

Mathematical Structures in Computer Science

<http://journals.cambridge.org/MSC>

Additional services for *Mathematical Structures in Computer Science*:

Email alerts: [Click here](#)

Subscriptions: [Click here](#)

Commercial reprints: [Click here](#)

Terms of use : [Click here](#)

Computing with relational machines

GÉRARD HUET and BENOÎT RAZET

Mathematical Structures in Computer Science / *FirstView* Article / July 2015, pp 1 - 20
DOI: 10.1017/S0960129515000390, Published online: 15 July 2015

Link to this article: http://journals.cambridge.org/abstract_S0960129515000390

How to cite this article:

GÉRARD HUET and BENOÎT RAZET Computing with relational machines. *Mathematical Structures in Computer Science*, Available on CJO 2015 doi:10.1017/S0960129515000390

Request Permissions : [Click here](#)

Computing with relational machines

GÉRARD HUET[†] and BENOÎT RAZET[‡]

[†]*INRIA Paris-Rocquencourt Center, France*

Email: gerard.huet@inria.fr

[‡]*Computer Science Department, Bucknell University, Lewisburg, PA 17837, USA*

Email: benoit.razet@gmail.com

Received September 2012; revised January 2015

We propose a relational computing paradigm based on Eilenberg machines, an effective version of Eilenberg's X-machines suitable for general non-deterministic computation. An Eilenberg machine generalizes a finite-state automaton, seen as its control component, with a computation component over a data domain specified as a relational algebra, its actions being interpreted as binary relations over the data domain. We show various strategies for the sequential simulation of our relational machines, using variants of the *reactive engine*. In a particular case of *finite machines*, we show that bottom-up search yields an efficient complete simulator.

Relational machines may be composed in a modular fashion, since atomic actions of one machine can be mapped to the characteristic relation of other relational machines acting as its parameters.

The control components of machines can be compiled from regular expressions. Several such translations have been proposed in the literature, which we briefly survey.

1. Introduction

Programming languages evolved from imperative notation for sequences of computer operations, operating on states of the machine and its memory, to mathematical definitions of computations. This reduced progressively the gap between applicative programming and constructive mathematics.

Functional programming actually started before the introduction of computers, with Kleene's systems of recursive definitions in arithmetic, and with Church's λ -calculus. When programmable computers appeared, computer science precursors used λ -calculus, and to a lesser extent Curry's combinatory algebra, as a foundation for programming languages. Thus, Corrado Böhm introduced the CUCH abstract programming model in 1966 (Böhm 1966).

The same year, Peter Landin's introduced in his landmark paper 'The next 700 programming languages' (Landin 1966) a crucial extension ISWIM of pure λ -calculus for use as a programming language, first by proposing the *let* notation for explicit redexes, and secondly by adding primitive control operators for conditional and recursion and primitive data structures such as booleans and integers for direct higher-order recursions. ISWIM was complemented by a polymorphic type system and the operators of cartesian product by Robin Milner, yielding a very clean functional language, ML, fit for deterministic programming.

Computing in a non-deterministic way has not been adapted to a uniform elegant programming language in the same manner so far. Of course, Dijkstra's Guarded Command Language (GCL) allows non-deterministic choice in pattern matching, but it stayed as a pencil and paper language usable for proofs in predicate transformer semantics, rather than production programming. Prolog is an attempt at logic programming seen as the non-deterministic search for the satisfiability of Horn clauses, but it is notoriously weak at controlling the search for solutions, in the absence of programmable tactics. The Prolog user can control this search solely by reordering the clauses of the program, and inserting 'cut' primitives that break the applicative semantics of the language. Later constraint programming languages suffer similar problems, and the user of such 'black-box computers' is torn between his fear of ad-hoc programming of backtracking processes in a conventional deterministic language, and his loathing of being the hostage of generic tactics ill-fitted to his particular problem.

Actually non-deterministic computation is fairly straightforward if one has a clear notion of the search space of the problem, if fairness is enforced (either by termination or by non-starvation tactics), and if higher-order applicative programming is used. Higher-order gives you continuations and streams for your enumerative processes, and when you resume search at backtracking points you retrieve the stored state of the computing thread, without risk of corruption by side effects. This style of programming may be systematized within a very general framework for relational computation, using ML as its core deterministic engine. Furthermore, the control of such non-deterministic relational search fits well within a very elegant proposal of Samuel Eilenberg for generalizing finite-state automata with more abstract X-machines (Eilenberg 1974).

We propose here this relational programming paradigm, formally defined as an ML library of parametric modules, as a tribute to Corrado Böhm, Peter Landin and Samuel Eilenberg.

2. Machines

2.1. Relational machines

We shall define a notion of abstract machine inspired from the work of Eilenberg (X-machines, presented in his *magnum opus* on automata theory (Eilenberg 1974)). Our machines are nondeterministic in nature. They comprise a *control component*, similar to the transitions state diagram of a (non-deterministic) automaton. These transitions are labelled by action generators[†]. Action expressions over free generators, generalizing regular expressions from the theory of languages, provide a specification language for the control component of machines. A program, or action expression, compiles into control components according to various translations. Control components in their turn may be compiled further into transition matrices or other representations.

[†] Our terminology of 'actions' for the computable relations interpreting the machine operations is inspired by the *actions* algebras of Pratt (1991), which we think are the natural models for our machines' semantics.

Our machines also comprise a *data component*, endowed with a relational semantics (typically, an action algebra in the sense of Pratt). That is, we interpret action generators by semantic attachments to binary relations over the data domain. These relations are themselves represented as functions from data elements to streams of data elements. This applicative apparatus replaces the imperative components of traditional automata (tapes, reading head, counters, stacks, etc) by clear mathematical notions. On the other hand, they give a concrete computational framework in which to interpret the non-constructive mathematical model of Eilenberg's original *X*-machines.

