

Five Languages Symposium: Rascal

Introduction

Rascal is programming language for source code analysis and transformation. The primary application areas are (legacy) system renovation, reverse engineering and reengineering, and the implementation of domain specific languages (DSLs). DSLs are languages tailored to a specific application domain. Examples include SQL, Excel, Make, LaTeX, VHDL, etc. Today we are going to use Rascal to implement a small language for statemachines. The example is derived from Martin Fowler's book on Domain Specific Languages. The relevant chapter is published online:

<http://www.informit.com/articles/article.aspx?p=1592379>

State machines are useful for describing state dependent behaviour, for instance to control machines or workflow engines. DSLs are a particular form of model-driven development (MDE), where software is specified using high-level models, from which then the implementation is generated.

For more information on the theory of this kind of state machines, you may want to consult Wikipedia.

In the tutorial we will explore the Rascal language and environment by implementing the following facets of the DSL:

- A context-free grammar to describe the syntax of state machines
- An algebraic data type (ADT) for describing state machine abstract syntax trees (ASTs)
- Reset events (see the link above) are syntactic sugar: they can be *desugared* into an equivalent statemachine that does not use them.
- Extraction of relations from a state machine. This allows easier analysis of state machines. The relations can be connected to the state machine visualizer (provided by us).
- A consistency checker for state machines. This component, for instance, highlights use of undefined states or events, marks duplicate states, commands or events and detects unreachable states.
- A code generator that produces a Java code consuming and producing tokens. (This is called a Model-to-Text transformation)
- A transformation that takes two state machines and produces a new state machine that runs the two original machines in parallel. (This is called a Model-to-Model transformation). The resulting state machine can be input to the original code generator and visualizer.

- A simple evaluator to simulate a state machine. This can be connected to the visualization to interactively step through a state machine.
- A simple source-to-source transformation to implement a rename refactoring for states, events and/or commands.
- Provide domain-specific IDE features for state machines: context-menus to invoke the code generator, outline views, folding, error marking.

Some of these assignments are more complicated than others. This is not a problem: we will see how far we get!

To get up and running, download the following zipfile:

5lang-rascal.zip

Unzip it into a dedicated directory. The Zip file contains pre-built Eclipse workspace including a Rascal project for state machines. The project contains some setup code and example state machines.

Warming Up

Syntax Definition for State Machines

Rascal has builtin support for context-free grammars which can be used to define the syntax of programming languages. A syntax rule consists of the following parts:

```
syntax NonTerminal = label_1: Element_1i ... Element_1m1
                    | ...
                    | label_n: Element ... Element_nmn
```

This defines a rule named “NonTerminal” with n alternatives. Each alternative has a label and a sequence of Elements. The elements of an alternative define the syntax to be recognized by this rule. An element can be one of the following symbols:

- “a literal”
- ANonTerminal
- a regular symbol: $X?$ for optional, X^* for zero-or-more, X^+ for one-or-more, and the separated list operators: $\{X \text{ “sep”}\}^*$, and $\{X \text{ “sep”}\}^+$. X can be any symbol, but typically will be a non-terminal

To define lexical rules (e.g., for identifiers) character classes are used (similar to regular expressions):

- char class `[a-z]`: recognize a character between a and z. Or: `[\t\n\r\]`: recognize a whitespace character
- `![a-z]`: recognize any non-lowercase-alphabetic character
- `?`, `*`, and `+` can be used on character classes as well

To create a grammar for state machines, create a new Rascal module, and add rules to recognize the syntax invented by Martin Fowler (page 3 of the article cited above).

To get up and running add the following definitions for layout (whitespace and comments) and identifiers:

```
syntax Id = lex [a-zA-Z][a-zA-Z0-9_]* - Reserved # [a-zA-Z0-9_];
syntax Reserved = "events" | "end" | "resetEvents" | "state" | "actions" ;
syntax LAYOUT = lex whitespace: [\t\n\r\ ] | lex Comment ;
layout LAYOUTLIST = LAYOUT* # [\t\n\r\ ] # "/*" ;
syntax Comment = lex @category="Comment" "/*" CommentChar* "*/" ;
syntax CommentChar = lex ![*] | lex Asterisk ;
syntax Asterisk = lex [*] # [/] ;
```

Quiz: what is the meaning of the `#` and `-` constructs?

