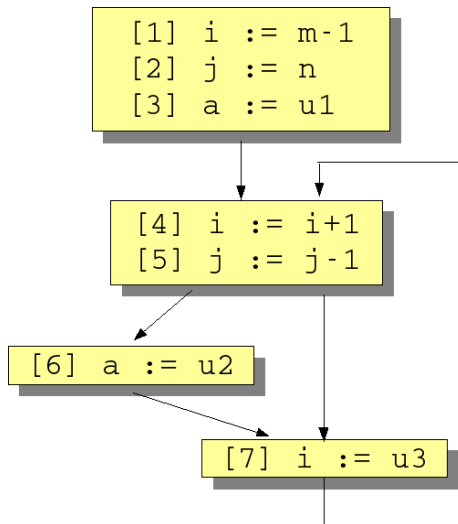
**Figure 1.7. Dominator tree**



## Reaching Definitions

We illustrate the calculation of reaching definitions using the example in Figure 1.8, “Flow graph for various dataflow problems” [31] which was inspired by [ASU86] (Example 10.15).

**Figure 1.8. Flow graph for various dataflow problems**



As before, we assume the following basic relations PRED, DEFS and USES about the program:

```

type stat = int
type var = str
rel[stat,stat] PRED = { <1,2>, <2,3>, <3,4>, <4,5>, <5,6>,
                        <5,7>, <6,7>, <7,4>
                      }
rel[stat, var] DEFS = { <1, "i">, <2, "j">, <3, "a">, <4, "i">,
                        <5, "j">, <6, "a">, <7, "i">
                      }
rel[stat, var] USES = { <1, "m">, <2, "n">, <3, "u1">, <4, "i">,
                        <5, "j">, <6, "u2">, <7, "u3">
                      }

```

```
}
```

For convenience, we introduce a notion `def` that describes that a certain statement defines some variable and we revamp the basic relations into a more convenient format using this new type:

```
type def = <stat theStat, var theVar>
```

```
rel[stat, def] DEF = {<S, <S, V>> | <stat S, var V> : DEFS}
```

```
rel[stat, def] USE = {<S, <S, V>> | <stat S, var V> : USES}
```

The new DEF relation gets as value:

```
{ <1, <1, "i">>, <2, <2, "j">>, <3, <3, "a">>, <4, <4, "i">>,
  <5, <5, "j">>, <6, <6, "a">>, <7, <7, "i">>
}
```

and USE gets as value:

```
{ <1, <1, "m">>, <2, <2, "n">>, <3, <3, "u1">>, <4, <4, "i">>,
  <5, <5, "j">>, <6, <6, "u2">>, <7, <7, "u3">>
}
```

Now we are ready to define an important new relation KILL. KILL defines which variable definitions are undone (killed) at each statement and is defined as follows:

```
rel[stat, def] KILL =
  {<S1, <S2, V>> | <stat S1, var V> : DEFS,
                    <stat S2, V> : DEFS,
                    S1 != S2
  }
```

In this definition, all variable definitions are compared with each other, and for each variable definition all *other* definitions of the same variable are placed in its kill set. In the example, KILL gets the value

```
{ <1, <4, "i">>, <1, <7, "i">>, <2, <5, "j">>, <3, <6, "a">>,
  <4, <1, "i">>, <4, <7, "i">>, <5, <2, "j">>, <6, <3, "a">>,
  <7, <1, "i">>, <7, <4, "i">>
}
```

and, for instance, the definition of variable `i` in statement 1 kills the definitions of `i` in statements 4 and 7. Next, we introduce the collection of statements

```
set[stat] STATEMENTS = carrier(PRED)
```

which gets as value  $\{1, 2, 3, 4, 5, 6, 7\}$  and two convenience functions to obtain the predecessor, respectively, the successor of a statement:

```
set[stat] predecessor(stat S) = PRED[-, S]
set[stat] successor(stat S) = PRED[S, -]
```

After these preparations, we are ready to formulate the reaching definitions problem in terms of two relations IN and OUT. IN captures all the variable definitions that are valid at the entry of each statement and OUT captures the definitions that are still valid after execution of each statement. Intuitively, for each statement `S`, `IN[S]` is equal to the union of the OUT of all the predecessors of `S`. `OUT[S]`, on the other hand, is equal to the definitions generated by `S` to which we add `IN[S]` minus the definitions that are killed in `S`. Mathematically, the following set of equations captures this idea for each statement:

## Warning

Fix expression

$$[ \text{IN}[S] = \bigcup \{ \text{OUT}[P] \mid P \text{ in predecessor of } S \} \text{ OUT}[S] \setminus \text{KILL}[S] ]$$



$$\setminus [OUT[S] = DEF[S] \setminus \cup (IN[S] - KILL[S]) \setminus ]$$

This idea can be expressed in Rascal quite literally:

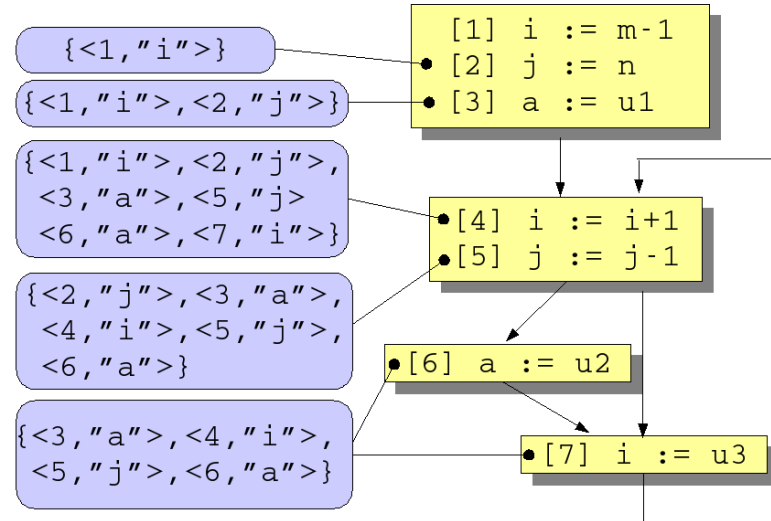
```

equations
  initial
    rel[stat,def] IN init {}
    rel[stat,def] OUT init DEF
  satisfy
    IN = {<S, D> | stat S : STATEMENTS,
                  stat P : predecessor(S),
                  def D : OUT[P]}
    OUT = {<S, D> | stat S : STATEMENTS,
                   def D : DEF[S] union (IN[S] \ KILL[S])}
end equations

```

First, the relations IN and OUT are declared and initialized. Next, two equations are given that very much resemble the ones given above.

**Figure 1.9. Reaching definitions for flow graph in Figure 1.8, “Flow graph for various dataflow problems” [31]**



For our running example (Figure 1.9, “Reaching definitions for flow graph in Figure 1.8, “Flow graph for various dataflow problems”[33]) the results are as follows (see Figure 1.9, “Reaching definitions for flow graph in Figure 1.8, “Flow graph for various dataflow problems”[33]). Relation IN has as value:

```

{ <2, <1, "i">, <3, <2, "j">, <3, <1, "i">, <4, <3, "a">,
  <4, <2, "j">, <4, <1, "i">, <4, <7, "i">, <4, <5, "j">,
  <4, <6, "a">, <5, <4, "i">, <5, <3, "a">, <5, <2, "j">,
  <5, <5, "j">, <5, <6, "a">, <6, <5, "j">, <6, <4, "i">,
  <6, <3, "a">, <6, <6, "a">, <7, <5, "j">, <7, <4, "i">,
  <7, <3, "a">, <7, <6, "a">>
}

```

If we consider statement 3, then the definitions of variables *i* and *j* from the preceding two statements are still valid. A more interesting case are the definitions that can reach statement 4:

- The definitions of variables *a*, *j* and *i* from, respectively, statements 3, 2 and 1.
- The definition of variable *i* from statement 7 (via the backward control flow path from 7 to 4).
- The definition of variable *j* from statement 5 (via the path 5, 7, 4).

- The definition of variable *a* from statement 6 (via the path 6, 7, 4).