We shall now formalize these notions in a way that will exhibit the symmetry between control and data. First of all, we postulate a finite set Σ of parameters standing for the names of the primitive operations of the machine, called (action) *generators*.

For the control component, we postulate a finite set S of states and a *transition relation map* interpreting each generator as a (binary) relation over S . This transition relation interpretation is usually presented as a curried *transition function* δ mapping each state in S to a finite set of pairs (a, q) with a a generator and q a state. This set can in turn be implemented as a finite list of such pairs.

Finally, we select in S a set of *initial states* and a set of *accepting states*.

For the data component, we postulate a set D of data values and a *computation relation map* interpreting each generator as a (binary) relation over D . Similarly as for the control component, we shall curry this relation map as a *computation function* mapping each generator a in Σ to a function $\rho(a)$ in $D \rightarrow \wp(D)$. We call this computation function ρ the *semantics*. Now the situation is different from control, since D and thus $\wp(D)$ may be infinite. In order to have a constructive characterization, we shall assume that D is recursively enumerable, and that each $\rho(a)$ maps $d \in D$ to a recursively enumerable subset of $\wp(D)$. We shall represent such subsets as progressively computed streams of values, as we shall explain in due time.

2.2. Progressive relations as streams

We recall that a recursively enumerable subset of \mathbb{N} is the range of a partial recursive function in $\mathbb{N} \rightarrow \mathbb{N}$, or equivalently it is either empty or the range of a (total) recursive function in $\mathbb{N} \rightarrow \mathbb{N}$. None of these two definitions is completely satisfactory, since in the first definition we may loop on some values of the parameter, obliging us to dovetail the computations in order to obtain a sequence of elements that completely enumerates the set, and in the second we may stutter enumerating the same element in multiple ways. This stuttering cannot be totally eliminated without looping, for instance for finite sets. Furthermore, demanding total functions is a bit illusory. It means either we restrict ourselves to a non-Turing complete algorithmic description language (such as primitive recursive programs), or else we cannot decide the totality of algorithms demanded by the definition.

We shall here assume that our algorithmic description language is ML, in other words, typed lambda-calculus evaluated in call by value with a recursion operator, inductive types and parametric modules. More precisely, we shall present all our algorithms in Pidgin ML, actually the so-called 'revised syntax' of Objective Caml.

In this framework, we can define computable streams over a parametric datatype data as follows:

```
type stream 'data =
  [ Void
  | Stream of 'data and delay 'data
  ]
and delay 'data = unit → stream 'data;
```

This definition expresses that a stream of data values is either Void, representing the empty set, or else a pair Stream(d,f) with d of type data, and f a frozen stream value, representing the set $\{d\} \cup F$, where F may be computed as the stream f(), where () is syntax for the canonical element in type unit. We overload the set membership notation with stream membership and simply write $d \in str$ when d is a data value enumerated by the stream str.

Using this inductive parametric datatype, we may now define progressive relations by the following type:

```
type relation 'data = 'data → stream 'data;
```

2.3. Kernel machines

We now have all the ingredients to define the module signature of *kernel machines*:

```
module type Kernel = sig
  type generator;
  type data;
  type state;
  value transition: state → list (generator × state);
  value initial: list state;
  value accept: state → bool;
  value semantics: generator → relation data;
end;
```

In the following, we shall continue to use Σ (resp. D, S, δ, ρ) as shorthand for generator (resp. data, state, transition, semantics). We also write I for initial and T for the set of accepting states (for which the predicate accept is true).

A machine \mathcal{M} is like a black box, which evolves through series of non-deterministic computation steps. At any point of the computation, its status is characterized by the pair (s, d) of its *current state* $s \in S$ and its *current data value* $d \in D$. Such a pair is called a *cell*.

A computation step issued from cell (s, d) consists in choosing a transition $(a, s') \in \delta(s)$ and a value $d' \in \rho(a)(d)$. If any of these choices fails, because the corresponding set is empty, the machine is said to be *blocked*; otherwise, the computation step succeeds, and the machine has as status the new cell (s', d') . We write $(s, d) \xrightarrow{a} (s', d')$.

A *computation path* is a finite sequence of such computations steps:

$$(s_0, d_0) \xrightarrow{a_1} (s_1, d_1) \xrightarrow{a_2} (s_2, d_2) \cdots \xrightarrow{a_n} (s_n, d_n).$$

The computation is said to be *accepting* whenever $s_0 \in I$ and $s_n \in T$, in which case we say that the machine *accepts* input d_0 and *computes* output d_n . We may also say that d_n is a *solution* of d_0 by the machine \mathcal{M} and write the corresponding predicate $Solution(d_0, d_n)$. Note that (d_0, d_n) belongs to the graph of the composition of relations labelling the path: $\rho(a_1) \circ \rho(a_2) \circ \dots \circ \rho(a_n)$.

Now, we define the *characteristic relation* of a machine \mathcal{M} as the set of pairs of accepted input with corresponding computed output, which we denote $\|\mathcal{M}\|$:

$$\|\mathcal{M}\| = \bigcup \{(d_0, d_n) \mid (s_0, d_0) \xrightarrow{a_1} \dots \xrightarrow{a_n} (s_n, d_n) \text{ is an accepting computation path}\}.$$

Characteristic relations are the relational interpretation over the data domain D of the action language recognized by the underlying automaton. They allow us to compose our machines in modular fashion.

We have thus a very general model of relational calculus. Our machines compute relations over the data domain D , and we shall thus speak of *D-machines*. The ‘machine language’ has the action generators for instructions. Actions compose by computation. Furthermore, a high-level programming language for relational calculus may be designed as an action calculus. The obvious point of departure for this calculus is to consider regular expressions, in other words the free Kleene algebra generated by the set of generators. We know from automata theory various translations from regular expressions to finite-state automata. Every such translation gives us a compiler of our relational algebra into the control components of our machines: S , δ , I and T . The data components, D and ρ , offer a clean mathematical abstraction over the imperative paraphernalia of classical automata: reading heads, tapes, etc. And we get immediately a programming language enriching the machine language of primitive actions by composition, iteration and choice.

