



**HAL**  
open science

# A Classification of Viruses through Recursion Theorems

Guillaume Bonfante, Matthieu Kaczmarek, Jean-Yves Marion

► **To cite this version:**

Guillaume Bonfante, Matthieu Kaczmarek, Jean-Yves Marion. A Classification of Viruses through Recursion Theorems. *Computability in Europe*, Jun 2007, Sienna, Italy. pp.73-82, 10.1007/978-3-540-73001-9\_8. inria-00175301

**HAL Id: inria-00175301**

**<https://inria.hal.science/inria-00175301>**

Submitted on 27 Sep 2007

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# A Classification of Viruses through Recursion Theorems

Guillaume Bonfante, Matthieu Kaczmarek and Jean-Yves Marion

Guillaume.Bonfante@loria.fr, Matthieu.Kaczmarek@loria.fr and  
Jean-Yves.Marion@loria.fr

Loria - INPL - Ecole Nationale Supérieure des Mines de Nancy  
B.P. 239, 54506 Vandœuvre-lès-Nancy Cédex, France

**Abstract.** We study computer virology from an abstract point of view. Viruses and worms are self-replicating programs, whose definitions are based on Kleene's second recursion theorem. We introduce a notion of delayed recursion that we apply to both Kleene's second recursion theorem and Smullyan's double recursion theorem. This leads us to define four classes of viruses, two of them being polymorphic. Then, we work on a simple imperative programming language in order to show how those theoretical constructions can be implemented. In particular, we propose a general virus builder, and distribution engines.

**Topics covered.** Computability theoretic aspects of programs, computer virology.

**Keywords.** Computer viruses, polymorphism, propagation, recursion theorem, iteration theorem.

## 1 Theoretical Computer Virology

An important information security breach is computer virus infections. Following Filiol's book [10], we do think that theoretical studies should help to design new defenses against computer viruses. The objective of this paper is to pursue a theoretical study of computer viruses initiated in [4]. It is worth to cite von Neumann [26].

Can an automaton be constructed, i.e., assembled and built from appropriately "raw material", by an other automaton? [...] Can the construction of automata by automata progress from simpler types to increasingly complicated types?

Since viruses are essentially self-replicating programs, we see that virus programming methods are an attempt to answer to von Neumann's question.

A long term practical motivation of this work is to suggest new research directions in order to improve anti-virus techniques. Of course, we have to keep in mind that Cohen demonstrated that finding a virus is undecidable in the general case. However, as the problem of defense still remains, various detection techniques have been developed. They are based on the combination of signature pattern matching and behavioural detection methods, see [11]. We think that a

first step to reach an immunization technique is to understand the key concepts of virus writing. According to Thompson [25] and Ludwig [17], this is not trivial.

Abstract computer virology was initiated in the 80's by the seminal works of Cohen and Adleman [7]. The latter coined the term *virus*. Cohen defined viruses with respect to Turing Machines [8]. Later [1], Adleman took a more abstract point of view in order to have a definition independent from any particular computational model. Then, only a few theoretical studies followed those seminal works. Chess and White refined the mutation model of Cohen in [6]. Zuo and Zhou formalized polymorphism from Adleman's work [27] and they analyzed the time complexity of viruses [28].

From [23, 19], it is commonly accepted that a programming language enjoys self-reference as long as it satisfies the recursion theorem, self-reference being the central issue in computer virology. Moreover, implementations using the recursion theorem are efficient as was shown by Jones [15] (up to a linear constant). Suppose that a programming environment provides a way to get its code to the currently executed program. For example, in `bash` this is given by the variable `$0`. This feature does not offer more efficient self-referential codes. Our constructions show moreover that viruses are not easier (at least theoretically) to write with than without built-in self-reference. That is why we focus on Kleene's recursion theorem, and consequently, we have chosen the core language `WHILE+` which fulfils the theorem.

Despite the works [12, 13], the recursion theorem is used essentially to prove "negative" results such as the constructions of undecidable or inseparable sets, see [22] for a general reference, or such as Blum's speed-up theorem [2]. Here, we show that the recursion theorem plays a key role in the *construction* of viruses.

Recently, we tried [3, 4] to formalize inside computability the notion of viruses. This formalization captures previous definitions that we have mentioned above. We also have characterized two kinds of viruses, blueprint and smith viruses, and we proved constructively their existence. The present work proposes to go further, introducing a notion of distribution to model polymorphism or metamorphism and a notion of virus distribution builder, a virus generator. Furthermore, we switch to a simple programming language named `WHILE+`, introduced by Jones, to illustrate the effectivity of our constructions.

The concept of distribution is independent from any mutation technique and consequently a good model to study polymorphism or metamorphism. We have defined four kinds of viruses, the blueprint, the smith and their polymorphic variants. We show that each category is closely linked to a corresponding form of the recursion theorem. Kleene's recursion theorem corresponds to blueprint viruses, a model of duplication processes. A variation of Kleene's theorem introduces a code of fixed point generators. This result provides a construction of polymorphic viruses based on a blueprint duplication. A second variation presents the double recursion theorem, which defines the class of smith viruses. This class of viruses models self-propagating infections. Lastly, the third variation is a delayed double recursion theorem, which leads to polymorphic smith viruses.

All our existential results are constructive. We give for each theoretical construction a code in  $\text{WHILE}^+$ . Actually, we follow the ideas of the experimentation of the iteration theorem and of the recursion theorem, which are developed in [12, 13] by Jones et al. and very recently by Moss in [18].

