

Towards a Certified Petri Net Model-Checker

Lukasz Fronc and Franck Pommereau

IBISC, University of Évry, Tour Évry 2
523 place des terrasses de l'Agora, 91000 Évry, France
{fronc,pommereau}@ibisc.univ-evry.fr

Abstract. *Petri nets* are widely used in the domain of automated verification through *model-checking*. In this approach, a Petri Net model of the system of interest is produced and its reachable states are computed, searching for erroneous executions. *Model compilation* can accelerate this analysis by generating code to explore the reachable states. This avoids the use of a fixed exploration tool involving an “interpretation” of the Petri net structure. In this paper, we show how to compile Petri nets targeting the *LLVM language* (a high-level assembly language) and formally prove the *correctness* of the produced code. To this aim, we define a *structural operational semantics* for the fragment of LLVM we use.

Keywords: explicit model-checking, model compilation, LLVM, SOS.

1 Introduction

Verification through *model-checking* [1] consists in defining a formal model of the system to be analysed and then using automated tools to check whether the expected properties are met or not. We consider here more particularly the widely adopted setting in which models are expressed using *coloured Petri nets* [9] and there states are explored using *explicit model-checking* that enumerates them all (contrasting with symbolic model-checking that works with sets of states). *Model compilation* is one of the numerous techniques to speedup explicit model-checking, it relies on generating source code (then compiled into machine code) to produce a high-performance implementation of the state space exploration primitives. For instance, this approach is successfully used in the well known SPIN tool [7], or in Helena coloured Petri net model-checker [14,3].

In this paper, we propose a way to prove the correctness of such an approach. More precisely, we focus on the produced code and prove that the object computed by its execution is an actual representation of the state space of the compiled model. We consider the *Low-Level Virtual Machine (LLVM)* language as our target language for compilation, which reconciles two otherwise contradictory objectives: on the one hand, this is a typed language with reasonably high-level operations allowing to express algorithms quite naturally; on the other hand, it is a low-level language that can be equipped with a formal semantics allowing to formally prove the programs correctness. To do so, we define a structural operational semantics of the fragment of LLVM we need and use it to establish the properties of the programs generated by our compiler.

To the best of our knowledge, this is the first attempt to provide a formal semantics for LLVM. Moreover, if model-checkers are widely used tools, there exists surprisingly few attempts to prove them at the implementation level [19], contrasting with the domain of proof assistants [2,15] for which “proving the prover” is a common expectation.

The rest of the paper is organised as follows. We first recall the main notions about coloured Petri nets. Then, we present the LLVM framework, in particular the syntax of the language and its intuitive semantics, and how it can be embedded LLVM into a Petri net as a concrete colour domain. In section 4, we present algorithms and data structures for state space exploration. We then formally define an operational semantics for LLVM, including an explicit memory model. Finally we present our correctness results. Due to the limited number of pages, many definitions and intermediary results have been omitted, as well as the detailed proofs. This material can be found in [5,4]. Notice also that our compilation approach is evaluated from a performance point of view in [6].

2 Coloured Petri Nets

A (coloured) Petri net involves objects defined by a *colour domain* that provides data values, variables, operators, a syntax for expressions, possibly typing rules, etc. Usually, elaborated colour domains are used to ease modelling; in particular, one may consider a functional programming language [9,17] or the functional fragment (expressions) of an imperative programming language [14,16]. In this paper we will consider LLVM as a concrete colour domain.

All these can be seen as implementations of a more general *abstract colour domain* providing \mathbb{D} the set of *data values*, \mathbb{V} the set of *variables* and \mathbb{E} the set of *expressions*. Let $e \in \mathbb{E}$, we denote by $\text{vars}(e)$ the set of variables from \mathbb{V} involved in e . Moreover, variables or values may be considered as (simple) expressions, *i.e.*, we assume $\mathbb{D} \cup \mathbb{V} \subseteq \mathbb{E}$.