Don't forget to add a start syntax rule. This will instruct the parser what to recognize when a file is parsed in an editor. For instance, like this:

```
start syntax Controller = controller: Events ResetEvents? Commands? State+ states;
```

After you've created the syntax module. Create a module `Plugin.rsc` in the `src` directory of the project. Add the following lines to make the syntax available to the environment:

Open up a Rascal console, enter `"import Plugin;"`, then `"main();"`. If all is well, you can now open `".ctl"` files (see the input directory of the project) and get syntax highlighting.

```
module Plugin

import <you syntax module>;
import util::IDE;
import ParseTree;
```

```

public void main() {
    registerLanguage("Controller", "ctl", <Your start symbol>(str input, loc org) {
        return parse(#<Your start symbol>, input, org);
    });
}

```

You can also parse from within the Rascal console:

```

import <Your syntax module>;
import ParseTree;
parse(#<Your start symbol>, "....")

```

The result will be a concrete syntax tree: a tree representing the structure of the input string. As you can see, it is very verbose: it includes *all* information about the source, including whitespace and comments. Sometimes it can be tedious to work with such trees. For this reason, one often defines an *abstract* syntax tree which omits details that are irrelevant and leaves only the structure of the source. This is the next step.

Abstract Syntax

Rascal uses algebraic data types (ADTs) for describing abstract syntax, as is common in functional programming language. The Rascal standard library defines a function “implode” that turns a concrete syntax tree into an abstract syntax tree.

To make this work, define an ADT that corresponds to the syntax definition in the following way:

- Every non-terminal maps to an ADT type: for non-terminal X , define “data $X =$ ”
- Every alternative of a non-terminal X maps to a constructor alternative of the ADT X , where the name of the constructor corresponds to the label of alternative.
- Every element in an alternative that is not a literal, maps to a constructor argument.
- Regular symbols $X?$, X , $X+$, $\{X \text{ “sep”}\}$, $\{X \text{ “sep”}\}+$ map to $\text{list}[T]$ where T is the type corresponding to X .
- Lexical symbols (e.g., identifiers) maps to the `str` datatype.

You can now convert concrete syntax trees to ASTs as follows:

```
pt = ... parse tree of previous snippet ....
import <Your AST module>;
implode(#<Your AST module>::<Root ADT type>, pt);
```

Not that you'll have to qualified the root type of the ADT with the module name since the name will otherwise clash with the start symbol of your grammar (Exercise: make a dedicated implode module which hides this complexity, so that you can just use implode(pt) without specifying the types.).

From now on, we will work with the AST values only.

Simple source to source transformation

Consider the following paragraph in Fowler's text:

In particular, you should note that reset events aren't strictly necessary to express Miss Grant's controller. As an alternative, I could just add a transition to every state, triggered by doorOpened, leading to the idle state. The notion of a reset event is useful because it simplifies the diagram.

What this means is, that it is possible to construct an equivalent state machine that does not depend on reset events. In this assignment you are to write a transformation that *desugars* reset events according to the quote above.

Tip: use the Rascal visit construct to transform a state machine AST.

Fact extraction

For some application, especially analysis, the tree structure of the AST is not ideal. In this assignment you will write a function that extracts a relational representation of state machines. Relations in Rascal are natural representations for graphs. And a state machine can be considered as a special kind of graph.

In order to connect the resulting analysis to the state machine visualizer, we use the following types as interfaces:

```
alias TransRel = rel[str from, str token, str to]
alias ActionRel = rel[str state, str token]
```

The first relation captures the transition structure of a state machine: it contains tuples <s, t, s'>, where s is the source state, t the triggering token, and s' the target state. The second relation captures which tokens should be output upon

entering a certain state. Note that both relations use tokens and not the names of events or actions. This means that in your extraction you should take care of looking up event/command names to find the associated tokens.

NB: You may assume that reset events have been desugared as in the previous step.

Well-formedness checking of state machines

Many programming languages have type checkers. In the case of state machines, there isn't really a notion of types. Nevertheless, it is still possible to make mistakes. In this assignment the goal is to make a checker function that detects such mistakes. This function will return a collection of error or warning messages. The data type to be used for this can be found in the standard library module "Message".