## 2 A Virus Definition

### 2.1 The $\text{WHILE}^+$ language

The domain of computation  $\mathbb{D}$  is the set of binary trees generated from an atom  $\text{nil}$  and a (computable) pairing mechanism  $\langle \cdot, \cdot \rangle$ . We suppose we are given a countable set  $\mathbb{V}$  of variables and an other (countable) set of names of programs, say  $\mathbb{N}$ . The syntax of  $\text{WHILE}^+$  is given by the following grammar.

Expressions:  $\mathbb{E} \rightarrow \mathbb{D} \mid \mathbb{V} \mid \text{cons}(\mathbb{E}_1, \mathbb{E}_2) \mid \text{hd}(\mathbb{E}) \mid \text{tl}(\mathbb{E}) \mid$   
 $\text{exec}(\mathbb{E}_0, \mathbb{E}_1, \dots, \mathbb{E}_n) \mid \text{spec}_n(\mathbb{E}_0, \mathbb{E}_1 \dots, \mathbb{E}_n)$  with  $n \geq 1$

Commands:  $\mathbb{C} \rightarrow \mathbb{V} := \mathbb{E} \mid \mathbb{C}_1; \mathbb{C}_2 \mid \text{while}(\mathbb{E})\{\mathbb{C}\} \mid \text{if}(\mathbb{E})\{\mathbb{C}_1\}\text{else}\{\mathbb{C}_2\}$

Programs:  $\mathbb{P} \rightarrow \mathbb{N}(\mathbb{V}_1, \dots, \mathbb{V}_n)\{\mathbb{C}; \text{return } \mathbb{E};\}$

Elements of  $\mathbb{P}$  are  $\text{WHILE}^+$  programs. We suppose that we are given a concrete syntax of  $\text{WHILE}^+$ , that is an embedding of  $\mathbb{P}$  into  $\mathbb{D}$ , see Appendix A. So, from now on, when the context is clear, we do not make any distinction between a program and its concrete syntax. Moreover, as a shorthand, to denote a program, we will use its name.

The semantics of  $\text{WHILE}^+$  is the usual one and we postpone it to the Appendix B.  $\llbracket \mathbf{p} \rrbracket$  denotes the function computed by the program  $\mathbf{p}$ . For convenience, we have chosen to have a built-in interpreter  $\text{exec}(e, e_1, \dots, e_n)$  which executes the evaluation of the expression  $e$  on the evaluation of  $e_1, \dots, e_n$ . We also use a built-in specializer. Let us introduce the programs:

```

specn ( $x_0, \dots, x_n$ ) {
  r := spec( $x_0, \dots, x_n$ );
  return r;
}

```

We have  $\llbracket \llbracket \text{spec}_n \rrbracket(\mathbf{p}, x_1 \dots, x_n) \rrbracket(y) = \llbracket \mathbf{p} \rrbracket(x_1, \dots, x_n, y)$ .

The use of an interpreter and of a specializer is justified by Jones who showed in [15] that programs with these constructions can be simulated up to a linear constant time by programs without them.

If  $f$  and  $g$  designate the same function, we write  $f \approx g$ . A function  $f$  is *semi-computable* if there is a program  $\mathbf{p}$  such that  $\llbracket \mathbf{p} \rrbracket \approx f$ , moreover, if  $f$  is total, we say that  $f$  is *computable*.

## 2.2 What is a Computer Virus?

We do think that Kleene’s recursion theorem is the mathematical backbone of computer virology. Several authors made the same observation [21, 19, 1, 10], but as far as we know, none have developed this idea into a framework.

To model viruses, we propose the following scenario. When a program  $\mathbf{p}$  is executed within an environment  $x$ , we compute  $\llbracket \mathbf{p} \rrbracket(x)$  and if the evaluation halts, we replace  $x$  by the result of the computation. The entry  $x$  is thought of as a finite sequence  $\langle x_1, \dots, x_n \rangle$  which represents the file system and all accessible parameters.

Let us provide some examples to fix the intuition. Typically, a program **copy** which copies a file satisfies  $\llbracket \mathbf{copy} \rrbracket(\mathbf{p}, x) = \langle \mathbf{p}, \mathbf{p}, x \rangle$ . The original environment is  $\langle \mathbf{p}, x \rangle$ . After the evaluation of **copy**, we have the environment  $\langle \mathbf{p}, \mathbf{p}, x \rangle$ . The particular element  $\mathbf{p}$  has been copied.

Our second example models a *parasitic virus* [10]. Parasitic viruses insert themselves into existing files. When an infected host is executed, first the virus infects a new host, then it gives the control back to the original host. We write  $x; y$  to denote the composition of  $x$  and  $y$  and  $\mathbf{id}$  is a code of the identity. Let us define  $B(y, \mathbf{p}) = \llbracket \mathbf{spec}_1 \rrbracket(y, \mathbf{id}); \mathbf{p}$ . The program  $B(y, \mathbf{p})$  is the “infected form” of  $\mathbf{p}$  by  $y$ . Suppose that a program  $\mathbf{v}$  satisfies:

$$\llbracket \mathbf{v} \rrbracket(\mathbf{p}, \mathbf{q}, x) = \llbracket \mathbf{p} \rrbracket(B(\mathbf{v}, \mathbf{q}), x) \quad (1)$$

The following equations hold:

$$\llbracket B(\mathbf{v}, \mathbf{p}) \rrbracket(\mathbf{q}, x) = \llbracket \mathbf{p} \rrbracket(B(\mathbf{v}, \mathbf{q}), x) \quad \text{by specialization} \quad (2)$$

$$= \llbracket \mathbf{v} \rrbracket(\mathbf{p}, \mathbf{q}, x) \quad \text{by 1.} \quad (3)$$

Equation (2) shows the propagation process of the scenario. Equation (3) expresses the fact that an infected form can be seen as a specialized form of the virus  $\mathbf{v}$  for the host  $\mathbf{p}$ .