Indeed, a finite automaton over alphabet Σ is readily emulated by the machine with generator set Σ having its state transition graph as its control component, and admitting the free monoid of actions Σ^* for data domain. Each generator a is interpreted in the semantics as the (functional) relation $\rho(a) = L_a^{-1} =_{def} \{(a \cdot w, w) \mid w \in \Sigma^*\}$ which ‘reads the input tape.’ And indeed the language recognized by the automaton is retrieved as the composition of actions along all accepting computations. Here, the data computation is merely a trace of the different states of the ‘input tape.’

This example is a particularly simple one, in which data computation is deterministic, since in this case $\rho(a)$ is a partial function. We may say that such a machine is ‘data driven.’ Control will be deterministic too, provided the underlying automaton is deterministic, since every $\delta(s)$ will then have a unique non-blocking transition. But note that the same control component could be associated with different semantics. For instance, with $\rho(a) = R_a =_{def} \{(w, w \cdot a) \mid w \in \Sigma^*\}$, the machine will enumerate with its accepting computations the regular language recognized by the automaton.

2.4. Modular construction of machines

Now that we understand that a *D-machine* implements a relation over D , we may compose machines vertically as follows. Let \mathcal{A} be a (non-deterministic) automaton over alphabet Σ , and for every $a \in \Sigma$ let \mathcal{N}_a be a *D-machine* over some generator set Σ_a . We may now

turn \mathcal{A} into a D -machine over generator set Σ by taking \mathcal{A} as its control component, and extending it by a data component having as semantics the function mapping $a \in \Sigma$ to $\|\mathcal{N}_a\|$.

We may thus construct large machines from smaller ones computing on the same data domain. A typical example of application for computational linguistics is to do morphological treatment (such as segmentation and tagging of some corpus) in a lexicon-directed way. The alphabet Σ defines the lexical categories or parts of speech, each machine \mathcal{N}_a implements access to the lexicon of category a , the automaton \mathcal{A} defines the morphological geometry, and the composite machine \mathcal{M} implements a lexicon-directed parts-of-speech tagger. By appropriate extension of the lexicon machines \mathcal{N}_a , morphophonemic treatment at the junction of the words may be effected, such as complete sandhi analysis for Sanskrit. This was the motivating example for which the Zen toolkit was designed (Huet 2005; Huet and Razet 2006).

2.5. Interfaces

What we described so far is the Eilenberg machine *kernel*, consisting of its control and data elements. We may complete this description by an *interface*, composed of an input domain D_- , an output domain D_+ , an input relation $\phi_- : D_- \rightarrow D$ and an output relation $\phi_+ : D \rightarrow D_+$. A kernel machine \mathcal{M} completed by this interface I defines a relation $\phi(\mathcal{M}, I) : D_- \rightarrow D_+$ by composition:

$$\phi(\mathcal{M}, I) = \phi_- \circ \|\mathcal{M}\| \circ \phi_+.$$

3. Finite machines

We shall now present an important special case of machines that exhibit a finite behaviour.

The relation $\rho : D \rightarrow D'$ is said to be *locally finite* if for every $d \in D$ the set $\rho(d)$ of elements of D' related to d is finite. The machine \mathcal{M} is said to be *locally finite* if every generating relation $\rho(a)$ is locally finite (Huet 1980). The machine \mathcal{M} is said to be *n etherian* if all its computations are finite.

We remark that a machine is n etherian when its data domain D is a well-founded ordering for the order relation $>$ generated by

$$d > d' \iff \exists a \in \Sigma \ d' \in \rho(a)(d).$$

Indeed, if there existed an infinite computation, there would exist an infinite subsequence going through the same state. But the converse is not true, since a machine may terminate for reasons depending of its control.

Finally, we say that a machine is *finite* if it is locally finite and n etherian.

We say that a machine kernel is *deterministic* (Eilenberg 1974) iff $|I| \leq 1$ and for each cell value (s, d) occurring in a computation there exists at most one computation transition issued from it, i.e. if $\delta(s)$ is a set of pairs $\{(a_1, s_1), (a_2, s_2), \dots, (a_n, s_n)\}$ such that for at most one $1 \leq k \leq n$ the set $\rho(a_k)(d)$ is nonempty, and if such k exists then $\rho(a_k)(d)$ is a singleton. This condition demands that on one hand the transition relation of the underlying automaton be a partial function, that is the automaton must be deterministic,

and on the other hand that the relations leading out of a state s be partial functions over the subset of D that is reachable by computations leading to s . We extend this property to a machine with interface by requiring that its input relation ϕ_- and its output relation ϕ_+ be partial functions. We remark that a deterministic machine may nevertheless generate several solutions, since a terminal cell is not necessarily blocking further computation. Under certain extra conditions, such a machine computes a partial function (see chapter X, Section 8 of Eilenberg's treatise (Eilenberg 1974)).

3.1. Examples

3.1.1. *Non deterministic finite automata* Let us consider a non-deterministic automaton \mathcal{A} with parameters (S, δ, I, T) . We construct an Eilenberg machine \mathcal{M} solving the word problem for the regular language $|\mathcal{A}|$ recognized by the automaton. \mathcal{M} has Σ for generating set, and it takes \mathcal{A} for its control component. For the data component, we take $D = \Sigma^*$, and the semantics is defined as $\rho(a) = L_a^{-1} =_{def} \{(a \cdot w, w) \mid w \in \Sigma^*\}$, as explained above.

We may check that $\rho(w) = 1$ iff $w \in |\mathcal{A}|$. It is easy to check that \mathcal{M} is finite, since data decreases in length, and semantics is a partial function. When \mathcal{A} is a deterministic automaton, \mathcal{M} is a deterministic machine.