Relation OUT has as value:

```
{ <1, <1, "i">>, <2, <2, "j">>, <2, <1, "i">>, <3, <3, "a">>,
  <3, <2, "j">>, <3, <1, "i">>, <4, <4, "i">>, <4, <3, "a">>,
  <4, <2, "j">>, <4, <5, "j">>, <4, <6, "a">>, <5, <5, "j">>,
  <5, <4, "i">>, <5, <3, "a">>, <5, <6, "a">>, <6, <6, "a">>,
  <6, <5, "j">>, <6, <4, "i">>, <7, <7, "i">>, <7, <5, "j">>,
  <7, <3, "a">>, <7, <6, "a">>
}
```

Observe, again for statement 4, that all definitions of variable *i* are missing in OUT[4] since they are killed by the definition of *i* in statement 4 itself. Definitions for *a* and *j* are, however, contained in OUT[4]. The result of reaching definitions computation is illustrated in Figure 1.9, “Reaching definitions for flow graph in Figure 1.8, “Flow graph for various dataflow problems”” [33]. The above definitions are used to formulate the function reaching-definitions. It assumes appropriate definitions for the types *stat* and *var*. It also assumes more general versions of predecessor and successor. We will use it later on in the section called “*Program Slicing*” [35] when defining program slicing. Here is the definition of reaching-definitions:

```
type def = <stat theStat, var theVar>
type use = <stat theStat, var theVar>

set[stat] predecessor(rel[stat,stat] P, stat S) = P[-,S]

set[stat] successor(rel[stat,stat] P, stat S) = P[S,-]

rel[stat, def] reaching-definitions(rel[stat,var] DEFS,
                                   rel[stat,stat] PRED) =
  IN
where
  set[stat] STATEMENT = carrier(PRED)

  rel[stat,def] DEF = {<S,<S,V>> | <stat S, var V> : DEFS}

  rel[stat,def] KILL =
    {<S1, <S2, V>> | <stat S1, var V> : DEFS,
                  <stat S2, V> : DEFS,
                  S1 != S2
    }

  equations
    initial
      rel[stat,def] IN init {}
      rel[stat,def] OUT init DEF
    satisfy
      IN = {<S, D> | int S : STATEMENT,
                  stat P : predecessor(PRED,S),
                  def D : OUT[P]}
      OUT = {<S, D> | int S : STATEMENT,
                  def D : DEF[S] union (IN[S] \ KILL[S])}
    end equations
  end where
```

## Live Variables

The live variables of a statement are those variables whose value will be used by the current statement or some successor of it. The mathematical formulation of this problem is as follows:

## Warning

Fix expression

$$\backslash [ \text{IN}[S] = \text{USE}[S] \backslash \text{cup} (\text{OUT}[S] - \text{DEF}[S]) \backslash ]$$

$$\backslash [ \text{OUT}[S] = \backslash \text{bigcup\_} \{ S' \text{ in successor of } S \} \text{IN}[S'] \backslash ]$$

The first equation says that a variable is live coming into a statement if either it is used before redefinition in that statement or it is live coming out of the statement and is not redefined in it. The second equation says that a variable is live coming out of a statement if and only if it is live coming into one of its successors.

This can be expressed in Rascal as follows:

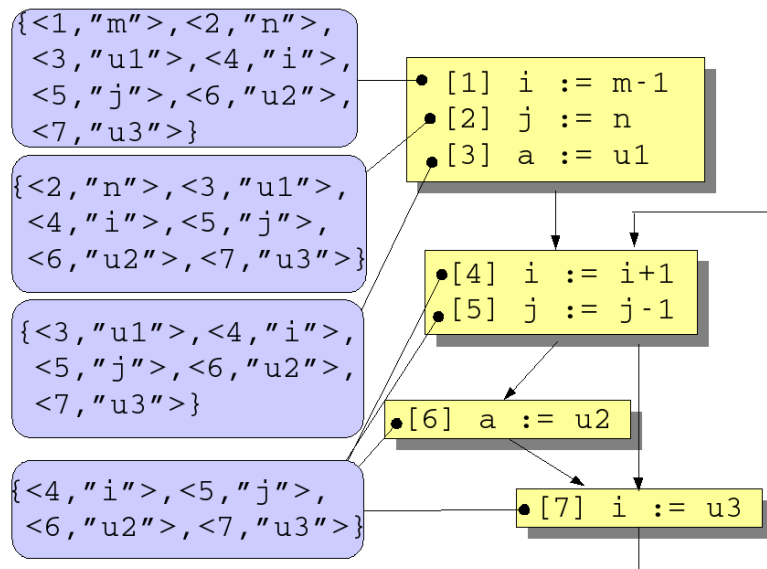
```

equations
  initial
    rel[stat,def] LIN init {}
    rel[stat,def] LOUT init DEF
  satisfy
    LIN = { < S, D> | stat S : STATEMENTS,
                    def D : USE[S] union (LOUT[S] \ (DEF[S]))
            }
    LOUT = { < S, D> | stat S : STATEMENTS,
                      stat Succ : successor(S),
                      def D : LIN[Succ]
            }
end equations

```

The results of live variable analysis for our running example are illustrated in Figure 1.10, “Live variables for flow graph in Figure 1.8, “Flow graph for various dataflow problems”” [35].

**Figure 1.10. Live variables for flow graph in Figure 1.8, “Flow graph for various dataflow problems” [31]**



## Program Slicing

Program slicing is a technique proposed by Weiser [Wei84] for automatically decomposing programs in parts by analyzing their data flow and control flow. Typically, a given statement in a program is selected as the *slicing criterion* and the original program is reduced to an independent subprogram,

called a *slice*, that is guaranteed to represent faithfully the behavior of the original program at the slicing criterion. An example will illustrate this:

[ 1] read(n)	[ 1] read(n)	[ 1] read(n)
[ 2] i := 1	[ 2] i := 1	[ 2] i := 1
[ 3] sum := 0	[ 3] sum := 0	
[ 4] product := 1		[ 4] product := 1
[ 5] while i<= n do	[ 5] while i<= n do	[ 5] while i<= n do
begin	begin	begin
[ 6] sum := sum + i	[ 6] sum := sum + i	
[ 7] product :=		[ 7] product :=
product * i		product * i
[ 8] i := i + 1	[ 8] i := i + 1	[ 8] i := i + 1
end	end	end
[ 9] write(sum)	[ 9] write(sum)	
[10] write(product)		[10] write(product)
<b>(a)</b> Sample program	<b>(b)</b> Slice for statement [9]	<b>(c)</b> Slice for statement [10]

The initial program is given as (a). The slice with statement [9] as slicing criterion is shown in (b): statements [4] and [7] are irrelevant for computing statement [9] and do not occur in the slice. Similarly, (c) shows the slice with statement [10] as slicing criterion. This particular form of slicing is called *backward slicing*. Slicing can be used for debugging and program understanding, optimization and more. An overview of slicing techniques and applications can be found in [Tip95]. Here we will explore a relational formulation of slicing adapted from a proposal in [JR94]J. The basic ingredients of the approach are as follows:

- We assume the relations PRED, DEFS and USES as before.
- We assume an additional set CONTROL-STATEMENT that defines which statements are control statements.
- To tie together dataflow and control flow, three auxiliary variables are introduced:
  - The variable TEST represents the outcome of a specific test of a conditional statement. The conditional statement defines TEST and all statements that are control dependent on this conditional statement will use TEST.
  - The variable EXEC represents the potential execution dependence of a statement on some conditional statement. The dependent statement defines EXEC and an explicit (control) dependence is made between EXEC and the corresponding TEST.
  - The variable CONST represents an arbitrary constant.

The calculation of a (backward) slice now proceeds in six steps:

- Compute the relation  $rel[use, def]$  use-def that relates all uses to their corresponding definitions. The function *reaching-definitions* as shown earlier in the section called “*Reaching Definitions*” [31] does most of the work.
- Compute the relation  $rel[def, use]$  def-use-per-stat that relates the *internal* definitions and uses of a statement.
- Compute the relation  $rel[def, use]$  control-dependence that links all EXECs to the corresponding TESTs.
- Compute the relation  $rel[use, def]$  use-control-def combines use/def dependencies with control dependencies.
- After these preparations, compute the relation  $rel[use, use]$  USE-USE that contains dependencies of uses on uses.

- The backward slice for a given slicing criterion (a use) is now simply the projection of USE-USE for the slicing criterion.