At this abstract level, we do not make any assumption about the typing or syntactical correctness of expressions; instead, we assume that any expression can be evaluated, possibly to $\perp \notin \mathbb{D}$ (undefined value) in case of any error. More precisely, a *binding* is a partial function $\beta : \mathbb{V} \rightarrow \mathbb{D} \cup \{\perp\}$. Then, let $e \in \mathbb{E}$ and β be a binding, we denote by $\beta(e)$ the evaluation of e under β ; if the domain of β does not include $\text{vars}(e)$ then $\beta(e) \stackrel{\text{def}}{=} \perp$. The application of a binding to evaluate an expression is naturally extended to sets and multisets of expressions.

Definition 1 (Petri nets). A Petri net is a tuple (S, T, ℓ) where S is the finite set of places, T , disjoint from S , is the finite set of transitions, and ℓ is a labelling function such that:

- for all $s \in S$, $\ell(s) \subseteq \mathbb{D}$ is the type of s , *i.e.*, the values that s may contain;
- for all $t \in T$, $\ell(t) \in \mathbb{E}$ is the guard of t , *i.e.*, a condition for its execution;
- for all $(x, y) \in (S \times T) \cup (T \times S)$, $\ell(x, y)$ is a multiset over \mathbb{E} and defines the arc from x toward y .

A marking of a Petri net is a map that associates to each place $s \in S$ a multiset of values from $\ell(s)$. From a marking M , a transition t can be fired using a binding β and yielding a new marking M' , which is denoted by $M[t, \beta]M'$, iff:

- there are enough tokens: for all $s \in S$, $M(s) \geq \beta(\ell(s, t))$;
- the guard is validated: $\beta(\ell(t)) = \text{true}$;
- place types are respected: for all $s \in S$, $\beta(\ell(t, s))$ is a multiset over $\ell(s)$;
- M' is M with tokens consumed and produced according to the arcs: for all $s \in S$, $M'(s) = M(s) - \beta(\ell(s, t)) + \beta(\ell(t, s))$.

Such a binding β is called a mode of t at marking M .

For a Petri net node $x \in S \cup T$, we define $\bullet x \stackrel{\text{df}}{=} \{y \in S \cup T \mid \ell(y, x) \neq \emptyset\}$ and $x^\bullet \stackrel{\text{df}}{=} \{y \in S \cup T \mid \ell(x, y) \neq \emptyset\}$ where \emptyset is the empty multiset. Finally, we extend the notation vars to a transition by taking the union of the variable sets in its guard and connected arcs.

In this paper, we assume that the considered Petri nets are such that, for all place $s \in S$ and all transition $t \in T$, $\ell(s, t)$ is either \emptyset or contains a single variable $x \in \mathbb{V}$. We also assume that $\text{vars}(t) = \bigcup_{s \in S} \text{vars}(\ell(s, t))$, i.e., all the variables involved in a transition can be bound using the input arcs. The second assumption is a classical one that allows to simplify the discovery of modes. The first assumption is made without loss of generality to simplify the presentation.

3 LLVM

The LLVM project (*Low Level Virtual Machine*) [11] is a modern and modular toolkit for compiler development used by a wide variety of commercial and open source projects as well as academic researches [13,12]. The LLVM-IR (*LLVM Intermediate Representation*) [10] is a part of the LLVM project and is a low-level, platform-independent, intermediate language. Every program written in this language can be run in a virtual machine or compiled to native code on all the platforms supported by the LLVM project. Importantly, the LLVM compiler runs a variety of optimisation passes on the LLVM-IR, which allows us to produce simple source code knowing it will be optimised by LLVM.

3.1 Syntax and Intuitive Semantics

A LLVM program is composed of a set of *blocks* (i.e., sequences of instructions) identified by *labels*. Entering or leaving a block is always explicit through branching instructions (jumps), subprograms calls or return instructions.

To define the syntax, we consider the following pairwise disjoint sets:

- \mathbb{P} is the set of *pointers*;
- \mathbb{T} is the set of *types*, defined inductively as the smallest set containing the *primitive types* in $\mathbb{T}_0 \stackrel{\text{df}}{=} \{\text{int}, \text{bool}, \dots\}$ (integers, Boolean values and other types defined by LLVM but not needed here) and such that if $t_0, \dots, t_n \in \mathbb{T}$ then $\text{struct}(t_0, \dots, t_n) \in \mathbb{T}$, which represents a data structure with $n + 1$ fields of types t_0 to t_n ;

- \mathbb{L} is the set of *labels*, it contains arbitrary names as well as some specific labels like $f(a_1, \dots, a_n)$, where $a_i \in \mathbb{V}$ for $1 \leq i \leq n$, that correspond to subprograms entry points (including the formal parameters). We define a set $\mathbb{L}_\perp \stackrel{\text{df}}{=} \mathbb{L} \cup \{\perp\}$ where $\perp \notin \mathbb{L}$ stands for an undefined label.