Generally speaking, the construction of viruses lies in the resolution of equations such as (1, 3) above with  $\mathbf{v}$  and  $B$  as unknowns. The latter equation links the virus to its infected forms while the former specifies the way the infection is propagated. The infection, expressed here by Equation (2), is then a consequence of the two equations (1, 3). We will show how this resolution can be done by use of Kleene’s recursion theorem and specialization.

With the previous discussion in mind, we present a formalization of viruses. A virus is a program which propagates its own code with possible mutations.

**Definition 1 (Computer Virus).** *Let  $B$  be a computable function. A virus w.r.t  $B$  is a program  $\mathbf{v}$  such that  $\forall \mathbf{p}, x : \llbracket \mathbf{v} \rrbracket(\mathbf{p}, x) = \llbracket B(\mathbf{v}, \mathbf{p}) \rrbracket(x)$ . Then,  $B$  is named a propagation function for the virus  $\mathbf{v}$ .*

Note that any virus  $\mathbf{v}$  accepts infinitely many different propagation functions. As a result, the fact that  $\mathbf{v}$  is a virus w.r.t some propagation function  $B$  does not imply that  $\mathbf{v}$  must infect through  $B$ ; it can use another propagation procedure.

Definition 1 includes the ones of Adleman and Cohen, and it handles more propagation and duplication features than the other models [4].

**Definition 2 (Virus Distribution).** A virus distribution is a pair  $(\mathbf{d}_V, \mathbf{d}_B)$  of programs such that for all  $i$ ,  $\llbracket \mathbf{d}_V \rrbracket(i)$  is a virus w.r.t  $\llbracket \llbracket \mathbf{d}_B \rrbracket(i) \rrbracket$ .  $\mathbf{d}_V$  is named a distribution engine and  $\mathbf{d}_B$  is named a propagation engine.

We refer to [24, 9] for practical details about virus distributions. As we will see later, a virus distribution can be used as a model of polymorphism. Instead of virus distributions, we speak about *virus builders* when  $i$  is intended to be a code of the behavior of the virus.

**Definition 3 (Distribution builder).** A Distribution builder is a pair of programs  $\mathbf{c}_V, \mathbf{c}_B$  such that for all  $\mathbf{h}$ ,  $(\llbracket \mathbf{c}_V \rrbracket(\mathbf{h}), \llbracket \mathbf{c}_B \rrbracket(\mathbf{h}))$  is a virus distribution.

### 3 Blueprint Duplication

#### 3.1 Blueprint Virus Builder

From [4], a *blueprint virus for a function  $g$*  is a program  $\mathbf{v}$  which computes  $g$  using its own code  $\mathbf{v}$  and its environment  $\mathbf{p}, x$ . The function  $g$  can be seen as the behavior of the virus. In other words, it is a program which satisfies

$$\begin{cases} \mathbf{v} \text{ is a virus w.r.t some propagation function} \\ \forall \mathbf{p}, x : \llbracket \mathbf{v} \rrbracket(\mathbf{p}, x) = g(\mathbf{v}, \mathbf{p}, x) \end{cases} \quad (4)$$

Note that a blueprint virus does not use any code of its propagation function. Clearly, the solutions of this system are provided by Kleene's recursion theorem.

**Theorem 4 (Kleene's Recursion Theorem [16]).** Let  $f$  be a semi-computable function. There is a program  $\mathbf{e}$  such that  $\llbracket \mathbf{e} \rrbracket(x) = f(\mathbf{e}, x)$ .

**Corollary 5.** There is a virus builder  $(\mathbf{d}_V, \mathbf{d}_B)$  such that for any program  $\mathbf{g}$ ,  $\llbracket \mathbf{d}_V \rrbracket(\mathbf{g})$  is a blueprint virus for  $\llbracket \mathbf{g} \rrbracket$ .

*Proof.* We use a construction for the recursion theorem due to Smullyan [23]. It provides a fixpoint which can be directly used as a virus builder. Let  $\mathbf{dg}$  and  $\mathbf{d}_V$  be the respective programs.

$$\begin{array}{lll} \mathbf{dg} (z,y,x)\{ & \mathbf{d}_V (\mathbf{g})\{ & \mathbf{cspec} (x)\{ \\ \quad r := \mathbf{exec}(z,\mathbf{spec}(y,z,y),x); & \quad r := \mathbf{spec}(\mathbf{dg},\mathbf{g},\mathbf{dg}); & \quad r := \mathbf{spec}_1; \\ \quad \mathbf{return} r; & \quad \mathbf{return} r; & \quad \mathbf{return} r; \\ \} & \} & \} \end{array}$$

We observe that  $\llbracket \llbracket \mathbf{d}_V \rrbracket(\mathbf{g}) \rrbracket(\mathbf{p}, x) = g(\llbracket \mathbf{d}_V \rrbracket(\mathbf{g}), \mathbf{p}, x)$ . Moreover,  $\llbracket \mathbf{d}_V \rrbracket(\mathbf{g})$  is clearly a virus w.r.t  $\llbracket \mathbf{spec}_1 \rrbracket$ . Then, we define  $\mathbf{d}_B = \mathbf{cspec}$ .  $\square$

We consider a typical example of blueprint duplication which looks like the real life virus `ILoveYou`. This program arrives as an e-mail attachment. Opening the attachment triggers the attack, the infection first scans the memory for passwords and sends them back to the attacker, then the virus self-duplicates

sending itself at every address of the local address book. Here, the propagation depends on an external actor who opens the attachment.