This informal description of backward slicing can now be expressed in Rascal:

```

set[use] BackwardSlice(
  set[stat] CONTROL-STATEMENT,
  rel[stat,stat] PRED,
  rel[stat,var] USES,
  rel[stat,var] DEFS,
  use Criterion)
= USE-USE[Criterion]

where
  rel[stat, def] REACH = reaching-definitions(DEFS, PRED)

  rel[use,def] use-def =
    { <<S1,V>, <S2,V>> | <stat S1, var V> : USES,
      <stat S2, V> : REACH[S1]
    }

  rel[def,use] def-use-per-stat =
    {<<S,V1>, <S,V2>> | <stat S, var V1> : DEFS,
      <S, var V2> : USES
    }
    union
    {<<S,V>, <S,"EXEC">> | <stat S, var V> : DEFS}
    union
    {<<S,"TEST">, <S,V>> | stat S : CONTROL-STATEMENT,
      <S, var V> : domainR(USES, {S})
    }

  rel[stat, stat] CONTROL-DOMINATOR =
    domainR(dominators(PRED), CONTROL-STATEMENT)

  rel[def,use] control-dependence =
    { <<S2, "EXEC">, <S1,"TEST">> |
      <stat S1, stat S2> : CONTROL-DOMINATOR
    }

  rel[use,def] use-control-def = use-def union control-dependence

  rel[use,use] USE-USE = (use-control-def o def-use-per-stat)*

endwhere

```

Let's apply this to the example from the start of this section and assume the following:

```

rel[stat,stat] PRED = { <1,2>, <2,3>, <3,4>, <4,5>, <5,6>, <5,9>,
  <6,7>, <7,8>,<8,5>, <8,9>, <9,10>
  }

rel[stat,var] DEFS = { <1, "n">, <2, "i">, <3, "sum">,
  <4,"product">, <6, "sum">, <7, "product">,
  <8, "i">
  }

rel[stat,var] USES = { <5, "i">, <5, "n">, <6, "sum">, <6,"i">,
  <7, "product">, <7, "i">, <8, "i">,

```

```
        <9, "sum">, <10, "product">
    }
set[int] CONTROL-STATEMENT = { 5 }
```

The result of the slice

```
BackwardSlice(CONTROL-STATEMENT, PRED, USES, DEFS, <9, "sum">)
```

will then be

```
{ <1, "EXEC">, <2, "EXEC">, <3, "EXEC">, <5, "i">, <5, "n">,
  <6, "sum">, <6, "i">, <6, "EXEC">, <8, "i">, <8, "EXEC">,
  <9, "sum"> }
```

Take the domain of this result and we get exactly the statements in (b) of the example.

## Built-in Operators and Functions

The built-in operators and functions can be subdivided in the following categories:

- Benchmark: measuring functions, see the section called “*Benchmark*” [39].
- Boolean: operators and functions on Boolean values, see the section called “*Boolean*” [39].
- Exception: data definition of all soft exceptions that can be caught by Rascal programs, see the section called “*Exception*” [40].
- Graph: graphs are a special kind of binary relation, see the section called “*Graph*” [40].
- Integer: operators and functions on integers, see the section called “*Integer*” [41].
- IO: simple print functions, see the section called “*IO*” [41].
- Labelled Graph: labelled graphs with addition edge information, see the section called “*Labelled Graph*” [42].
- List: operators and functions on lists, see the section called “*List*” [42].
- Location: operators and functions on source locations, see the section called “*Location*” [44].
- Map: operators and functions on maps, see the section called “*Map*” [44].
- Node: operators and functions on nodes, see the section called “*Map*” [44].
- Real: operators and functions on reals, see the section called “*Real*” [46].
- Relation: operators and functions on relations, see the section called “*Relation*” [46].
- Resource: functions to retrieve resources from an Eclipse workspace, see the section called “*Resource (Eclipse only)*” [48].
- RSF: function for reading files in Rigi Standard Format, see the section called “*RSF*” [48].
- Set: operators and functions on sets, see the section called “*Set*” [48].
- String: operators and functions on strings, see the section called “*String*” [51].
- Tuple: operators and functions on tuples, see the section called “*Tuple*” [52].
- UnitTest: functions for unit testing, see the section called “*UnitTest*” [52].
- ValueIO: functions for reading and writing Rascal values, both in textual and in binary form, see the section called “*ValueIO*” [53].

- View: functions for graphical display of values in Eclipse, see the section called “*View (Eclipse only)*” [53].
- Void: the type void, see

All operators are directly available for each program, but library functions have to be imported in each module that uses them.

## Benchmark

**Table 1.3. Benchmark Functions**

Function	Description
<code>real currentTimeMillis()</code>	current time in milliseconds since January 1, 1970 GMT.
p.m. benchmark	measure and report the execution time of name:void-closure pairs.

## Boolean

**Table 1.4. Boolean Operators**

Operator	Description
<code>bool<sub>1</sub> == bool<sub>2</sub></code>	true if both arguments are identical
<code>bool<sub>1</sub> != bool<sub>2</sub></code>	true if both arguments are not identical
<code>bool<sub>1</sub> &lt;= bool<sub>2</sub></code>	true if both arguments are identical or <code>bool<sub>1</sub></code> is false and <code>bool<sub>2</sub></code> is true
<code>bool<sub>1</sub> &lt; bool<sub>2</sub></code>	true if <code>bool<sub>1</sub></code> is false and <code>bool<sub>2</sub></code> is true
<code>bool<sub>1</sub> &gt;= bool<sub>2</sub></code>	true if both arguments are identical or <code>bool<sub>1</sub></code> is true and <code>bool<sub>2</sub></code> is false
<code>bool<sub>1</sub> &gt; bool<sub>2</sub></code>	true if <code>bool<sub>1</sub></code> is true and <code>bool<sub>2</sub></code> is false
<code>bool<sub>1</sub> &amp;&amp; bool<sub>2</sub></code>	yields true if both arguments have the value true and false otherwise
<code>bool<sub>1</sub>    bool<sub>2</sub></code>	yields true if either argument has the value true and false otherwise
<code>bool<sub>1</sub> ==&gt; bool<sub>2</sub></code>	yields false if <code>bool<sub>1</sub></code> has the value true and <code>bool<sub>2</sub></code> has value false, and true otherwise
<code>! bool</code>	yields true if bool is false and true otherwise
<code>bool<sub>1</sub> ? bool<sub>2</sub> : bool<sub>3</sub></code>	if <code>bool<sub>1</sub></code> is true then <code>bool<sub>2</sub></code> else <code>bool<sub>3</sub></code>

**Table 1.5. Boolean Functions**

Function	Description
<code>bool arbBool()</code>	arbitrary boolean value
<code>bool fromInt(int i)</code>	convert an integer to a bool
<code>bool fromString(str s)</code>	convert the strings "true" or "false" to a bool
<code>int toInt(bool b)</code>	convert a boolean value to integer
<code>real toReal(bool b)</code>	convert a boolean value to a real value
<code>str toString(bool b)</code>	convert a boolean value to a string

## Exception

The following "soft" exceptions are defined:

```
data RuntimeException =
  EmptyList
  | EmptyMap
  | EmptySet
  | IndexOutOfBounds(int index)
  | AssertionFailed
  | AssertionFailed(str label)
  | NoSuchElement(value v)
  | IllegalArgument(value v)
  | IllegalArgument
  | IO(str message)
  | FileNotFound(str filename)
  | LocationNotFound(loc location)
  | PermissionDenied
  | PermissionDenied(str message)
  | ModuleNotFound(str name)
  | NoSuchKey(value key)
  | NoSuchAnnotation(str label)
  | Java(str message)
;
```

## Graph

The graph datatype is a special form of binary relation defined as follows:

```
alias graph[&T] = rel[&T from, &T to];
```

**Table 1.6. Graph Functions**

Function	Description
set[&T] bottom(graph[&T] G)	bottom nodes of a graph
set[&T] top(graph[&T] G)	top nodes of a graph
set[&T] reach(graph[&T] G, set[&T] Start)	Reachability from set of start nodes.
set[&T] reachR(graph[&T] G, set[&T] Start, set[&T] Restr)	Reachability from set of start nodes with restriction to certain nodes.
set[&T] reachX(graph[&T] G, set[&T] Start, set[&T] Excl)	Reachability from set of start nodes with exclusion of certain nodes
list[&T] shortestPathPair(graph[&T] G, &T From, &T To)	Shortest path between pair of nodes

The following examples illustrate these functions:

```
rascal> top(<{<1,2>, <1,3>, <2,4>, <3,4>}>);
{1}
rascal> bottom(<{<1,2>, <1,3>, <2,4>, <3,4>}>);
{4}
rascal> reachR({1}, {1, 2, 3}, <{<1,2>, <1,3>, <2,4>, <3,4>}>);
{2, 3}
```



```
rascal> reachX({1}, {2}, {<1,2>, <1,3>, <2,4>, <3,4>});
{3, 4}
```

## Integer

Rascal integers are unbounded in size.