Another machine with the same control component may be defined to enumerate all the words in set $|\mathcal{A}|$. In general it will be neither finite, nor deterministic.

3.1.2. *Rational transducers* Let Σ and Γ be two finite alphabets. A transducer $\mathcal{A} : \Sigma \Rightarrow \Gamma$ is similar to a (non-deterministic) automaton, whose transitions are labelled with pairs of words in $D = \Sigma^* \times \Gamma^*$. Let Ω be the (finite) set of labels occurring as labels of the transitions of \mathcal{A} . The transition graph of \mathcal{A} may thus be considered as an ordinary non-deterministic automaton over generator alphabet Ω , and constitutes the control component of the machines we shall define to solve various transduction tasks.

We recall that a transducer 'reads its input' on an input tape representing a word in Σ^* and 'prints its output' on an output tape representing a word in Γ^* . On transition (w, w') it reads w off the input tape, and if successful appends w' to its output tape. If by a succession of transitions starting from an initial state with input i and empty output it reaches an accepting state with empty input and output o , we say that (i, o) belongs to the *rational relation* in $\Sigma \Rightarrow \Gamma$ recognized by the transducer \mathcal{A} , which we shall write $|\mathcal{A}|$. We shall now solve various decision problems on $|\mathcal{A}|$ using machines that use \mathcal{A} for control and D for data, but replace the tapes by various semantic functions:

1. Recognition. Given $(w, w') \in D$, decide whether $(w, w') \in |\mathcal{A}|$.
2. Synthesis. Given $w \in \Sigma^*$, compute its image $|\mathcal{A}|(w) \subset \Gamma^*$.
3. Analysis. Given $w \in \Gamma^*$, compute the inverse image $|\mathcal{A}^{-1}|(w) \subset \Sigma^*$.

Recognition. The semantics ρ is defined by $\rho(\sigma, \gamma) = L_\sigma^{-1} \times L_\gamma^{-1}$. Like for ordinary automata we obtain a finite machine, provided the transducer has no transition labelled (ϵ, ϵ) , since at least one of the two lengths decreases. We choose as interface $D_- = \Sigma^* \times \Gamma^*$, $\phi_- = Id_{\Sigma^* \times \Gamma^*}$, $D_+ = 0, 1$, and $\phi_+(w, w') = 1$ iff $w = w' = \epsilon$.

Synthesis. The semantics ρ is defined by $\rho(\sigma, \gamma) = L_\sigma^{-1} \times R_\gamma$, with $R_\gamma =_{def} \{(w, w \cdot \gamma) \mid w \in \Gamma^*\}$. We choose as interface $D_- = \Sigma^*$, $\phi_- = \{(w, (w, \epsilon)) \mid w \in \Sigma^*\}$, $D_+ = \Gamma^*$, and

$\phi_+ = \{(\epsilon, w'), w' \mid w' \in \Gamma^*\}$. We get $|\mathcal{A}| = \phi_- \circ \|\mathcal{M}\| \circ \phi_+$. Such a machine is locally finite, since relations L_σ^{-1} and R_γ are partial functions. However, it may not be noetherian, since there may exist transitions labelled with actions (ϵ, w) . Actually the machine is noetherian iff cycles of such transitions do not occur, i.e. iff the set $|\mathcal{A}|(w)$ is finite for every $w \in \Sigma^*$ – see Razet (2008).

Analysis. Symmetric to synthesis, replacing L_σ^{-1} by R_σ and R_γ by L_γ^{-1} .

3.2. Reactive engine

We may simulate the computations of a finite Eilenberg machine by adapting the notion of *reactive engine* from the Zen library (Huet 2002, 2005; Huet and Razet 2006; Razet 2008). The engine is a deterministic simulator of the non-deterministic machine.

We start with a simple depth-first search engine, appropriate for finite machines. We define the engine as an ML functor that is a module taking as parameter a kernel machine.

```

module Engine (Machine: Kernel) = struct
open Machine;

type choice = list (generator  $\times$  state);

(* We stack backtrack choice points in a resumption *)
type backtrack =
  [ React of data and state
  | Choose of data and choice and delay data and state
  ]
and resumption = list backtrack;

(* The 3 internal loops of the reactive engine (using terminal calls) *)
(* react : data  $\rightarrow$  state  $\rightarrow$  resumption  $\rightarrow$  stream data *)
value rec react d q res =
  let ch = transition q in
  (* We need to compute 'choose d ch res' but first we deliver data d
  to the stream of solutions when state q is accepting *)
  if accept q
  then Stream d (fun ()  $\rightarrow$  choose d ch res) (* Solution d found *)
  else choose d ch res
(* choose : data  $\rightarrow$  choice  $\rightarrow$  resumption  $\rightarrow$  stream data *)
and choose d ch res =
  match ch with
  [ []  $\rightarrow$  resume res
  | [(g, q') :: rest]  $\rightarrow$  match semantics g d with
    [ Void  $\rightarrow$  choose d rest res
    | Stream d' del  $\rightarrow$  react d' q' [ Choose d rest del q' :: res ]
    ]
  ]

```

```

]
(* The scheduler that backtracks in depth-first exploration *)
(* resume: resumption → stream data *)
and resume res =
  match res with
  [ [] → Void
  | [ React d q :: rest ] → react d q rest
  | [ Choose d ch del q' :: rest ] →
    match del () with (* We thaw the delayed stream of solutions *)
    [ Void → choose d ch rest (* And we look for next pending choice *)
    | Stream d' del' → react d' q' [ Choose d ch del' q' :: rest ]
    ]
  ]
;
(* Simulating the characteristic relation : relation data *)
value simulation d =
  let rec init_res l acc =
    match l with
    [ [] → acc
    | [ q :: rest ] → init_res rest [ React d q :: acc ]
    ] in
  resume (init_res initial [])
;
end; (* module Engine *)

```

This reactive engine has a very simple management of pending choices, since the backtrack choices are stored on a resumption stack, last-in first-out. It is very fast, since the ML compiler replaces the terminal calls by jumps. It is the workhorse of the original motivating application, Sanskrit sentence segmentation (Huet 2005).