The list of things you could check for includes (but might not be limited by) the following:

- Duplicate definitions of events/commands and their tokens.
- Duplicate state definitions.
- Reset events that are used in a transition.
- Non-determinism (two transitions from the same state that fire on the same token).
- Undeclared reset events, actions, events or states.
- Unreachable states.
- Unused commands or events.

The Message data type accepts source locations. They can be retrieved from AST nodes if you add an annotation declarations to your AST module for each ADT type:

```
anno loc <ADT type>@location;
```

Now you can obtain an AST node's source location using "n@location". The locations in the Message data type are used by the IDE to do error marking.

Note: you may put the facts you've extracted from the previous assignment to good use in the analysis.

Tip: for reachability analysis use the builtin Rascal operator for transitive closure (post-fix +).

Connecting the checker to the IDE

To hook up your checker to the IDE you have to add the following line to main in Plugin.rsc:

```
registerAnnotator("Controller", check);
```

The check function (which you'll have to provide), is expected to take a (concrete!) parse tree and return an annotated parse tree. Concrete parse trees can be annotated with a set of Messages. In the check function you'll have to invoke your checker, get the collection of error message, and annotate the parse tree with it. You can annotate a tree using `x[@a=v]`, where `x` is the tree, `a` is the annotation name and `v` is the annotated value.

NB: error marking is still very experimental. You may encounter incorrect marking in some cases.

Code-generation to Java

State machine have to be executed in code somehow. One approach is to generate (Java) code. In this assignment you are to write a function that generates Java code using Rascal's built-in string templates. String templates are string literals with advanced mechanisms for interpolation:

- Expression interpolation: "Hello !"
- For loop interpolation: "abc<for (x <- [\"c\", \"d\", \"e\"]) {><}>fgh"
- If statement interpolation: "abc<if (x > 0) {><} else {><x + 1><}>"

There are multiple ways for generating code from a state machine. You may consider one of the following alternatives:

- Generate a single switch statement, which dispatches on integer constants defined for each state. Upon transition, a current-state variable is updated.
- Generate methods for each state which call other methods upon transitioning.
- Generate object instantiations according to Fowler's text.

Note that upon entering states the actions (if any) should be executed. To be able to run the code, assume the input stream is a `java.util.Scanner` object, and use `nextLine` to obtain the next token. Actions print tokens onto an output-stream, which can be any `java.io.Writer`.

Connecting the code generator to the IDE

To be able to invoke the code generator from the state machine editor, add the following lines to main in Plugin.rsc.

```
contribs = {popup(menu("Controller",[action("Generate Java", generate)]))});
registerContributions("Controller", contribs);
```

You'll have to write a generate function that actually invokes your code generator and write the result to file. Take a look in util::IDE to find out the expected signature for generate. Function for input/output can be found in the standard library IO.

Parallel merge of two state machines

The desugaring reset events (cf. above) is an instance of a simple model-to-model transformation. In this assignment we will engage in a model-to-model transformation that is slightly more complex. The goal is to take two state machines and produce a new one that runs the two original state machines in parallel.

The states in the machine resulting from merging machines S1 and S2 are identified by tuples of the states of both machines. Execution thus starts in a the initial state $\langle s_0, u_0 \rangle$ where s_0 and u_0 are the initial states of S1 and S2 respectively. Running S1 and S2 in parallel then entails the following:

- If in state $\langle s, u \rangle$, on event e , both S1 and S2 have transitions to s' , and u' , the combined machine transitions to $\langle s', u' \rangle$.
- If in state $\langle s, u \rangle$, on event e , only S1 has a transition to s' , the combined machine transitions to $\langle s', u \rangle$.
- If in state $\langle s, u \rangle$, on event e , only S2 has a transition to u' , the combined machine transitions to $\langle s, u' \rangle$.

Note: you have to decide how commands, events and reset events are combined and how, upon entering a combined state $\langle s', u' \rangle$, the actions of both s' and u' are combined.

Since the result of the parallel merge transformation is again just an ordinary state machine, you can reuse the code generator of the previous assignment *as is* to run two state machines in parallel.

Simulation of a state machine