**Table 1.7. Integer Operators**

Operator	Description
$int_1 == int_2$	true if both arguments are numerically equal and false otherwise
$int_1 != int_2$	true if both arguments are numerically unequal and false otherwise
$int_1 \leq int_2$	true if $int_1$ is numerically less than or equal to $int_2$ and false otherwise
$int_1 < int_2$	true if $int_1$ is numerically less than $int_2$ and false otherwise
$int_1 \geq int_2$	true if $int_1$ is numerically greater than or equal than $int_2$ and false otherwise
$int_1 > int_2$	true if $int_1$ is numerically greater than $int_2$ and false otherwise
$int_1 + int_2$	sum of $int_1$ and $int_2$
$int_1 - int_2$	difference of $int_1$ and $int_2$
$int_1 * int_2$	$int_1$ multiplied by $int_2$
$int_1 / int_2$	$int_1$ divided by $int_2$
$int_1 \% int_2$	remainder of dividing $int_1$ by $int_2$
$- int$	negate sign of $int$
$bool ? int_1 : int_2$	if $bool$ is true then $int_1$ else $int_2$

**Table 1.8. Integer Functions**

Function	Description
<code>int abs(int N)</code>	absolute value of integer N
<code>int arbInt()</code>	arbitrary integer value
<code>int arbInt(int limit)</code>	arbitrary integer value in the interval [0, limit)
<code>int max(int n, int m)</code>	largest of two integers
<code>int min(int n, int m)</code>	smallest of two integers
<code>real toReal(int n)</code>	convert an integer value to a real value
<code>str toString(int n)</code>	convert an integer value to a string

## IO

**Table 1.9. IO Functions**

Function	Description
<code>void java println(value V...)</code>	print a list of values on the output stream
<code>list[str] readFile(str filename)</code> <code>throws NoSuchFileError(str msg),</code> <code>IOError(str msg)</code>	read a named file as list of strings

## Labelled Graph

The labelled graph datatype is a special form of binary relation with labelled edges and is defined as follows:

```
alias lgraph[&T,&L] = rel[&T from, &L label, &T to];
```

**Table 1.10. Labelled Graph Functions**

Function	Description
set[&T] bottom(lgraph[&T] G)	bottom nodes of a labelled graph
set[&T] top(lgraph[&T] G)	top nodes of a labelled graph
set[&T] reach(lgraph[&T] G, set[&T] Start)	Reachability from set of start nodes.
set[&T] reachR(lgraph[&T] G, set[&T] Start, set[&T] Restr)	Reachability from set of start nodes with restriction to certain nodes.
set[&T] reachX(lgraph[&T] G, set[&T] Start, set[&T] Excl)	Reachability from set of start nodes with exclusion of certain nodes
list[&T] shortestPathPair(lgraph[&T] G, &T From, &T To)	Shortest path between pair of nodes

### Warning

shortestPath not yet implemented for lgraph.

## List

**Table 1.11. List Operators**

Operator	Description
$list_1 == list_2$	true if both arguments have the same elements in the same order
$list_1 != list_2$	true if both arguments have different elements
$list_1 <= list_2$	true if both lists are equal or $list_1$ is a sublist of $list_2$
$list_1 < list_2$	true if $list_1$ is a sublist of $list_2$
$list_1 >= list_2$	true if both lists are equal or $list_2$ is a sublist of $list_1$
$list_1 > list_2$	true if $list_2$ is a sublist of $list_1$
$list_1 + list_2$	concatenation of $list_1$ and $list_2$
$list_1 - list_2$	list consisting of all elements in $list_1$ that do not occur in $list_2$
$list_1 * list_2$	$list_1$ multiplied by $list_2$
$elm \text{ in } list$	true if $elm$ occurs as element in $list$
$elm \text{ not in } list$	true if $elm$ does not occur as element in $list$
$bool ? list_1 : list_2$	if $bool$ is true then $list_1$ else $list_2$
$list [ int ]$	element at position $int$ in $list$

**Table 1.12. List Functions**

Function	Description
<code>&amp;T average(list[&amp;T] lst, &amp;T zero)</code>	average of elements of a list
<code>list[&amp;T] delete(list[&amp;T] lst, int n)</code>	delete nth element from list
<code>list[int] domain(list[&amp;T] lst)</code>	a list of all legal index values for a list
<code>&amp;T java head(list[&amp;T] lst) throws EmptyListError</code>	get the first element of a list
<code>list[&amp;T] head(list[&amp;T] lst, int n) throws IndexOutOfBoundsException</code>	get the first n elements of a list
<code>&amp;T getOneFrom(list[&amp;T] lst)</code>	get an arbitrary element from a list
<code>list[&amp;T] insertAt(list[&amp;T] lst, int n, &amp;T elm) throws IndexOutOfBoundsException</code>	add an element at a specific position in a list
<code>bool isEmpty(list[&amp;T] lst)</code>	is list empty?
<code>list[&amp;T] mapper(list[&amp;T] lst, &amp;T (&amp;T) fn)</code>	apply a function to each element of a list
<code>&amp;T max(list[&amp;T] lst)</code>	largest element of a list
<code>&amp;T min(list[&amp;T] lst)</code>	smallest element of a list
<code>&amp;T multiply(list[&amp;T] lst, &amp;T unity)</code>	multiply the elements of a list
<code>list[list[&amp;T]] permutations(list[&amp;T] lst)</code>	all permutations of a list
<code>&amp;T reducer(list[&amp;T] lst, &amp;T (&amp;T, &amp;T) fn, &amp;T unit)</code>	apply function F to successive elements of a list
<code>list[&amp;T] reverse(list[&amp;T] lst)</code>	elements of a list in reverse order
<code>int size(list[&amp;T] lst)</code>	number of elements in a list
<code>list[&amp;T] slice(list[&amp;T] lst, int start, int len)</code>	sublist from start of length len
<code>list[&amp;T] sort(list[&amp;T] lst)</code>	sort the elements of a list
<code>&amp;T sum(list[&amp;T] lst, &amp;T zero)</code>	add elements of a List
<code>list[&amp;T] tail(list[&amp;T] lst)</code>	all but the first element of a list
<code>list[&amp;T] tail(list[&amp;T] lst, int len) throws IndexOutOfBoundsException</code>	last n elements of a list
<code>tuple[&amp;T, list[&amp;T]] takeOneFrom(list[&amp;T] lst)</code>	remove an arbitrary element from a list, returns the element and the modified list
<code>map[&amp;A, &amp;B] toMap(list[tuple[&amp;A, &amp;B]] lst)</code>	convert a list of tuples to a map
<code>set[&amp;T] toSet(list[&amp;T] lst)</code>	convert a list to a set
<code>str toString(list[&amp;T] lst)</code>	convert a list to a string

## Location

**Table 1.13. Operations on Locations**

Operator	Description
$loc_1 == loc_2$	true if both arguments are identical and false otherwise
$loc_1 != loc_2$	true if both arguments are not identical and false otherwise
$loc_1 \leq loc_2$	true if $loc_1$ is textually contained in or equal to $loc_2$ and false otherwise
$loc_1 < loc_2$	true if $loc_1$ is strictly textually contained in $loc_2$ and false otherwise
$loc_1 \geq loc_2$	true if $loc_1$ is textually encloses or is equal to $loc_2$ and false otherwise
$loc_1 > loc_2$	true if $loc_1$ is textually encloses $loc_2$ and false otherwise
$loc . field$	retrieve one of the fields of location value

The field names for locations are:

- url
- offset
- length
- beginLine, beginColumn
- endLine, endColumn.