A program is represented as a partial function P from \mathbb{L} to the set of blocks, *i.e.*, that associates each label in its domain to a sequence of instructions.

For our purpose, we need to consider a fragment of LLVM that is formed by three main syntactic classes: sequences in *seq*, commands in *cmd* (*i.e.*, instructions) and expressions in *expr*. A sequence is a list of commands which may end with an expression, in which case it is considered as an expression itself (which is not reflected on the grammar in figure 1 to keep it simpler).

We assume that programs are syntactically correct and well typed, so that we can simplify the syntax by forgetting all types in LLVM source code. The resulting syntax is presented in figure 1. To allow for writing one-line sequences, we introduce the sequencing operator “;” that corresponds to the line endings. We also introduce the *skip* command that denotes the empty sequence. It may be noted that *pcall* (procedure call) and *fcall* (function call) do not exist in LLVM but are different instances of the *call* instruction. This distinction can be easily made in LLVM because the instruction contains the return type of the subprogram (function or procedure). Instruction *store* (resp. *load*) is the action of storing (resp. loading) data into (resp. from) the memory through a pointer. Instruction *icmp* compares two integers. Instruction *phi* is used to access variables assigned in previously executed blocks. Instruction *gep* corresponds to pointer arithmetic, we freeze the second argument to 0, which is enough to access fields in structures by their indexes.

3.2 LLVM-Labelled Petri Nets

To compile Petri nets as defined previously into LLVM programs, we need to consider a variant where the colour domain explicitly refers to a LLVM program.

Definition 2 (LLVM labelled Petri nets). *A LLVM labelled Petri net is a tuple $N \stackrel{\text{df}}{=} (S, T, \ell, P)$, where P is a LLVM program, and such that (S, T, ℓ) is a coloured Petri net with the following changes:*

- for all place $s \in S$, $\ell(s)$ is a LLVM type in \mathbb{T} , interpreted as a subset of \mathbb{D} ;
- for all transition $t \in T$, $\ell(t)$ is a call to a Boolean function in P whose parameters are the elements of $\text{vars}(t)$;
- for all $s \in t^\bullet$, $\ell(t, s)$ is a singleton multiset whose unique element is a call to a $\ell(s)$ -typed function in P whose parameters are the elements of $\text{vars}(t)$.

We assume that all the functions involved in the annotations terminate.

With respect to the previous definition, we have concretized the place types and each expression is now implemented as a LLVM function called from the corresponding annotation. To simplify the presentation, we have also restricted the output arcs to be singleton multisets, but this can be easily generalised.

$seq ::= cmd$	(statement)
$expr$	(expression)
$cmd; seq$	(sequence of instructions)
$cmd ::= br label$	(unconditional branching)
$br rvalue, label, label$	(conditional branching)
$pcall label(rvalue, \dots, rvalue)$	(procedure call)
ret	(return from a procedure)
$var = expr$	(variable assignment)
$store rvalue, rvalue$	(assignment through a pointer)
$skip$	(empty sequence)
$expr ::= add rvalue, rvalue$	(addition)
$load rvalue$	(read a value through a pointer)
$gep rvalue, 0, nat$	(get a pointer to a structure field)
$icmp op, rvalue, rvalue,$	(integers comparison)
$phi (rvalue, label), \dots, (rvalue, label)$	(get a value after branching)
$fcall label(rvalue, \dots, rvalue)$	(function call)
$alloc type$	(memory allocation)
$ret rvalue$	(return a value from a function)
$rvalue$	(variable or value)