To give a model of this scenario we need to deal with mailing processes. A mail  $m = \langle @, y \rangle$  is an association of an address  $@$  and data  $y$ . Then, we consider that the environment contains a mailbox  $mb = \langle m_1, \dots, m_n \rangle$  which is a sequence of mails. To send a mail  $m$ , we add it to the mailbox, that is to assign  $mb := \text{cons}(m, mb)$ . We suppose that an external process deals with mailing.

In the following,  $x$  denotes the local file structure, and  $@bk = \langle @_1, \dots, @_n \rangle$  denotes the local address book, a sequence of addresses. We finally introduce a WHILE<sup>+</sup> program **find** which searches its input for passwords and which returns them as its evaluation. The virus behavior for the scenario of ILoveYou is given by the following program.

```

g (v,mb,@bk,x) {
  pass := exec(find,x);
  mb := cons(cons("badguy@dom.com",pass),mb);
  y := @bk;
  while (y) {
    mb := cons(cons(hd(y),v),mb);
    y := tl(y);
  }
  return cons(mb,cons(@bk,x));
}

```

The corresponding blueprint virus is yield by the virus builder of Corollary 5.

### 3.2 Blueprint Distributions

A *blueprint distribution* is a virus distribution which achieves blueprint duplication. The viruses of such a distribution can use a code of the distribution engine. As a result, any member can generate the entire distribution. This idea conveys an evolving ability. More formally, a blueprint distribution  $(\mathbf{d}_v, \mathbf{d}_B)$  for a function  $g$  is defined by the following system.

$$\begin{cases} (\mathbf{d}_v, \mathbf{d}_B) \text{ is a virus distribution} \\ \forall i, \mathbf{p}, x : \llbracket \mathbf{d}_v \rrbracket(i)(\mathbf{p}, x) = g(\mathbf{d}_v, i, \mathbf{p}, x) \end{cases} \quad (5)$$

The function  $g$  is thought of as the behavior of the whole distribution. The existence of blueprint distributions corresponds to a stronger form of the recursion theorem due to Case [5].

**Theorem 6 (Delayed Recursion [5]).** *Let  $f$  be a semi-computable function. There exists a computable function  $e$  such that  $\forall x, y : \llbracket e(x) \rrbracket(y) = f(\mathbf{e}, x, y)$  where  $\mathbf{e}$  computes  $e$ .*

**Corollary 7.** *There is a distribution builder  $(\mathbf{c}_v, \mathbf{c}_B)$  such that for any program  $\mathbf{g}$ ,  $(\llbracket \mathbf{c}_v \rrbracket(\mathbf{g}), \llbracket \mathbf{c}_B \rrbracket(\mathbf{g}))$  is a blueprint distribution for  $\llbracket \mathbf{g} \rrbracket$ .*

*Proof.* We use a construction close to the proof of Corollary 5. It provides a fixpoint which can be directly used as a distribution builder. We define:

<pre> <b>edg</b> (z,t,x,y) {   e := <b>spec</b>(<b>spec</b><sub>3,t,z,t</sub>);   <b>return exec</b>(z,e,x,y); } </pre>	<pre> <b>c<sub>v</sub></b> (<b>g</b>){   r := <b>spec</b>(<b>spec</b><sub>3,edg,g,edg</sub>);   <b>return r</b>; } </pre>	<pre> <b>c<sub>B</sub></b> (<b>g</b>){   r := <b>cspec</b>;   <b>return r</b>; } </pre>
---	---	---

We observe that for any  $i$ ,  $\llbracket \llbracket \llbracket \mathbf{c}_v \rrbracket(\mathbf{g}) \rrbracket(i) \rrbracket(\mathbf{p}, x) = g(\llbracket \mathbf{c}_v \rrbracket(\mathbf{g}), i, \mathbf{p}, x)$ . Moreover,  $\llbracket \llbracket \mathbf{c}_v \rrbracket(\mathbf{g}) \rrbracket(i)$  is a virus w.r.t  $\llbracket \llbracket \mathbf{c}_B \rrbracket(\mathbf{g}) \rrbracket(i)$ .  $\square$

To illustrate Theorem 7, we come back to the scenario of the virus `ILoveYou`, and we add to it mutation abilities. We introduce a `WHILE+` program `poly` which is a polymorphic engine. This program takes a program `p` and a key  $i$ , and it rewrites `p` according to  $i$  and conserving the semantics of `p`. Formally, `poly` satisfies  $\llbracket \llbracket \mathbf{poly} \rrbracket(\mathbf{p}, i) \rrbracket$  is one-one on  $i$  and  $\llbracket \llbracket \mathbf{poly} \rrbracket(\mathbf{p}, i) \rrbracket \approx \llbracket \mathbf{p} \rrbracket$ . From a recursion theoretic point of view, `poly` is a padding function.

We build a virus which self-duplicates sending mutated forms of itself. With the notations of the Sect. 3.1, we consider a behavior described by the following `WHILE+` program.

```

g (dv,i,mb,@bk,x) {
  pass := exec(find,x);
  mb := cons(cons("badguy@dom.com",pass),mb);
  next_key := cons(nil,i)
  virus := exec(dv,next_key);
  mutation := exec(poly,virus,i);
  y := @bk;
  while (y) {
    mb := cons(cons(hd(y),mutation),mb);
    y := tl(y);
  }
  return cons(mb,cons(@bk,x));
}

```

We apply Corollary 7 to transform this program into a code of the corresponding distribution engine.

## 4 Smith Reproduction

Blueprint viruses do not use their propagation function. Now, we present a more evolved model which captures the notion of infection. Here, viruses include their propagation mechanism within their own code. As parasitic viruses do.