## Map

**Table 1.14. Map Operators**

Operator	Description
$map_1 == map_2$	true if both arguments consist of the same pairs
$map_1 != map_2$	true if both arguments have different pairs
$map_1 \leq map_2$	true if all pairs in $map_1$ occur in $map_2$ or $map_1$ and $map_2$ are equal
$map_1 < map_2$	true if all pairs in $map_1$ occur in $map_2$ but $map_1$ and $map_2$ are not equal
$map_1 \geq map_2$	true if all pairs in $map_2$ occur in $map_1$ or $map_1$ and $map_2$ are equal
$map_1 > map_2$	true if all pairs in $map_2$ occur in $map_1$ but $map_1$ and $map_2$ are not equal
$map_1 + map_2$	union of $map_1$ and $map_2$
$map_1 - map_2$	difference of $map_1$ and $map_2$
$key \text{ in } map$	true if $key$ occurs in a key:value pair in map
$key_1 \text{ not in } map_2$	true if $key$ does not occur in a key:value pair in map
$bool ? map_1 : map_2$	if $bool$ is true then $map_1$ else $map_2$
$map [ key ]$	the value associated with $key$ in $map$ if that exists, undefined otherwise

**Table 1.15. Map Functions**

Function	Description
<code>set[&amp;K] domain(map[&amp;K, &amp;V] M)</code>	the domain (keys) of a map
<code>&amp;K getOneFrom(map[&amp;K, &amp;V] M)</code>	arbitrary key of a map
<code>map[&amp;V, &amp;K] invert(map[&amp;K, &amp;V] M)</code>	map with key and value inverted
<code>bool isEmpty(map[&amp;K, &amp;V] M)</code>	is map empty?
<code>map[&amp;K, &amp;V] mapper(map[&amp;K, &amp;V] M, &amp;K (&amp;K) F, &amp;V (&amp;V) G)</code>	apply two functions to each key/value pair in a map.
<code>set[&amp;V] range(map[&amp;K, &amp;V] M)</code>	the range (values) of a map
<code>int size(map[&amp;K, &amp;V] M)</code>	number of elements in a map.
<code>list[tuple[&amp;K, &amp;V]] toList(map[&amp;K, &amp;V] M)</code>	convert a map to a list
<code>rel[&amp;K, &amp;V] toRel(map[&amp;K, &amp;V] M)</code>	convert a map to a relation
<code>str toString(map[&amp;K, &amp;V] M)</code>	convert a map to a string.

## Node

**Table 1.16. Node Operators**

Operator	Description
<code>node<sub>1</sub> == node<sub>2</sub></code>	true if both arguments are identical
<code>node<sub>1</sub> != node<sub>2</sub></code>	true if both arguments are not identical
<code>node<sub>1</sub> &lt;= node<sub>2</sub></code>	
<code>node<sub>1</sub> &lt; node<sub>2</sub></code>	
<code>node<sub>1</sub> &gt;= node<sub>2</sub></code>	
<code>node<sub>1</sub> &gt; node<sub>2</sub></code>	
<code>bool ? node<sub>1</sub> : node<sub>2</sub></code>	if <i>bool</i> is true then <i>node<sub>1</sub></i> else <i>node<sub>2</sub></i>
<code>node [ int ]</code>	child of <i>node</i> at position <i>int</i>

**Table 1.17. Node Functions**

Function	Description
<code>int arity(node T)</code>	number of children of a node
<code>list[value] getChildren(node T)</code>	get the children of a node
<code>str getName(node T)</code>	get the function name of a node
<code>node makeNode(str N, value V...)</code>	create a node given its function name and arguments

## Real

**Table 1.18. Real Operators**

Operator	Description
$real_1 == real_2$	true if both arguments are numerically equal and false otherwise
$real_1 != real_2$	true if both arguments are numerically unequal and false otherwise
$real_1 \leq real_2$	true if $real_1$ is numerically less than or equal to $real_2$ and false otherwise
$real_1 < real_2$	true if $real_1$ is numerically less than $real_2$ and false otherwise
$real_1 \geq real_2$	true if $real_1$ is numerically greater than or equal than $real_2$ and false otherwise
$real_1 > real_2$	true if $real_1$ is numerically greater than $real_2$ and false otherwise
$real_1 + real_2$	sum of $real_1$ and $real_2$
$real_1 - real_2$	difference of $real_1$ and $real_2$
$real_1 * real_2$	$real_1$ multiplied by $real_2$
$real_1 / real_2$	$real_1$ divided by $real_2$
$-real$	negate sign of $real$
$real_1 \% real_2$	remainder of dividing $real_1$ by $real_2$

**Table 1.19. Real Functions**

Function	Description
<code>real arbReal()</code>	an arbitrary real value in the interval [0.0,1.0).
<code>real max(real n, real m)</code>	largest of two reals
<code>int toInteger(real d)</code>	convert a real to integer.
<code>str toString(real d)</code>	convert a real to a string.

## Relation

Relation are sets of tuples, therefore all set operators (see, Table 1.24, “Set Operators”[49]) apply to relations as well

**Table 1.20. Operations on Relations**

Operator	Description
$rel_1 \circ rel_2$	yields the relation resulting from the composition of the two arguments
$set_1 \times set_2$	yields the relation resulting from the Cartesian product of the two arguments
$rel +$	yields the relation resulting from the transitive closure of $rel$
$rel *$	yields the relation resulting from the reflexive transitive closure of $rel$
$rel [ elem ]$	yields the right image of $rel$
$rel [ set ]$	yields the right image of $rel$
$rel < index_1, index_2, \dots >$	

Examples:

```
{<1,10>, <2,20>, <3,15>} o {<10,100>, <20,200>} yields {<1,100>, <2,200>}.
```

```
{1, 2, 3} x {9} yields {<1, 9>, <2, 9>, <3, 9>}.
```

Rel has value {<1,10>, <2,20>, <1,11>, <3,30>, <2,21>} in the following example.

```
Rel[1] yields {10, 11}.
```

```
Rel[{1}] yields {10, 11}.
```

```
Rel[{1, 2}] yields {10, 11, 20, 21}.
```

**Table 1.21. Relation Functions**

Function	Description
set[&T] carrier (rel[&T,&T] R)	all elements in any tuple in a relation
rel[&T,&T] carrierR (rel[&T,&T] R, set[&T] S)	relation restricted to tuples with elements in a set S
rel[&T,&T] carrierX (rel[&T,&T] R, set[&T] S)	relation excluded tuples with some element in S
rel[&T0, &T1] complement (rel[&T0, &T1] R)	complement of relation
set[&T0] domain (rel[&T0,&T1] R)	first element of each tuple in binary relation
rel[&T0,&T1] domainR (rel[&T0,&T1] R, set[&T0] S)	restriction of a relation to tuples with first element in S
rel[&T0,&T1] domainX (rel[&T0,&T1] R, set[&T0] S)	relation excluded tuples with first element in S
ident?	
rel[&T1, &T0] invert (rel[&T0, &T1] R)	inverse the tuples in a relation
rel[&T2, &T1, &T0] invert (rel[&T0, &T1, &T2] R)	all but the first element of each tuples in binary relation
rel[&T0,&T1] rangeR (rel[&T0,&T1] R, set[&T2] S)	restriction of a binary relation to tuples with second element in set S

Examples:

```
id({1,2,3}) yields {<1,1>, <2,2>, <3,3>}.
```

```
id({"mon", "tue", "wed"}) yields {<"mon","mon">, <"tue","tue">, <"wed","wed">}.
```

```
inv({<1,10>, <2,20>}) yields {<10,1>,<20,2>}.
```

```
compl({<1,10>}) yields {<1, 1>, <10, 1>, <10, 10>}.
```

```
domain({<1,10>, <2,20>}) yields {1, 2}.
```

```
domain({<"mon", 1>, <"tue", 2>}) yields {"mon", "tue"}.
```

```

range(<{<1,10>, <2,20>}>) yields {10, 20}.

range(<{"mon", 1>, <"tue", 2>}>) yields {1, 2}.

carrier(<{<1,10>, <2,20>}>) yields {1, 10, 2, 20}

domainR(<{<1,10>, <2,20>, <3,30>}>, {3, 1});
{<1,10>, <3,30>}.

rangeR(<{<1,10>, <2,20>, <3,30>}>, {30, 10}) yields {<1,10>, <3,30>}.

carrierR(<{<1,10>, <2,20>, <3,30>}>, {10, 1, 20}) yields {<1,10>}.

domainX(<{<1,10>, <2,20>, <3,30>}>, {3, 1}) yields {<2, 20>}.
rangeX(<{<1,10>, <2,20>, <3,30>}>, {30, 10}) yields {<2, 20>}.

carrierX(<{<1,10>, <2,20>, <3,30>}>, {10, 1, 20}) yields {<3,30>}.