Fig. 1. Our fragment of the LLVM syntax, where $label \in \mathbb{L}$, $rvalue \in \mathbb{D} \cup \mathbb{P} \cup \mathbb{V}$, $var \in \mathbb{V}$, $type \in \mathbb{T}$, $nat \in \mathbb{N}$ and $op \in \{<, \leq, =, \neq, \geq, >\}$

Moreover, the definitions of binding and modes are extended to LLVM. A *LLVM binding* is a partial function $\beta : \mathbb{V} \rightarrow \mathbb{D} \cup \mathbb{P}$ that maps each variable from its domain to a pointer or a value, and that is widened to \mathbb{D} by the identity function. \mathbb{B} is the set of all LLVM bindings. A *LLVM mode* is thus a LLVM binding enabling a transition in a LLVM labelled Petri net.

4 Implementing State Space Exploration

Given an *initial marking* M_0 , the state space we want to compute in this paper is the set R of *reachable marking*, *i.e.*, the smallest set such that $M_0 \in R$ and, if $M \in R$ and $M[t, \beta]M'$ then $M' \in R$ also. The correctness and termination of the implementation presented in this section are addressed in section 6.

Our algorithms are implemented on the basis of data structures (multisets, places, markings, and sets) that must respect some interfaces. An interface is presented as a set of procedures or functions that manipulates a data structure through a pointer (C-like interfaces). Moreover, each such subprogram has a formal specification of its behaviour that relies on an explicit interpretation of the data structure before and after the subprogram call. Precise examples are given in section 5.3, after the definition of LLVM formal semantics.

A multiset structure to store values from a type d is assumed and we call the set of values from d having non-zero occurrences in the multiset its *domain*. The multiset interface contains in particular: two procedures $add_{mset}(p_{mset}, elt)$ and $rem_{mset}(p_{mset}, elt)$ to respectively add or remove an element elt in p_{mset} ; a

function $size_{mset}(p_{mset})$ to return the domain size; a function $nth_{mset}(p_{mset}, n)$ to return the n^{th} element from the domain (for an arbitrary fixed order).

As a container of tokens, a place can be basically implemented as a multiset of tokens. So the place interface is exactly the multiset interface but annotated by the place name, for instance add_s is like add_{mset} but for place s .

The markings interface contains for each place s a function $get_s(p_{mrk})$ that returns a pointer to the corresponding place structure, as well as a function $copy_{mrk}(p_{mrk})$ that returns a copy of the marking structure.

Finally, the set interface contains a function $cons_{set}()$ that builds a new empty set and a procedure $add_{set}(p_{set}, elt)$ that adds an element elt to p_{set} .

Transitions Firing. Let $t \in T$ be a transition such that $\bullet t = \{s_1, \dots, s_n\}$ and $t^\bullet = \{s'_1, \dots, s'_m\}$. Then, function $fire_t$, that computes the marking M' reachable from M by firing t given a valuation of its variables, can be expressed as shown on the left of figure 2. This function simply creates a copy M' of M , removes from it the consumed tokens and adds the produced tokens before to return M' . One could remark that it avoids a loop over the Petri net places but instead it executes a sequence of statements. This is generally more efficient (no branching penalties, no loop overhead, no lookup of functions f_{t,s'_j}, \dots) and the code is simpler to generate. Let now x_{mrq} be a pointer to a marking structure implementing M . The firing algorithm can be implemented as shown on the right of figure 2.

Successors Computation. To discover all the possible modes for transition t , function $succ_t$ enumerates all the combinations of tokens from the input places. If a combination corresponds to a mode then the suitable transition firing function is called to produce a new marking. This algorithm is shown in figure 3. Note the nesting of loops that avoids an iteration over $\bullet t$, which saves from querying the Petri net structure and avoids the explicit construction of a binding. Moreover, since g_t is written in the target language, we avoid an interpretation of the