#### 4.1 Smith Viruses

We define a *smith virus* as two programs  $\mathbf{v}, \mathbf{B}$  which compute a behavior  $g$  according to the following system.

$$\begin{cases} \mathbf{v} \text{ is a virus w.r.t } \llbracket \mathbf{B} \rrbracket \\ \forall \mathbf{p}, x : \llbracket \mathbf{v} \rrbracket(\mathbf{p}, x) = g(\mathbf{B}, \mathbf{v}, \mathbf{p}, x) \end{cases} \quad (6)$$

This class corresponds to the double recursion theorem due to Smullyan [20].

**Theorem 8 (Double Recursion Theorem [20]).** *Let  $f_1$  and  $f_2$  be two semi-computable functions. There are two programs  $\mathbf{e}_1$  and  $\mathbf{e}_2$  such that*

$$\llbracket \mathbf{e}_1 \rrbracket(x) = f_1(\mathbf{e}_1, \mathbf{e}_2, x) \quad \llbracket \mathbf{e}_2 \rrbracket(x) = f_2(\mathbf{e}_1, \mathbf{e}_2, x) \quad (7)$$

**Corollary 9.** *There is a virus builder  $(\mathbf{d}_\mathbf{v}, \mathbf{d}_\mathbf{B})$  such that for any program  $\mathbf{g}$ ,  $\llbracket \mathbf{d}_\mathbf{v} \rrbracket(\mathbf{g}), \llbracket \mathbf{d}_\mathbf{B} \rrbracket(\mathbf{g})$  is a smith virus for  $\llbracket \mathbf{g} \rrbracket$ .*

*Proof.* We use a double fixpoint. We define the following programs.

$$\begin{array}{lll} \mathbf{dg1} (z1, z2, y1, y2, x) \{ & \mathbf{dg2} (z1, z2, y1, y2, x) \{ & \mathbf{pispec} (\mathbf{g}, \mathbf{B}, \mathbf{v}, y, \mathbf{p}) \{ \\ \quad \mathbf{e1} := \text{spec}(y1, z1, z2, y1, y2); & \quad \mathbf{e1} := \text{spec}(y1, z1, z2, y1, y2); & \quad \mathbf{r} := \text{spec}(\mathbf{g}, \mathbf{B}, \mathbf{v}, \mathbf{p}); \\ \quad \mathbf{e2} := \text{spec}(y2, z2, z2, y1, y2); & \quad \mathbf{e2} := \text{spec}(y2, z2, z2, y1, y2); & \quad \text{return } \mathbf{r}; \\ \quad \text{return } \text{exec}(z1, \mathbf{e1}, \mathbf{e2}, x); & \quad \text{return } \text{exec}(z2, \mathbf{e1}, \mathbf{e2}, x); & \} \\ \} & \} & \} \end{array}$$

Let  $\mathbf{d}_\mathbf{v}$  and  $\mathbf{d}_\mathbf{B}$  be the following programs.

$$\begin{array}{ll} \mathbf{d}_\mathbf{v} (\mathbf{g}) \{ & \mathbf{d}_\mathbf{B} (\mathbf{g}) \{ \\ \quad \mathbf{r} := \text{spec}(\mathbf{pispec}, \mathbf{g}); & \quad \mathbf{r} := \text{spec}(\mathbf{pispec}, \mathbf{g}); \\ \quad \text{return } \text{spec}(\mathbf{dg2}, \mathbf{r}, \mathbf{g}, \mathbf{dg1}, \mathbf{dg2}); & \quad \text{return } \text{spec}(\mathbf{dg1}, \mathbf{r}, \mathbf{g}, \mathbf{dg1}, \mathbf{dg2}); \\ \} & \} \end{array}$$

We observe that for any program  $\mathbf{g}$

$$\llbracket \llbracket \mathbf{d}_\mathbf{v} \rrbracket(\mathbf{g}) \rrbracket(\mathbf{p}, x) = \llbracket \llbracket \llbracket \mathbf{d}_\mathbf{B} \rrbracket(\mathbf{g}) \rrbracket(\llbracket \mathbf{d}_\mathbf{v} \rrbracket(\mathbf{g}), \mathbf{p}) \rrbracket(x) = g(\llbracket \mathbf{d}_\mathbf{B} \rrbracket(\mathbf{g}), \llbracket \mathbf{d}_\mathbf{v} \rrbracket(\mathbf{g}), \mathbf{p}, x) \quad \square$$

We present how to build the parasitic virus of Sect. 2. We define

$$\mathbf{g} (\mathbf{B}, \mathbf{v}, \mathbf{p}, \mathbf{q}, x) \{ \\ \quad \text{infected\_form} := \text{exec}(\mathbf{B}, \mathbf{v}, \mathbf{q}); \\ \quad \text{return } \text{exec}(\mathbf{p}, \text{infected\_form}, x); \\ \}$$

The behavior  $\mathbf{g}$  of the virus is the following. First, it infects a new host  $\mathbf{q}$  with the virus  $\mathbf{v}$  using the propagation procedure  $\mathbf{B}$ . Then, it executes the original host  $\mathbf{p}$ . This corresponds to the behavior of a parasitic virus. We obtain a smith virus using the builder of Corollary 9.

Concerning the construction of Sect. 2, the attentive reader would notice  $\llbracket \llbracket \mathbf{B} \rrbracket(\mathbf{v}, \mathbf{p}) \rrbracket \approx \llbracket \llbracket \text{spec}_1 \rrbracket(\mathbf{v}, \mathbf{id}); \mathbf{p} \rrbracket$  which justifies the definition of  $B$  in the example.