```

## RSF

**Table 1.22. RSF Functions**

Function	Description
map[str, rel[str,str]] readRSF(str nameRSFFile)	read a file in Rigi Standard Format (RSF).

## Resource (Eclipse only)

```

data Resource = root(set[Resource] projects)
                | project(str name, set[Resource] contents)
                | folder(str name, set[Resource] contents)
                | file(str name, str extension);

```

**Table 1.23. Resource Functions**

Function	Description
Resource java root()	The root of the Eclipse workspace
set[str] java projects()	The projects in the Eclipse workspace
set[str] java references(str project)	The project references of a given project
loc java location(str project)	Source location of given project
set[loc] java files(str project)	The files contained in a project

## Set



**Table 1.24. Set Operators**

Operator	Description
$set_1 == set_2$	true if both arguments are equal sets and false otherwise
$set_1 != set_2$	true if both arguments are unequal sets and false otherwise
$set_1 \leq set_2$	true if $set_1$ is a subset of $set_2$ and false otherwise
$set_1 < set_2$	true if $set_1$ is a strict subset of $set_2$ and false otherwise
$set_1 \geq set_2$	true if $set_1$ is a superset of $set_2$ and false otherwise
$set_1 > set_2$	true if $set_1$ is a strict superset of $set_2$ and false otherwise
$set_1 + set_2$	set resulting from the union of the two arguments
$set_1 - set_2$	the set resulting from the difference of the two arguments
$set_1 * set_2$	set resulting from the product of the two arguments
$set_1 \& set_2$	set resulting from the intersection of the two arguments
$elm \text{ in } set$	true if $elm$ occurs as element in $set$ and false otherwise
$elm \text{ not in } set$	false if $elm$ occurs as element in $set$ and false otherwise
$set \text{ join } set$	
$bool ? set_1 : set_2$	

Examples:

```

rascal> {1, 2, 3} + {4, 5, 6};
{1, 2, 3, 4, 5, 6}

rascal> {1, 2, 3} + {1, 2, 3};
{1, 2, 3}

rascal> {1, 2, 3, 4} - {1, 2, 3};
{4}

rascal> {1, 2, 3} - {4, 5, 6};
{1, 2, 3}

rascal> {1, 2, 3} & {4, 5, 6};
{ }

rascal> {1, 2, 3} & {1, 2, 3};
{1, 2, 3}

rascal> 3 in {1, 2, 3};
true

rascal> 4 in {1, 2, 3};
false

rascal> 3 notin {1, 2, 3};
false

rascal> 4 notin {1, 2, 3};
true

rascal> <2,20> in {<1,10>, <2,20>, <3,30>};
true

```

```
rascal> <4,40> notin {<1,10>, <2,20>, <3,30>};
true
```

**Table 1.25. Set Functions**

Function	Description
&T average(set[&T] st, &T zero)	average of the elements of a set
&T getOneFrom(set[&T] st)	pick a random element from a set
bool isEmpty(set[&T] st)	Is set empty?
set[&T] mapper(set[&T] st, &T (&T,&T) fn)	apply a function to each element of a set
&T max(set[&T] st)	largest element of a set
&T min(set[&T] st)	smallest element of a set
&T multiply(set[&T] st, &T unity)	multiply the elements of a set
set[set[&T]] power(set[&T] st)	all subsets of a set
set[set[&T]] power1(set[&T] st)	all subsets (excluding empty set) of a set
&T reducer(set[&T] st, &T (&T,&T) fn, &T unit)	apply function F to successive elements of a set
int size(set[&T] st)	number of elements in a set
&T sum(set[&T] st, &T zero)	add the elements of a set
tuple[&T, set[&T]] takeOneFrom(set[&T] st)	remove an arbitrary element from a set, returns the element and the modified set
list[&T] toList(set[&T] st)	convert a set to a list
map[&A,&B] toMap(rel[&A, &B] st)	convert a set of tuples to a map
str toString(set[&T] st)	convert a set to a string

Examples:

```
rascal> power({1, 2, 3, 4});
{ {}, {1}, {2}, {3}, {4}, {1,2}, {1,3}, {1,4}, {2,3}, {2,4},
  {3,4}, {1,2,3}, {1,2,4}, {1,3,4}, {2,3,4}, {1,2,3,4}
}

rascal> power1({1, 2, 3, 4});
{ {1}, {2}, {3}, {4}, {1,2}, {1,3}, {1,4}, {2,3}, {2,4},
  {3,4}, {1,2,3}, {1,2,4}, {1,3,4}, {2,3,4}, {1,2,3,4}
}

rascal> size({1,2,3});
3

rascal> size({<1,10>, <2,20>, <3,30>});
3
```

# String

**Table 1.26. Operations on Strings**

Operator	Description
<code>str<sub>1</sub> == str<sub>2</sub></code>	yields true if both arguments are equal and false otherwise
<code>str<sub>1</sub> != str<sub>2</sub></code>	yields true if both arguments are unequal and false otherwise
<code>str<sub>1</sub> &lt;= str<sub>2</sub></code>	yields true if <code>str<sub>1</sub></code> is lexicographically less than or equal to <code>str<sub>2</sub></code> and false otherwise
<code>str<sub>1</sub> &lt; str<sub>2</sub></code>	yields true if <code>str<sub>1</sub></code> is lexicographically less than <code>str<sub>2</sub></code> and false otherwise
<code>str<sub>1</sub> &gt;= str<sub>2</sub></code>	yields true if <code>str<sub>1</sub></code> is lexicographically greater than or equal to <code>str<sub>2</sub></code> and false otherwise
<code>str<sub>1</sub> &gt; str<sub>2</sub></code>	yields true if <code>str<sub>1</sub></code> is lexicographically greater than <code>str<sub>2</sub></code> and false otherwise
<code>str<sub>1</sub> + str<sub>2</sub></code>	concatenates <code>str<sub>1</sub></code> and <code>str<sub>2</sub></code>

**Table 1.27. String Functions**

Function	Description
<code>int charAt(str s, int i)</code> throws <code>out_of_range(str msg)</code>	character at position <code>i</code> in string <code>s</code> .
<code>bool endsWith(str s, str suffix)</code>	true if string <code>s</code> ends with given string suffix.
<code>str center(str s, int n)</code>	center <code>s</code> in string of length <code>n</code> using spaces
<code>str center(str s, int n, str pad)</code>	center <code>s</code> in string of length <code>n</code> using a pad character
<code>bool isEmpty(str s)</code>	is string empty?
<code>str left(str s, int n)</code>	left align <code>s</code> in string of length <code>n</code> using spaces
<code>str left(str s, int n, str pad)</code>	left align <code>s</code> in string of length <code>n</code> using pad character
<code>str right(str s, int n)</code>	right align <code>s</code> in string of length <code>n</code> using spaces
<code>str reverse(str s)</code>	string with all characters in reverse order.
<code>int size(str s)</code>	the length of string <code>s</code> .
<code>bool startsWith(str s, str prefix)</code>	true if string <code>s</code> starts with the string prefix.
<code>str toLowerCase(str s)</code>	convert all characters in string <code>s</code> to lowercase.
<code>str toUpperCase(str s)</code>	convert all characters in string <code>s</code> to uppercase.