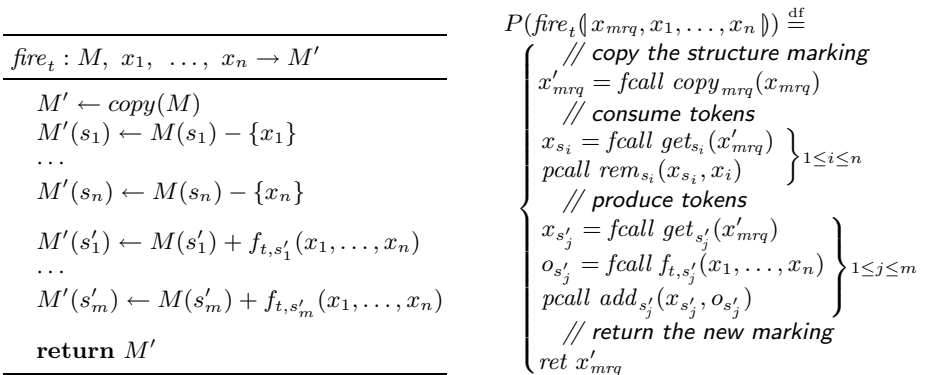


Fig. 2. On the left, the transition firing algorithm, where x_i is the variable in $\ell(s_i, t)$ for all $1 \leq i \leq n$, and f_{t,s'_j} is the function called in $\ell(t, s'_j)$ for all $1 \leq j \leq m$. On the right, its LLVM implementation.

corresponding expression. For the LLVM version, let x_{mrq} be a pointer to a marking structure and x_{next} be a pointer to a marking set structure. Then, the algorithm from figure 3 can be implemented as shown in figure 4. Each iteration over x_k is implemented as a set of blocks subscribed by t, k (for $n \geq k \geq 1$); blocks subscribed by $t, 0$ corresponds to the code inside the innermost loop. Note the *phi* instruction to update the value of index i_{s_i} (used to enumerate the tokens in place s_i): when the program enters block $loop_{t,i}$ for the first time, it comes from block $header_{t,i}$, so we initialise the value of i_{s_i} to the last index in the domain of s_i ; later, the program comes back to block $loop_{t,i}$ from block $footer_{t,i}$, so it assigns i'_{s_i} to i_{s_i} that is the value of $i_{s_i} - 1$ (i.e., the previous index).

A function *succ* is also defined to compute the set of all the successors of a marking, which is made by calling all the transition specific successor functions and accumulating the discovered markings into the same set. This algorithm and its translation in LLVM are shown in figure 5.

5 A Formal Semantics of LLVM

5.1 Memory Model

To start with, we define a memory model for LLVM, including *heaps* to store dynamically allocated pointers as well as *stacks* to store local variables and arguments for subprograms calls.

A *heap* is a partial function $H : \mathbb{P} \rightarrow \mathbb{T} \times (\mathbb{D} \cup \mathbb{P} \cup \{\perp\})^*$ with a finite domain. Each heap maps every pointer in its domain to a type and a tuple of values or pointers. The set of all heaps is \mathbb{H} . A heap is well formed if every pair in its image is type-consistent, for instance if $H(p) = (int, d)$ then d is indeed an integer or is \perp (uninitialised).

The set of all the pointers accessible starting from a pointer p in a heap H is denoted by $p \downarrow_H$ and is defined for all p in \mathbb{P} as:

$$\begin{aligned} p \downarrow_H &\stackrel{\text{df}}{=} \{ \} && \text{if } p \notin \text{dom}(H) \\ p \downarrow_H &\stackrel{\text{df}}{=} \{ p \} && \text{if } H(p) = (t, v) \text{ and } t \in \mathbb{T}_0 \\ p \downarrow_H &\stackrel{\text{df}}{=} \{ p \} \cup p_0 \downarrow_H \cup \dots \cup p_n \downarrow_H && \text{if } H(p) = (\text{struct}(t_0, \dots, t_n), (p_0, \dots, p_n)) \end{aligned}$$

It can be shown that if a heap H is well formed then $p \downarrow_H \subseteq \text{dom}(H)$ for every $p \in \text{dom}(H)$, and more generally that $\text{dom}(H) = \bigcup_{p \in \mathbb{P}} p \downarrow_H$.