3.3. Correctness, completeness, certification

A proof of correctness and completeness of this simulator was given by the first author in the case of segmentation transducers (Huet 2005). Razet generalized the proof to the general case of finite Eilenberg machines, and formalized it in the Coq proof assistant (Razet 2011). From the formal proof object it is possible to extract mechanically ML algorithms identical to the ones we showed above.

Here is the exact correctness and completeness result formally proved using the Coq proof assistant:

Theorem 1. Let \mathcal{M} be a machine, the simulation function computes its characteristic relation $\|\mathcal{M}\|$. That is, for all $d \in D$ we have:

$$\text{Solution}(d, d') \Leftrightarrow d' \in (\text{simulation } d)$$

4. A general reactive engine, driven by a strategy

When a machine is not finite, and in particular when there are infinite computation sequences, the above bottom-up engine may loop, and the simulation is not complete. In order to remedy this problem, we shall change the specific last-in first-out policy of resumption management, and replace it by a more general strategy, given as an extra parameter of the machine.

```

module Engine (Machine: Kernel) = struct
open Machine;

type choice = list (generator × state);
(* We separate the control choices and the data relation choices *)
type backtrack =
  [ React of data and state
  | Choose of data and choice
  | Relate of stream data and state
  ];

```

In the previous reactive engine using a depth-first search we had a resumption datatype for the backtrack stack. Let us generalize this stack and specify resumption as an abstract data type, encoded in a module signature Resumption:

```

module type Resumption = sig
  type resumption;
  value empty: resumption;
  value pop: resumption → option (backtrack × resumption);
  value push: backtrack → resumption → resumption;
end;

```

We now modify the reactive engine, so that resumption management is governed by the given strategy. The reactive engine is a functor Strategy parameterized by a module of type Resumption.

```

module Strategy (R: Resumption) = struct
open R;

(* react : data → state → resumption → stream data *)
value rec react d q res =
  let ch = transition q in
  if accept q (* Solution d found? *)
  then Stream d (fun () → resume (push (Choose d ch) res))
  else resume (push (Choose d ch) res)
(* choose : data → choice → resumption → stream data *)
and choose d ch res =
  match ch with
  [ [] → resume res

```

```

| [ (g, q') :: rest ] →
  resume (push (Relate (semantics g d) q') (push (Choose d rest) res))
]
(* relate : stream data → state → resumption → stream data *)
and relate str q res =
  match str with
  [ Void → resume res
  | Stream d del →
    resume (push (React d q) (push (Relate (del ()) q) res))
  ]
(* resume : resumption → stream data *)
and resume res =
  match pop res with
  [ None → Void
  | Some (b, rest) →
    match b with
    [ React d q → react d q rest
    | Choose d ch → choose d ch rest
    | Relate str q → relate str q rest
    ]
  ]
;
(* simulation : relation data *)
value simulation d =
  let rec init_res l acc =
    match l with
    [ [] → acc
    | [ q :: rest ] → init_res rest (push (React d q) acc)
    ] in
  resume (init_res initial empty)
;
end; (* module Strategy *)

```

4.1. A few typical strategies

We now give a few variations on search strategies. First of all, we show how the original depth-first reactive engine may be obtained by a `DepthFirst` strategy module, adequate for finite Eilenberg machines.

```

module DepthFirst: Resumption = struct
  type resumption = list backtrack;
  value empty = [];
  value push b res = [ b :: res ];
  value pop res = match res with

```

```

  [ [] → None
  | [ b :: rest ] → Some (b,rest)
  ];
end; (* module DepthFirst *)

```

We remark that there is a cost involved in reversing the *input* list above, and that suppressing this reversing operation yields a more efficient machine working in a ‘boustrophedon’ manner, while preserving fairness:

```

module Fair: Resumption = struct
type resumption = (list backtrack × list backtrack);
value empty = ([], []);
value push b res = let (left,right) = res in (left, [ b :: right ]);
value pop res =
  let (left,right) = res in
  match left with
  [ [] → match right with
    [ [] → None
    | [ r :: rrest ] → Some (r, (rrest, []))
    ]
  | [ l :: lrest ] → Some (l, (lrest,right))
  ]
;
end; (* module Fair *)

```

Finally, we examine the special case of deterministic machines. The following simple Det tactic is adapted to this case.

```

module Det: Resumption = struct
type resumption = list backtrack;
value empty = [];
value push b res = match b with
  [ React _ _ → [ b :: res ]
  | Choose _ _ → [ b ] (* cut : the list contains only one element *)
  | Relate _ _ → res (* no other delay *)
  ];
value pop res = match res with
  [ [] → None
  | [ b :: rest ] → Some (b,rest)
  ];
end; (* module Det *)

```

Now, we may build the various modules encapsulating the various strategies.

```

module FEM = Strategy DepthFirst; (* The bottom-up engine *)
module Fair_Engine = Strategy Fair; (* A fair engine *)
module Deterministic_Engine = Strategy Det; (* The deterministic one *)
end; (* module Engine *)

```

Nevertheless, a complete evaluation strategy in the general case demands a more complex stream definition, where the computation is sliced into provably terminating states. This extension is discussed in the second author's thesis (Razet 2009).

5. From regular expressions to automata

Our motivation here is the design of a language for describing the control part of Eilenberg machines. The control part of Eilenberg machines is a finite automaton. It naturally leads us to *regular expressions* and their translations into finite automata.

There have been more than 50 years of research on the problem of compilation (or translation) of regular expressions into automata. It started with Kleene who stated the equivalence between the class of languages recognized by finite automata and the class of languages defined by regular expressions. This topic is particularly fruitful because it has applications to string-search algorithms, circuits, synchronous languages, computational linguistics, etc. This wide range of applications leads one to several automata and regular-expressions variants.