## Tuple

**Table 1.28. Tuple Operators**

Operator	Description
$tuple_1 == tuple_2$	true if both arguments are identical
$tuple_1 != tuple_2$	true if both arguments are not identical
$tuple_1 <= tuple_2$	true if both arguments are identical or if the leftmost element in $tuple_1$ that differs from the corresponding in $tuple_2$ is smaller than that element in $tuple_2$
$tuple_1 < tuple_2$	true if both arguments are not identical and the leftmost element in $tuple_1$ that differs from the corresponding in $tuple_2$ is smaller than that element in $tuple_2$
$tuple_1 >= tuple_2$	true if both arguments are identical or if the leftmost element in $tuple_1$ that differs from the corresponding in $tuple_2$ is greater than that element in $tuple_2$
$tuple_1 > tuple_2$	true if both arguments are not identical and the leftmost element in $tuple_1$ that differs from the corresponding in $tuple_2$ is greater than that element in $tuple_2$
$tuple_1 + tuple_2$	concatenates $tuple_1$ and $tuple_2$
$bool ? tuple_1 : tuple_2$	if $bool$ is true then $tuple_1$ else $tuple_2$
$tuple . name$	select field <i>name</i> from <i>tuple</i>
$tuple [ int ]$	select field at position <i>int</i> from <i>tuple</i>

## UnitTest

We provided a very rudimentary library for unit testing that will certainly evolve over time:

**Table 1.29. UnitTest Functions**

Function	Description
<code>void assertTrue(bool outcome)</code>	check that outcome is true
<code>void assertEquals(value V1, value V2)</code>	check that two values are equal
<code>bool report()</code>	print unit test summary
<code>bool report(str msg)</code>	print unit test summary, including msg

## Value

**Table 1.30. Value Operators**

Operator	Description
<code>value<sub>1</sub> == value<sub>2</sub></code>	true if both arguments are identical
<code>value<sub>1</sub> != value<sub>2</sub></code>	true if both arguments are not identical
<code>value<sub>1</sub> &lt;= value<sub>2</sub></code>	
<code>value<sub>1</sub> &lt; value<sub>2</sub></code>	
<code>value<sub>1</sub> &gt;= value<sub>2</sub></code>	
<code>value<sub>1</sub> &gt; value<sub>2</sub></code>	
<code>bool ? value<sub>1</sub> : value<sub>2</sub></code>	if <i>bool</i> is true then <i>value<sub>1</sub></i> else <i>value<sub>2</sub></i>

## ValueIO

**Table 1.31. ValueIO Functions**

Function	Description
<code>value readValueFromBinaryFile(str namePBFFile)</code>	read a value from a binary file in PBF format
<code>value readValueFromTextFile(str namePBFFile)</code>	read a value from a text file
<code>void writeValueToBinaryFile(str namePBFFile, value val)</code>	write a value to a binary file in PBF format
<code>void writeValueToTextFile(str namePBFFile, value val)</code>	write a value to a binary file in PBF format

## View (Eclipse only)

**Table 1.32. View Functions**

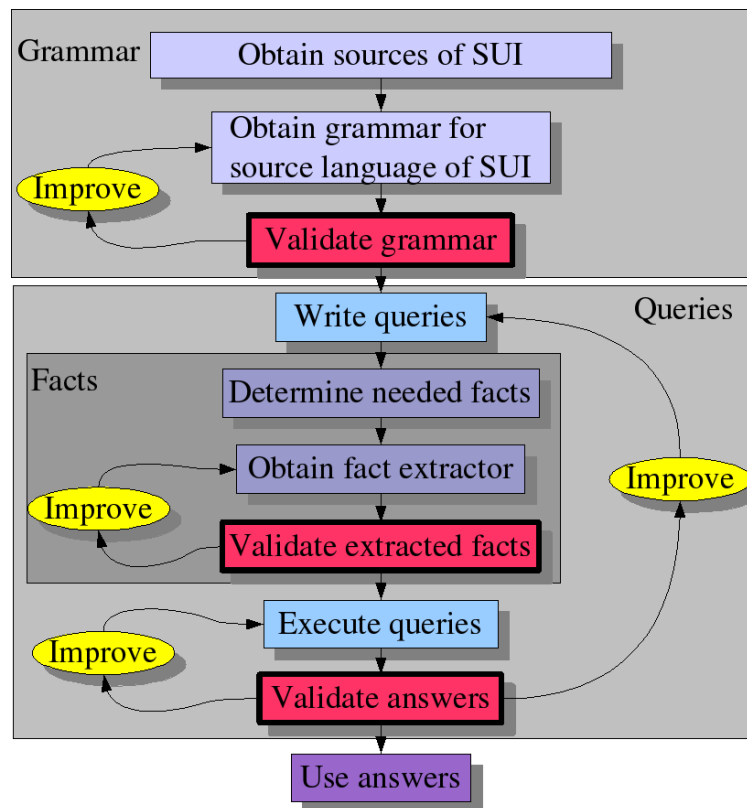
Function	Description
<code>void show(value v)</code>	Show value <i>v</i> in a graphical viewer
<code>void browse(value v)</code>	Show value <i>v</i> a graphical browser

## Void

There are no operators or functions defined on the type `void`.

## Extracting Facts from Source Code (TODO)

In this tutorial we have, so far, concentrated on querying and enriching facts that have been extracted from source code. As we have seen from the examples, once these facts are available, a concise Rascal program suffices to do the required processing. But how is fact extraction achieved and how difficult is it? To answer these questions we first describe the workflow of the fact extraction process (the section called “*Workflow for Fact Extraction*” [54]) and give strategies for fact extraction (the section called “*Strategies for Fact Extraction*” [55]).

**Figure 1.11. Workflow for fact extraction**

## Workflow for Fact Extraction

Figure 1.11, “Workflow for fact extraction”[54] shows a typical workflow for fact extraction for a *System Under Investigation* (SUI). It assumes that the SUI uses only *one* programming language and that you need only one grammar. In realistic cases, however, several such grammars may be needed. The workflow consists of three main phases:

- Grammar: Obtain and improve the grammar for the source language of the SUI.
- Facts: Obtain and improve facts extracted from the SUI.
- Queries: Write and improve queries that give the desired answers.

Of course, it may happen that you have a lucky day and that extracted facts are readily available or that you can reuse a good quality fact extractor that you can apply to the SUI. On ordinary days you have the above workflow as fall-back. It may come as a surprise that there is such a strong emphasis on validation in this workflow. The reason is that the SUI is usually a huge system that defeats manual inspection. Therefore we must be very careful that we validate the outcome of each phase.

**Grammar.** In many cases there is no canned grammar available that can be used to parse the programming language dialect used in the SUI. Usually an existing grammar can be adjusted to that dialect, but then it is then mandatory to validate that the adjusted grammar can be used to parse the sources of the SUI.

**Facts.** It may happen that the facts extracted from the source code are *wrong*. Typical error classes are:

- Extracted facts are *wrong*: the extracted facts incorrectly state that procedure P calls procedure Q but this is contradicted by a source code inspection.

- Extracted facts are *incomplete*: the inheritance between certain classes in Java code is missing.

The strategy to validate extracted facts differ per case but here are three strategies:

- Postprocess the extracted facts (using Rascal, of course) to obtain trivial facts about the source code such as total lines of source code and number of procedures, classes, interfaces and the like. Next validate these trivial facts with tools like `wc` (word and line count), `grep` (regular expression matching) and others.
- Do a manual fact extraction on a small subset of the code and compare this with the automatically extracted facts.
- Use another tool on the same source and compare results whenever possible. A typical example is a comparison of a call relation extracted with different tools.

**Queries.** For the validation of the answers to the queries essentially the same approach can be used as for validating the facts. Manual checking of answers on random samples of the SUI may be mandatory. It also happens frequently that answers inspire new queries that lead to new answers, and so on.

## Strategies for Fact Extraction

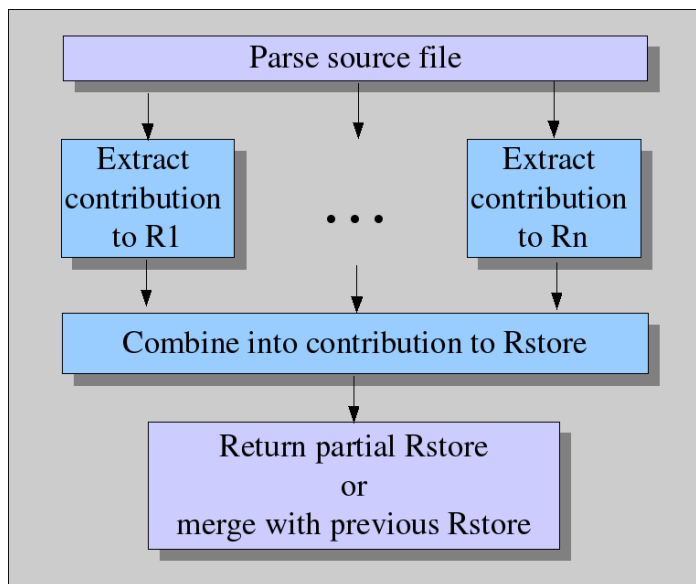
The following global scenario's are available when writing a fact extractor:

- *Dump-and-Merge*: Parse each source file, extract the relevant facts, and return the resulting (partial) Rstore. In a separate phase, merge all the partial Rstores into a complete Rstore for the whole SUI. The tool `{\tt merge-rstores}` is available for this.
- *Extract-and-Update*: Parse each source file, extract the relevant facts, and add these directly to the partial Rstore that has been constructed for previous source files.