## 4.2 Smith Distributions

Smith distributions model viruses which are able to mutate their code and their propagation mechanism. A *smith distribution*  $(\mathbf{d}_V, \mathbf{d}_B)$  for the behavior  $g$  satisfies

$$\begin{cases} (\mathbf{d}_V, \mathbf{d}_B) \text{ is a virus distribution} \\ \forall i, \mathbf{p}, x : \llbracket \mathbf{d}_V \rrbracket(i)(\mathbf{p}, x) = g(\mathbf{d}_B, \mathbf{d}_V, i, \mathbf{p}, x) \end{cases} \quad (8)$$

This class of distributions corresponds to the following theorem.

**Theorem 10 (Delayed Double Recursion).** *Let  $f_1$  and  $f_2$  be two semi-computable functions. There are two computable functions  $e_1$  and  $e_2$  such that for all  $x$  and  $y$*

$$\begin{aligned} \llbracket e_1(x) \rrbracket(y) &= f_1(\mathbf{e}_1, \mathbf{e}_2, x, y) & \llbracket e_2(x) \rrbracket(y) &= f_2(\mathbf{e}_1, \mathbf{e}_2, x, y) \\ &\text{where } \mathbf{e}_1 \text{ and } \mathbf{e}_2 \text{ respectively compute } e_1 \text{ and } e_2 \end{aligned}$$

**Corollary 11.** *There is a distribution builder  $(\mathbf{c}_V, \mathbf{c}_B)$  such that for any program  $g$ ,  $(\llbracket \mathbf{c}_V \rrbracket(g), \llbracket \mathbf{c}_B \rrbracket(g))$  is a smith distribution for  $\llbracket g \rrbracket$ .*

*Proof.* We define the following programs.

<pre> <b>edg1</b> (z1,z2,t1,t2,x,y) {   e1 := spec(spec<sub>5</sub>,t1,z1,z2,t1,t2);   e2 := spec(spec<sub>5</sub>,t2,z1,z2,t1,t2);   return exec(z1,e1,e2,x,y); } <b>pispec'</b> (g,d<sub>B</sub>,d<sub>V</sub>,i,y,p) {   r := spec(g d<sub>V</sub>,d<sub>V</sub>,i,p);   return r; } </pre>	<pre> <b>edg2</b> (z1,z2,t1,t2,x,y) {   e1 := spec(spec<sub>5</sub>,t1,z1,z2,t1,t2);   e2 := spec(spec<sub>5</sub>,t2,z1,z2,t1,t2);   return exec(z2,e1,e2,x,y); } </pre>
--	---

Let  $\mathbf{c}_V$  and  $\mathbf{c}_B$  be the following programs.

<pre> <b>c<sub>V</sub></b> (g){   r := spec(<b>pispec'</b>,g)   return spec(spec<sub>5</sub>,<b>edg2</b>,r,g,<b>edg1</b>,<b>edg2</b>); } </pre>	<pre> <b>c<sub>B</sub></b> (g){   r := spec(<b>pispec'</b>,g)   return spec(spec<sub>5</sub>,<b>edg1</b>,r,g,<b>edg1</b>,<b>edg2</b>); } </pre>
---	---

We observe that for any program  $g$

$$\begin{aligned} \llbracket \llbracket \mathbf{c}_V \rrbracket(g) \rrbracket(i)(\mathbf{p}, x) &= \llbracket \llbracket \llbracket \mathbf{c}_B \rrbracket(g) \rrbracket(i) \rrbracket(\llbracket \mathbf{c}_V \rrbracket(g) \rrbracket(i), \mathbf{p})(x) \\ &= g(\llbracket \mathbf{c}_B \rrbracket(g), \llbracket \mathbf{c}_V \rrbracket(g), i, \mathbf{p}, x) \end{aligned} \quad \square$$

We enhance the virus of Sect. 4.1, adding some polymorphic abilities. Any virus of generation  $i$  infects a new host  $\mathbf{q}$  with a virus of generation  $\langle i, \mathbf{nil} \rangle$  using the propagation procedure of generation  $i$ . Then it gives the control back to the original host  $\mathbf{p}$ . This behavior is illustrated by the following program.

```

g (dB, dv, i, p, q, x) {
  B := exec(dB, i);
  v := exec(dv, cons(i, nil));
  mutation := exec(poly, v, i);
  infected_form := exec(B, mutation, q);
  return exec(p, infected_form, x);
}

```

Then, we obtain the smith distribution by the builder of Corollary 11.