Usually, an algorithm compiling regular expressions into automata is described in an imperative programming style for managing states and edges: states are allocated, merged or removed and so on concerning the edges. However, and this may seem somewhat surprising, it is possible to describe each of the well-known algorithms in an applicative manner, while preserving its computational complexity. This methodology leads to formal definitions of the algorithms exhibiting important invariants, an essential step towards their formal verification.

We focus on fast translations, whose time complexity is linear or quadratic with respect to the size of the regular expression. First, we present *Thompson's algorithm* (Thompson 1968) and then we review other algorithms that may be put to use by our methodology.

Let us mention Brzozowski's algorithm (Brzozowski 1964), which translates a regular expression (even with Boolean operators) into a *deterministic* automaton. Unfortunately, its complexity is theoretically exponential in space and time. Nevertheless, it introduced the notion of regular expression *derivative* which is a fundamental idea pervading other algorithms.

5.1. Thompson's algorithm

Thompson presented his algorithm in 1968 and it is one of the most famous translations. It computes a finite non-deterministic automaton with ϵ -moves in linear time.

Let us first define regular expressions as the following datatype:

```

type regexp 'a =
  [ One
  | Symb of 'a
  | Union of regexp 'a and regexp 'a
  | Conc of regexp 'a and regexp 'a
  | Star of regexp 'a
  ];

```

The constructor `One` of arity 0 corresponds to the identity relation. The following constructor `Symb` of arity 1 is the node for the ‘a’ generator. The type for the generator is abstract as expressed by the type parameter ‘a’ in the definition. The two following constructors are `Union` and `Conc` of arity 2 and describe union and concatenation operations. The last constructor `Star` is for the iteration or Kleene’s star operator.

Now that we have given the datatype for the input of our algorithm, let us present the datatype for the output (automata). We choose to implement states of the automaton with integers:

```
type state = int;
```

Automata obtained by Thompson’s algorithm are nondeterministic and furthermore may contain ϵ -moves. We shall implement the control graph of such non-deterministic automata as a list of fanout pairs associating a list of labelled transitions to a state. This method amounts to encoding a set of edges $s \xrightarrow{a} s'$ or triples (s, a, s') as an association list.

```
type fanout 'a = (state × list (label 'a × state))
and label 'a = option 'a
and transitions 'a = list (fanout 'a)
;
type automaton 'a = (state × transitions 'a × state);
```

A label is of type `option 'a` because it may either be an ϵ -move of value `None` or a generator `a` of value `Some a`. Note that even if they are nondeterministic, the automata we consider have only one initial and one accepting state.

We shall instantiate the transition function of the control component of our machines by composing the transitions list component of the constructed automaton with the primitive `List.assoc`, as we shall show later in Section 6.

The algorithm performs a recursive traversal of the expression. It is presented in the order of the datatype definition: 1, generator, union, concatenation and Kleene’s star.

```
(* thompson: regexp 'a → automaton 'a *)
value thompson e =
  let rec analyse e t n =
    (* e is current regexp, t accumulates the state space,
       n is the latest created location *)
    match e with
    [ One → let n1=n+1 and n2=n+2 in
      (n1, [ (n1, [ (None, n2) ]) ] :: t ], n2)
    | Symb s → let n1=n+1 and n2=n+2 in
      (n1, [ (n1, [ (Some s, n2) ]) ] :: t ], n2)
    | Union e1 e2 →
      let (i1,t1,f1) = analyse e1 t n in
      let (i2,t2,f2) = analyse e2 t1 f1 in
      let n1=f2+1 and n2=f2+2 in
      (n1, [ (n1, [ (None, i1); (None, i2) ]) ] ::
        [ (f1, [ (None, n2) ]) ] ::
```

```

                [ (f2, [ (None, n2) ]) :: t2 ] ] ], n2)
| Conc e1 e2 →
  let (i1,t1,f1) = analyse e1 t n in
  let (i2,t2,f2) = analyse e2 t1 f1 in
  (i1, [ (f1, [ (None, i2) ]) :: t2 ], f2)
| Star e1 →
  let (i1,t1,f1) = analyse e1 t n in
  let n1=f1+1 and n2=f1+2 in
  let t1' = [ (f1, [ (None, i1); (None, n2) ]) :: t1 ] in
  (n1, [ (n1, [ (None, i1); (None, n2) ]) :: t1' ], n2)
] in
analyse e [] 0
;

```

The algorithm constructs the automaton from the regular expression with a single recursive traversal of the expression. States are created at each node encountered in the expression: each constructor creates two states except the concatenation `Conc` that does not create any state. Notice the invariant of the recursion: each regular subexpression builds an automaton (i, fan, f) with $0 < i < f$ and $dom(fan) = [k..f - 1]$. States are allocated so that disjoint subexpressions construct disjoint segments $[i..f]$. This invariant of the thompson function implies that we have to add a last (empty) fanout for the final state.

(* thompson_alg: regexp 'a → automaton 'a *)

```

value thompson_alg e =
  let (i,t,f) = thompson e in
  (i, [(f, []) :: t], f)
;

```

The function `thompson_alg` implements Thompson's algorithm in linear time and space because it performs a unique traversal of the expression.