Then, we need to access and update the data stored onto a heap. For each heap H , we define a data structure traversal function $\cdot[\cdot]_H : \mathbb{P} \times \mathbb{N} \rightarrow \mathbb{P} \cup \mathbb{D}$ as:

$$p[i]_H \stackrel{\text{df}}{=} \begin{cases} p_i & \text{if } H(p) = (\text{struct}(t_0, \dots, t_n), (p_0, \dots, p_n)) \text{ and } 0 \leq i \leq n \\ \text{undefined} & \text{otherwise} \end{cases}$$

The overwriting function $\oplus : \mathbb{H} \times \mathbb{H} \rightarrow \mathbb{H}$ represents the writing into memory and is defined for each $p \in \mathbb{P}$ as:

$$(H \oplus H')(p) \stackrel{\text{df}}{=} \begin{cases} H'(p) & \text{if } p \in \text{dom}(H') \\ H(p) & \text{if } p \notin \text{dom}(H') \wedge p \in \text{dom}(H) \\ \text{undefined} & \text{otherwise} \end{cases}$$

 $succ_t : M, next \rightarrow \perp$

```

for  $x_n$  in  $M(s_n)$  do
  ...
  for  $x_1$  in  $M(s_1)$  do
    if  $g_t(x_1, \dots, x_n)$  then
       $next \leftarrow next \cup \{fire_t(M, x_1, \dots, x_n)\}$ 
    endif
  endfor
  ...
endfor

```

Fig. 3. Transition specific successors computation algorithm, where g_t is the function that evaluates the guard $\ell(t)$

$$\begin{aligned}
 P(succ_t(x_{mrq}, x_{next})) &\stackrel{\text{df}}{=} \{ br\ header_{t,n} \\
 P(header_{t,k}) &\stackrel{\text{df}}{=} \begin{cases} x_{s_k} = fcall\ get_{s_k}(x_{mrq}) \\ s_{s_k} = fcall\ size_{s_k}(x_{s_k}) \\ br\ loop_{t,k} \end{cases} \\
 P(loop_{t,k}) &\stackrel{\text{df}}{=} \begin{cases} i_{s_k} = phi(s_{s_k}, header_{t,k}), (i'_{s_k}, footer_{t,k}) \\ c_{s_k} = icmp\ >, i_{s_k}, 0 \\ br\ c_{s_k}, body_{t,k}, footer_{t,k+1} \end{cases} \\
 P(body_{t,k}) &\stackrel{\text{df}}{=} \begin{cases} x_k = fcall\ nth_{s_k}(x_{s_k}, i_{s_k}) \\ br\ header_{t,k-1} \end{cases} \\
 P(footer_{t,k}) &\stackrel{\text{df}}{=} \begin{cases} i'_{s_k} = add\ i_{s_k}, -1 \\ br\ loop_{t,k} \end{cases} \\
 P(header_{t,0}) &\stackrel{\text{df}}{=} \begin{cases} c_g = fcall\ g_t(x_1, \dots, x_n) \\ br\ c_g, body_{t,0}, footer_{t,1} \end{cases} \\
 P(body_{t,0}) &\stackrel{\text{df}}{=} \begin{cases} x'_{mrq} = fcall\ fire_t(x_{mrq}, x_1, \dots, x_n) \\ pcall\ add_{set}(x_{next}, x'_{mrq}) \\ br\ footer_{t,1} \end{cases} \\
 P(footer_{t,n+1}) &\stackrel{\text{df}}{=} \{ ret
 \end{aligned}$$

Fig. 4. LLVM transition specific successor function, for $1 \leq k \leq n$

$succ : M \rightarrow next$	
$next \leftarrow \emptyset$ $succ_{t_1}(M, next)$ $succ_{t_2}(M, next)$... $succ_{t_n}(M, next)$ return $next$	$P(succ(x_{mrq})) \stackrel{\text{df}}{=} \begin{cases} x_{next} = fcall\ cons_{set}() \\ pcall\ succ_{t_1}(x_{mrq}, x_{next}) \\ pcall\ succ_{t_2}(x_{mrq}, x_{next}) \\ \dots \\ pcall\ succ_{t_n}(x_{mrq}, x_{next}) \\ ret\ x_{next} \end{cases}$

Fig. 5. Computation of all successors (left) and its LLVM implementation (right), where x_{mrq} is a pointer to a structure marking

In order to compare heaps, a notion of structural equivalence needs to be defined. This relation ensures that two heaps contain the same data, accessible from distinct sets of pointers but with the same layout. More precisely we consider two heaps H, H' and two pointers p, p' and write $(H, p) =_{st} (H', p')$ whenever $H(p)$ and $H'(p')$ are structurally the same values.