## References

1. L. Adleman. An abstract theory of computer viruses. In *Advances in Cryptology – CRYPTO’88*, volume 403. Lecture Notes in Computer Science, 1988.
2. M. Blum. A machine-independent theory of the complexity of recursive functions. *Journal of the Association for Computing Machinery*, 14(2):322–336, 1967.
3. G. Bonfante, M. Kaczmarek, and J.-Y. Marion. Toward an abstract computer virology. In *ICTAC*, pages 579–593, 2005.
4. G. Bonfante, M. Kaczmarek, and J.-Y. Marion. On abstract computer virology from a recursion-theoretic perspective. *Journal in Computer Virology*, 1(3-4), 2006.
5. J. Case. Periodicity in generations of automata. *Theory of Computing Systems*, 8(1):15–32, 1974.
6. D. Chess and S. White. An undetectable computer virus. *Proceedings of the 2000 Virus Bulletin Conference (VB2000)*, 2000.
7. F. Cohen. *Computer Viruses*. PhD thesis, University of Southern California, January 1986.
8. F. Cohen. On the implications of computer viruses and methods of defense. *Computers and Security*, 7:167–184, 1988.
9. O. Drori, N. Pappo, and D. Yachan. New malware distribution methods threaten signature-based av. *Virus Bulletin*, pages 9–11, September 2005.
10. E. Filiol. *Computer Viruses: from Theory to Applications*. Springer-Verlag, 2005.
11. E. Filiol. Malware pattern scanning schemes secure against black-box analysis. *Journal of Computer Virology*, 2(1):35–50, 2006.
12. T. Hansen, T. Nikolajsen, J. Träff, and N. Jones. Experiments with implementations of two theoretical constructions. In *Lecture Notes in Computer Science*, volume 363, pages 119–133. Springer Verlag, 1989.
13. N. Jones. Computer implementation and applications of kleene’s S-m-n and recursive theorems. In Y. N. Moschovakis, editor, *Lecture Notes in Mathematics, Logic From Computer Science*, pages 243–263. Springer Verlag, 1991.
14. N. Jones. *Computability and Complexity: From a Programming Perspective*. MIT Press, Cambridge, MA, USA, 1997.
15. N. Jones. *Constant Time Factors Do Matter*. MIT Press, Cambridge, MA, USA, 1997.
16. S. Kleene. *Introduction to Metamathematics*. Van Nostrand, 1952.
17. M. Ludwig. *The Giant Black Book of Computer Viruses*. American Eagle Publications, 1998.
18. L. Moss. Recursion theorems and self-replication via text register machine programs. In *EATCS bulletin*, 2006.

19. P. Odifredi. *Classical Recursion Theory*. North-Holland, 1989.
20. H. Rogers. *Theory of Recursive Functions and Effective Computability*. McGraw Hill, New York, 1967.
21. M. Sipser. *Introduction to the Theory of Computation*. PWS Publishing Co. Boston, MA, USA, 1996.
22. R. Smullyan. *Recursion Theory for Metamathematics*. Oxford University Press, 1993.
23. R. Smullyan. *Diagonalization and Self-Reference*. Oxford University Press, 1994.
24. P. Ször. *The Art of Computer Virus Research and Defense*. Addison-Wesley Professional, 2005.
25. K. Thompson. Reflections on trusting trust. *Communications of the Association for Computing Machinery*, 27(8):761–763, 1984.
26. J. von Neumann. *Theory of Self-Reproducing Automata*. University of Illinois Press, Urbana, Illinois, 1966. edited and completed by A.W.Burks.
27. Z. Zuo and M. Zhou. Some further theoretical results about computer viruses. *The Computer Journal*, 47(6):627–633, 2004.
28. Z. Zuo, Q.-x. Zhu, and M.-t. Zhou. On the time complexity of computer viruses. *IEEE Transactions on information theory*, 51(8):2962–2966, August 2005.

## A A concrete syntax for WHILE<sup>+</sup>

Take  $\lambda x.x : \{\mathbf{nil}, \mathbf{quote}, \mathbf{v}, \mathbf{cons}, \mathbf{hd}, \mathbf{tl}, \mathbf{exec}, \mathbf{spec}, :=, ;, \mathbf{while}, \mathbf{if}\} \rightarrow \mathbb{D}$  one one. We suppose that we are given a one-one mapping from  $\mathbb{V}$  to  $\mathbb{D}$ . The concrete syntax of WHILE<sup>+</sup> is defined by the following extension of  $\underline{\ast}$ . For any  $d \in \mathbb{D}$ ,  $x \in \mathbb{V}$ ,  $e_1, e_2 \in \mathbb{E}$ ,  $c_1, c_2 \in \mathbb{E}$ ,

$$\begin{aligned}
 \underline{d} &= \langle \mathbf{quote}, d \rangle \\
 \underline{x} &= \langle \mathbf{v}, x \rangle \\
 \underline{\mathbf{cons}(e_1, e_2)} &= \langle \mathbf{cons}, \underline{e_1}, \underline{e_2} \rangle \\
 \underline{\mathbf{hd}(e_1)} &= \langle \mathbf{hd}, \underline{e_1} \rangle \\
 \underline{\mathbf{tl}(e_1)} &= \langle \mathbf{tl}, \underline{e_1} \rangle \\
 \underline{\mathbf{exec}(e_0, \dots, e_n)} &= \langle \mathbf{exec}, \underline{e_0}, \dots, \underline{e_n} \rangle \\
 \underline{\mathbf{spec}(e_0, \dots, e_n)} &= \langle \mathbf{spec}, \underline{e_0}, \dots, \underline{e_n} \rangle \\
 \underline{x := e_1} &= \langle :=, \underline{x}, \underline{e_1} \rangle \\
 \underline{c_1; c_2} &= \langle ;, \underline{c_1}, \underline{c_2} \rangle \\
 \underline{\mathbf{while}(e_1)\{c_1\}} &= \langle \mathbf{while}, \underline{e_1}, \underline{c_1} \rangle \\
 \underline{\mathbf{if}(e)\{c_1\}\mathbf{else}\{c_2\}} &= \langle \mathbf{if}, \underline{e}, \underline{c_1}, \underline{c_2} \rangle
 \end{aligned}$$

Let  $v_0$  be an element of  $\mathbb{V}$ . By convention for any element  $\mathbf{p} \in \mathbb{D}$  which is not a concrete syntax of any WHILE<sup>+</sup> program, we define  $\mathbf{p}$  as the concrete syntax of  $(v_0)\{v_0 := \mathbf{nil}; \mathbf{return nil}; \}$ .