5.2. Other algorithms

We have seen that Thompson's algorithm is linear, produces an automaton of size linear in the size of the regular expression, and can be implemented in an applicative manner. Let us mention also Berry and Sethi's algorithm (Berry and Sethi 1986) that computes a non-deterministic automaton (without ϵ -move), more precisely a *Glushkov* automaton. This construction is quadratic and we provided an implementation of it in ML (Huet and Razet 2006). In 2003, Ilie and Yu (Ilie and Yu 2003) introduced the Follow automata which are also non-deterministic automata. Actually, Champarnaud et al. (2006) showed that the Follow automaton is a quotient of the one produced by the Berry–Sethi algorithm (i.e. some states are merged together) and they provide an algorithm implementing the Follow construction in quadratic time. The applicative implementation of the Berry–Sethi algorithm may be extended to yield the Follow automaton (Razet 2009). Finally, in 1996

Antimirov proposed an algorithm (Antimirov 1996) that compiles even smaller automata than the ones obtained by the Follow construction, provided the input regular expression is presented in *star normal form* (as defined by Brüggemann-Klein (1993)). The algorithm presented originally was polynomial in $O(n^5)$ but Champarnaud and Ziadi (2001, 2002) proposed yet another implementation in quadratic time.

It is possible to validate these various compiling algorithms using some of the algebraic laws of Pratt's action algebras (Pratt 1991). In particular, using idempotency to collapse states will indicate that the corresponding construction does not preserve the notion of multiplicity of solutions. Furthermore, such a notion of multiplicity, as well as weighted automata modeling statistical properties, generalize to the treatment of valuation semi-rings, for which Allauzen and Mohri (Allauzen and Mohri 2006) propose extensions of the various algorithms. Recently, Fischer et al. (2010) presented a functional program implementing efficiently the matching problem for weighted regular expressions.

6. A worked-out example

We briefly discussed above how to implement as a machine a finite automaton recognizing a regular language. We may use for instance Thompson's algorithm to compile the automaton from a regular expression defining the language. This example will show that recognizing the language and generating the language are two instances of machines which share the same control component, and vary only on the data domain and its associated semantics. Furthermore, we show in the recognition part that we may compute the multiplicities of the analysed string. However, note that this is possible only because Thompson's construction preserves this notion of multiplicity.

Let us work out completely this method with the regular language defined by the regular expression $(a^*b + aa(b^*))^*$.

```

value exp =
  let a = Symb 'a' and b = Symb 'b' in
  let astarb = Conc (Star a) b and aabstar = Conc a (Conc a (Star b)) in
  Star (Union astarb aabstar)
;
value (i,fan,t) = thompson_alg exp;
value graph n = List.assoc n fan;
value delay_eos = fun () → Void;
value unit_stream x = Stream x delay_eos;

module AutoRecog = struct
  type data = list char;
  type state = int;
  type generator = option char;
  value transition = graph;
  value initial = [ i ];
  value accept s = (s = t);

```

```

value semantics c tape = match c with
  [ None → unit_stream tape
  | Some c → match tape with
    [ [] → Void
    | [ c' :: rest ] → if c = c' then unit_stream rest else Void
    ]
  ];
end (* AutoRecog *)
;
module LanguageDeriv = Engine AutoRecog
;
(* The Recog module controls the output of the sub-machine
  LanguageDeriv, insuring that its input is exhausted *)
module Recog = struct
  type data = list char;
  type state = [ S1 |S2 | S3 ];
  type generator = int;
  value transition = fun
    [ S1 → [ (1,S2) ] |S2 → [ (2,S3) ] |S3 → [ ] ];
  value initial = [ S1 ];
  value accept s = (s = S3);
  value semantics g tape = match g with
    [ 1 → LanguageDeriv.Fair_Engine.simulation tape
    | 2 → if tape = [] then unit_stream tape else Void
    | _ → assert False
    ];
end (* Recog *)
;
module WordRecog = Engine Recog (* The machine recognizing the language *)
;
module AutoGen = struct
  type data = list char;
  type state = int;
  type generator = option char;
  value transition = graph;
  value initial = [ i ];
  value accept s = (s = t);
  value semantics c tape = match c with
    [ None → unit_stream tape
    | Some c → unit_stream (List.append tape [c])
    ];
end (* AutoGen *)
;

```

```

module WordGen = Engine AutoGen (* The machine generating the language *)
;
(* Service functions on character streams for testing *)
(* print char list *)
value print_cl l = do { List.iter print_char l; print_string ",_" };
value iter_stream f str =
  let rec aux str = match str with
    [ Void → ()
    | Stream v del → let () = f v in aux (del ())
    ] in aux str
;
value cut str n =
  let rec aux i str = if i ≥ n then Void else match str with
    [ Void → Void
    | Stream v del → Stream v (fun () → aux (i+1) (del ()))
    ] in aux 0 str
;
value count s =
  let rec aux s n = match s with
    [ Void → n
    | Stream _ del → aux (del ()) (n+1)
    ] in aux s 0
;
print_string "Recognition of word 'aaa' with multiplicity: ";
print_int (count (WordRecog.FEM.simulation ['a' ; 'a' ; 'a' ; 'a']));
print_string "\nRecognition of word 'aab' with multiplicity: ";
print_int (count (WordRecog.FEM.simulation ['a' ; 'a' ; 'b']));
print_string "\nFirst 10 words in L in a complete enumeration: \n";
iter_stream print_cl (cut (WordGen.Fair_Engine.simulation []) 10);

```

We now show the output of executing the above code:

```

Recognition of word aaa with multiplicity: 1
Recognition of word aab with multiplicity: 3
First 10 words in L in a complete enumeration
, b, ab, aa, aab, aab, aaab, aabb, bb, aaaab,

```

The running-time of the reactive engine on these small examples is negligible. However, the reactive engine performs a backtracking search that has an exponential complexity (this exponential behaviour is observable using longer words in the recognition problem). Considering the generality of the relational machine model we propose, the backtracking search is an adequate technique for solving general problems. For specific problems, there might exist specific algorithms reducing the complexity; for instance, the recognition problem for automata on words can be solved in $O(mn)$ with m the size of the regular expression and n the length of the word (Fischer et al. 2010).