The experience is that the *Extract-and-Update* is more efficient. A second consideration is the scenario used for the fact extraction per file. Here there are again two possibilities:

- *All-in-One*: Write one function that extracts all facts in one traversal of the source file. Typically, this function has an Rstore as argument and returns an Rstore as well. During the visit of specific language constructs additions are made to named sets or relations in the Rstore.
- *Separation-of-Concerns*: Write a separate function for each fact you want to extract. Typically, each function takes a set or relation as argument and returns an updated version of it. At the top level all these functions are called and their results are put into an Rstore. This strategy is illustrated in Figure 1.12, “Separation-of-Concerns strategy for fact extraction” [56].

The experience here is that everybody starts with the *All-in-One* strategy but that the complexities of the interactions between the various fact extraction concerns soon start to hurt. The advice is therefore to use the *Separation-of-Concerns* strategy even if it may be seem to be less efficient since it requires a traversal of the source program for each extracted set or relation.

**Figure 1.12. Separation-of-Concerns strategy for fact extraction**

## Fact Extraction using ASF+SDF

Although facts can be extracted in many ways, ASF+SDF is the tool of preference to do this. Examples are given In XXX.

### Warning

Fix reference

## Concluding remarks

It is not unusual that the effort that is needed to write a fact extractor is much larger than the few lines of Rascal that are sufficient for the further processing of these facts. What can we learn from this observation? First, that even in simple cases fact extraction is more complicated than the processing of these facts. This may be due to the following:

- The facts we are interested in may be scattered over different language constructs. This implies that the fact extractor has to cover all these cases.
- The extracted facts are completely optimized for relational processing but places a burden on the fact extractor to perform this optimization.

Second, that several research questions remain unanswered:

- Is it possible to solve (parts of) the fact extraction in a language-parametric way. In other words, is it possible to define generic extraction methods that apply to multiple languages?
- Is a further integration of fact extraction with relational processing desirable? Is it, for instance, useful to bring some of the syntactic program domains like expressions and statements to the relational domain?

## Table of Built-in Operators



**Table 1.33. All Operators**

Operator	Description	See
$exp_1 [ name = exp_2 ]$	Field assignment	
$exp . name$	Field selection	
$exp < field, \dots >$	Field projection	
$exp_1 [ exp_2 ]$	Index	
$exp ?$	Isdefined: true is $exp$ has a well-defined value	
$! exp$	Negation	
$- exp$	Negation	
$exp +$	Transitive closure	
$exp *$	Reflexive transitive	
$exp @ name$	Attribute value	
$exp_1 [ @ name = exp_2 ]$	Assign attribute value	
$exp_1 \circ exp_2$	Composition	
$exp_1 / exp_2$	Division	
$exp_1 \% exp_2$	Modulo	
$exp_1 * exp_2$	Multiplication/product	
$exp_1 \& exp_2$	Intersection	
$exp_1 + exp_2$	Addition/concatenation/ union	
$exp_1 - exp_2$	Subtraction/difference	
$exp_1 \text{ join } exp_2$	Join	
$exp_1 \text{ in } exp_2$	Element of	
$exp_1 \text{ not in } exp_2$	Not element of	
$exp_1 \leq exp_2$	Less than/sublist /subset	
$exp_1 < exp_2$	Less than/strict sublist/ strict subset	
$exp_1 \geq exp_2$	Greater than/superlist/ superset	
$exp_1 > exp_2$	Greater than/strict superlist/ strict superset	
$pat := exp$	Match	
$pat !:= exp$	No Match	
$exp_1 == exp_2$	Equality	
$exp_1 != exp_2$	Inequality	
$exp_1 ? exp_2$	Ifdefined Otherwise	
$exp_1 ? exp_2 : exp_3$	Conditional expression	
$exp_1 ==> exp_2$	Implication	
$exp_1 <==> exp_2$	Equivalence	
$exp_1 \&\& exp_2$	Boolean and	
$exp_1    exp_2$	Boolean or	



# Table of Built-in Functions

		Real?, Set	
	multiply	List, Set	
	permutations	String with Rascal	
	power	Set	
	power1	Set	
	println	IO	
	projects	Resource	
	range	Map	
	reach	Graph, LabelledGraph	
	reachR	Graph, LabelledGraph	
	reachX	Graph, LabelledGraph	
	readFile	IO	
	readRSF	RSF	
	readValueFromBinaryFile	ValueIO	
	readValueFromTextFile	ValueIO	
	reducer	List, Set	
	references	Resource	
	report	UnitTest	
	reverse	List, String	
	right	String	
	root	Resource	
	shortestPathPair	Graph, LabelledGraph	
	show	View	
	size	List?, Map, Set, String	
	slice	List	
	sort	List	
	startsWith	String	
	sum	List, Set	
	tail	List	
	takeOneFrom	List, Set	
	toInt	Boolean	
	toInteger!	Real	
	toList	Map, Set	
	toLowerCase	String	
	toMap	List, Set	
	top	Graph, LabelledGraph	
	toReal	Boolean, Integer	
	toRel	Map	
	toSet	List	
	toString	Boolean, Integer, List, Map, Real, Set	
	toUpperCase	String	
	writeValueToBinaryFile	ValueIO	
	writeValueToTextFile	ValueIO	

# Bibliography

- [ASU86] A.V. Aho, R. Sethi, and J.D. Ullman. *Compilers: Principles, Techniques and Tools*. Addison-Wesley. 1986.
- [BNL03] D. Beyer, A Noack, and C. Lewerentz. *Simple and efficient relational querying of software structures*. Proceedings of the 10th IEEE Working Conference on Reverse Engineering (WCRE 2003). . 2003. To appear.
- [KN96] E. Koutsofios and S.C. North. *Drawing graphs with dot*. Technical report. AT&T Bell Laboratories. Murray Hill, NJ. 1996. See also [www.graphviz.org](http://www.graphviz.org).
- [FKO98] L.M.G. Feijs, R. Krikhaar, and R.C. Ommering. *A relational approach to support software architecture analysis*. 371--400. *Software Practice and Experience*. 28. 4. april 1998.
- [Hol96] R.C. Holt. *Binary relational algebra applied to software architecture*. CSRI345. University of Toronto. march 1996.
- [JR94] D.J. Jackson and E.J. Rollins. *A new model of program dependences for reverse engineering*. 2--10. Proceedings of the 2nd ACM SIGSOFT symposium on Foundations of software engineering. . ACM SIGSOFT Software Engineering Notes.  
<seriesvolnum>19</seriesvolnum>  
1994.
- [Kli03] P. Klint. *How understanding and restructuring differ from compiling---a rewriting perspective*. 2--12. Proceedings of the 11th International Workshop on Program Comprehension (IWPC03). . 2003. IEEE Computer Society.
- [Kri99] R.L. Krikhaar. *Software Architecture Reconstruction*. PhD thesis. University of Amsterdam. 1999.
- [McC76] T.J. McCabe. *A complexity measure*. 308--320. *IEEE Transactions on Software Engineering*. SE-12. 3. 1976.
- [MK88] H. Müller and K. Klashinsky. *Rigi -- a system for programming-in-the-large*. 80--86,. Proceedings of the 10th International Conference on Software Engineering (ICSE 10),. . April 1988.
- [Tip95] F. Tip. *A survey of program slicing techniques*. 121--189. *Journal of Programming Languages*. 3. 3. 1995.
- [Wei84] M. Weiser. *Program slicing*. 352--357. *IEEE Transactions on Software Engineering*. SE-10. 4. July 1984.