## B Semantics of WHILE<sup>+</sup>

We note  $\pi_1$  and  $\pi_2$  the projections associated to  $\langle , \rangle$ . By convention  $\pi_1(\mathbf{nil}) = \pi_2(\mathbf{nil}) = \mathbf{nil}$ . Finite sequences are built by repeated applications of the pairing function:  $\langle x_1, \dots, x_n \rangle = \langle x_1, \langle x_2, \langle \dots, x_n \rangle \dots \rangle \rangle$ .

A store is a mapping  $\sigma : \mathbb{V} \rightarrow \mathbb{D}$ . The set of stores is noted  $\mathbb{S}$ . For any  $x \in \mathbb{V}$  the notation  $\sigma[x \mapsto d]$  denotes the function  $\sigma'$  such that  $\sigma'(x) = d$  and for any  $y \in \mathbb{V}$  such that  $y \neq x$ ,  $\sigma'(y) = \sigma(y)$ . By extension, for any  $x_1, \dots, x_k$ ,  $[x_1 \mapsto d_1, \dots, x_k \mapsto d_k]$  denotes the store  $\sigma$  such that  $\forall i \leq k : \sigma(x_i) = d_i$  and  $\forall i > k : \sigma(x_i) = \mathbf{nil}$ .

The semantics is defined by three functions,  $\mathcal{E} : \mathbb{E} \rightarrow (\mathbb{S} \rightarrow \mathbb{D})$  for expressions,  $\mathcal{C} : \mathbb{C} \rightarrow (\mathbb{S} \rightarrow \mathbb{S})$  for commands and  $\mathcal{P} : \mathbb{P} \rightarrow (\mathbb{D} \rightarrow \mathbb{D})$  for WHILE<sup>+</sup> programs.

*Evaluation of expressions.* For any  $d \in \mathbb{D}$ ,  $x \in \mathbb{V}$ ,  $e_1, e_2 \in \mathbb{E}$ ,  $\sigma \in \mathbb{S}$

$$\begin{aligned}
\mathcal{E}_d(\sigma) &= d \\
\mathcal{E}_x(\sigma) &= \sigma(x) \\
\mathcal{E}_{\text{cons}(e_1, e_2)}(\sigma) &= \langle \mathcal{E}_{e_1}(\sigma), \mathcal{E}_{e_2}(\sigma) \rangle \\
\mathcal{E}_{\text{hd}(e_1)}(\sigma) &= \pi_1(\mathcal{E}_{e_1}(\sigma)) \\
\mathcal{E}_{\text{tl}(e_1)}(\sigma) &= \pi_2(\mathcal{E}_{e_1}(\sigma)) \\
\mathcal{E}_{\text{exec}(e, e_1, \dots, e_n)}(\sigma) &= \llbracket \mathcal{E}_e(\sigma) \rrbracket (\langle \mathcal{E}_{e_1}(\sigma), \dots, \mathcal{E}_{e_n}(\sigma) \rangle)
\end{aligned}$$

*Evaluation of commands.* For any  $e \in \mathbb{E}$ ,  $c_1, c_2 \in \mathbb{C}$ ,  $\sigma \in \mathbb{S}$

$$\begin{aligned}
\mathcal{C}_{x := e}(\sigma) &= \sigma[x \mapsto \mathcal{E}_e(\sigma)] \\
\mathcal{C}_{c_1; c_2}(\sigma) &= \mathcal{C}_{c_2}(\mathcal{C}_{c_1}(\sigma)) \\
\mathcal{C}_{\text{while}(e)\{c_1\}}(\sigma) &= \begin{cases} \sigma & \text{if } \mathcal{E}_e(\sigma) = \text{nil} \\ \mathcal{C}_{\text{while}(e)\{c_1\}}(\mathcal{C}_c(\sigma)) & \text{otherwise} \end{cases} \\
\mathcal{C}_{\text{if}(e)\{c_1\}\text{else}\{c_2\}} &= \begin{cases} \mathcal{C}_{c_2}(\sigma) & \text{if } \mathcal{E}_e(\sigma) = \text{nil} \\ \mathcal{C}_{c_1}(\sigma) & \text{otherwise} \end{cases}
\end{aligned}$$

*Evaluation of programs.* For any  $d \in \mathbb{D}$ ,  $x_1, \dots, x_k \in \mathbb{V}$ ,  $e \in \mathbb{E}$ ,  $c \in \mathbb{C}$

$$\begin{aligned}
\mathcal{P}_{(x_1, \dots, x_k)\{c; \text{return } e; \}}(d) &= \mathcal{E}_e(\mathcal{C}_c([x_1 \mapsto d_1, \dots, x_k \mapsto d_k])) \\
&\text{where } d_1 = \pi_1(d), \forall 1 < i < k : d_k = \pi_1(\pi_1^{k-1}(d)) \text{ and } d_K = \pi_1^k(d)
\end{aligned}$$

We define  $\llbracket \mathbf{p} \rrbracket \approx \mathcal{P}_{\mathbf{w}}$  where  $\mathbf{p}$  is the concrete syntax of  $\mathbf{w}$ . The programming language  $\llbracket \cdot \rrbracket$  can be shown acceptable, see [14]. Then, there is a computable function  $S_n$  such that  $\llbracket S_n(\mathbf{p}, x_1, \dots, x_n) \rrbracket(y) = \llbracket \mathbf{p} \rrbracket(x_1, \dots, x_n, y)$ . We define  $\mathcal{E}_{\text{spec}(e, e_1, \dots, e_n)}(\sigma) = S_n(\mathcal{E}_e(\sigma), \mathcal{E}_{e_1}(\sigma), \dots, \mathcal{E}_{e_n}(\sigma))$  and we extend  $\mathcal{E}$ ,  $\mathcal{C}$ ,  $\mathcal{P}$ ,  $\llbracket \cdot \rrbracket$ , accordingly.