7. Conclusion

We have presented a general model of non-deterministic computation based on a computable version of Eilenberg machines. Such relational machines complement a non-deterministic finite-state automaton over an alphabet of relation generators with a semantics function interpreting each relation functionally as a map from data elements to streams of data elements. We have surveyed several algorithms that permit to compile the control component of our machines from regular expressions. The data component is implemented as an ML module consistent with an `Kernel` interface. We have shown how to simulate our non-deterministic machines with a reactive engine, parameterized by a strategy. Under appropriate fairness assumptions of the strategy, the simulation is complete. An important special case is that of finite machines, for which the bottom-up strategy is complete, while being efficiently implemented as a flowchart algorithm.

We believe this applicative model of relational computing is a sound general basis for non-deterministic search processes. It encompasses the usual applications to parsing/recognition but also to generation of formal languages. It also applies to logic programming, constraints processing, database querying, and proof search for automated proof assistants. It provides a clean framework in which to develop applications to natural-language processing and similar ‘artificial intelligence’ problems. The flexible nature of the search strategy parameter allows one to account for statistical-optimisation techniques such as hidden Markov chains. Extensions of the action algebras to numerical operators (max, plus) should allow the adaptation of these techniques to important operations research applications such as optimisation. Finally, the ubiquitous nature of relations ought to allow the extension of this model to various models of distributed processing.

Acknowledgement

Most of this material is extracted from the doctoral research of the second author (Razet 2009). A preliminary version of this paper was presented as a tutorial at ICON’2008 in Pune, Maharashtra, in December 2008. The Zen toolkit is available as an open-source Ocaml library[†]. The Sanskrit Heritage platform[‡] exhibits a real-scale application demonstrating the effectiveness of this technology for Sanskrit segmentation.

References

- Allauzen, C. and Mohri, M. (2006). A unified construction of the Glushkov, follow, and Antimirov automata. In: Proceedings of the 31st International Symposium on Mathematical Foundations of Computer Science (MFCS 2006). *Springer Lecture Notes in Computer Science* **4162** 110–121.
- Antimirov, V. (1996). Partial derivatives of regular expressions and finite automaton constructions. *Theoretical Computer Science* **155** (2) 291–319.
- Berry, G. and Sethi, R. (1986). From regular expressions to deterministic automata. *Theoretical Computer Science* **48** (1) 117–126.

[†] <http://sanskrit.inria.fr/ZEN/>

[‡] <http://sanskrit.inria.fr/>

- Böhm, C. (1966). The CUCH as a formal and description language. In: Steel T. B. (ed.) *Formal Description Languages for Computer Programming*, North Holland 179–197.
- Brüggemann-Klein, A. (1993). Regular expressions into finite automata. *Theoretical Computer Science* **120** (2) 197–213.
- Brzozowski, J. A. (1964). Derivatives of regular expressions. *Journal of the ACM* **11** (4) 481–494.
- Champarnaud, J-M., Nicart, F. and Ziadi, D. (2006). From the ZPC structure of a regular expression to its follow automaton. *International Journal of Algebra and Computation* **16** (1) 17–34.
- Champarnaud, J-M. and Ziadi, D. (2001). From c-continuations to new quadratic algorithms for automaton synthesis. *International Journal of Algebra and Computation* **11** (6) 707–736.
- Champarnaud, J-M. and Ziadi, D. (2002). Canonical derivatives, partial derivatives and finite automaton constructions. *Theoretical Computer Science* **289** (1) 137–163.
- Eilenberg, S. (1974). *Automata, Languages, and Machines*, vol. A, Academic Press.
- Fischer, S., Huch, F. and Wilke, T. (2010). A play on regular expressions: Functional pearl. In: *Proceedings of the 15th ACM SIGPLAN International Conference on Functional Programming, ICFP'10*, New York, USA 357–368.
- Huet, G. (1980). Confluent reductions: Abstract properties and applications to term rewriting systems. *Journal of the ACM* **27** (4) 797–821.
- Huet, G. (2002). The Zen computational linguistics toolkit: Lexicon structures and morphology computations using a modular functional programming language. In: *Tutorial, Language Engineering Conference LEC'2002*.
- Huet, G. (2005). A functional toolkit for morphological and phonological processing, application to a Sanskrit tagger. *Journal of Functional Programming* **15** (4) 573–614.
- Huet, G. and Razet, B. (2006). The reactive engine for modular transducers. In: Futatsugi, K., Jouannaud, J.-P. and Meseguer, J. (eds.) *Algebra, Meaning and Computation, Essays Dedicated to Joseph A. Goguen*. Springer-Verlag Lecture Notes in Computer Science vol. 4060 355–374.
- Ilie, L. and Yu, S. (2003). Follow automata. *Information and Computation* **186** (1) 140–162.
- Kozen, D. (1994a). A completeness theorem for Kleene algebras and the algebra of regular events. *Information and Computation* **110** (2) 366–390.
- Kozen, D. (1994b). On action algebras. In: van Eijck, J. and Visser, A. (eds.) *Logic and Information Flow*, MIT Press 78–88.
- Landin, P. (1966). The next 700 programming languages. *CACM* **9** (3) 157–166.
- Pratt, V. (1991). Action logic and pure induction. In: *Proceedings of the European Workshop on Logics in AI (JELIA '90)*, Springer-Verlag 97–120.
- Razet, B. (2008). Finite Eilenberg machines. In: Ibarra, O. and Ravikumar, B. (eds.) *Proceedings of the CIAA 2008*. Springer-Verlag *Lecture Notes in Computer Science* **5148** 242–251.
- Razet, B. (2009). *Machines d'Eilenberg Effectives*, Ph.D. thesis, University Denis Diderot (Paris 7).
- Razet, B. (2011). Simulating finite Eilenberg machines with a reactive engine. In: *proceedings second workshop on mathematically structured functional programming (MSFP 2008)*. *Electronic Notes in Theoretical Computer Science* **229** (5) 119–134.
- Thompson, K. (1968). Programming techniques: Regular expression search algorithm. *CACM* **11** (6) 419–422.