We also need to define an operation $new : \mathbb{H} \times \mathbb{T} \rightarrow \mathbb{H} \times \mathbb{P}$ to build new heaps, which corresponds to a pointer allocation, using a helper function $alloc : 2^{\mathbb{P}} \times \mathbb{P} \times \mathbb{T} \rightarrow \mathbb{H}$ as follows:

$$new(H, t) \stackrel{\text{df}}{=} (alloc(dom(H) \cup \{p\}, p, t), p) \\ \text{for } p \notin dom(H) \text{ a "fresh" pointer}$$

$$alloc(d, p, t) \stackrel{\text{df}}{=} \{p \mapsto (t, \perp)\} \quad \text{for } t \in \mathbb{T}_0$$

$$alloc(d, p, struct(t_0, \dots, t_n)) \stackrel{\text{df}}{=} \{p \mapsto (struct(t_0, \dots, t_n), (p_0, \dots, p_n))\} \\ \oplus alloc(d \cup \{p_0, \dots, p_n\}, p_0, t_0) \\ \oplus \dots \\ \oplus alloc(d \cup \{p_0, \dots, p_n\}, p_n, t_n) \\ \text{for } p_0, \dots, p_n \notin d \text{ "fresh" pointers}$$

It can be shown that new always returns a well formed heap, and that calling new using equivalent heaps always returns equivalent heaps.

To define subprogram calls, our memory model also defines *stacks* that contain *frames* implicitly pushed onto the stack by the inference rules in the semantics. A *frame* is a tuple $F \in \mathbb{F} \stackrel{\text{df}}{=} \mathbb{L}_\perp \times \mathbb{L} \times \mathbb{B}$ whose elements are denoted by $(l_{p,F}, l_{c,F}, \beta_F)$, where $l_{p,F}$ is the label the block the program comes from (or undefined), $l_{c,F}$ is the label of the block currently executed, and β_F is a LLVM binding representing the current evaluation context. We widen the binding functional notation to the frames, so we denote by $F(x)$ the binding $\beta_F(x)$ of x by β_F .

Like for heaps we need operations to update frames. The same operator \oplus is used because the operations are very similar, but on distinct objects. The binding overwriting operation $\oplus : \mathbb{B} \times \mathbb{B} \rightarrow \mathbb{B}$ and the frame binding overwriting operation $\oplus : \mathbb{F} \times \mathbb{B} \rightarrow \mathbb{F}$ are defined as:

$$(\beta \oplus \beta')(p) \stackrel{\text{df}}{=} \begin{cases} \beta'(p) & \text{if } p \in dom(\beta') \\ \beta(p) & \text{if } p \notin dom(\beta') \wedge p \in dom(\beta) \\ \text{undefined} & \text{otherwise} \end{cases}$$

$$(l, l', \beta) \oplus \beta' \stackrel{\text{df}}{=} (l, l', \beta \oplus \beta')$$

The structural equivalence can be widened to pairs of heaps and frames and is denoted by $(H, F) =_{st} (H', F')$ for two heaps H, H' and two frames F, F' . Intuitively, it checks that all data accessible from the frame bindings are structurally equivalent. This holds also for values stored directly in the bindings (*i.e.*, without pointers) since the heap equivalence reduces to the equality on \mathbb{D} .

5.2 Inference Rules

The operational semantics is defined for a fixed and immutable program P , which means that no function nor block can be created nor modified during the execution. We denote the result of a computation by $\overline{\cdot}$, for example $\overline{2+3}$ is 5. The main objects handled by our inference rules are *configurations* that represent a state of the program during its execution. A configuration is a tuple (seq, H, F) , denoted by $(seq)_{H,F}$, where seq is a sequence of instructions, H is a heap and F is a frame.

The inference rules for expressions are shown in figure 6; expressions evaluate to values in the context of a frame. The inference rules for sequence and commands are shown in figure 7; sequences and commands evaluate to other sequences or commands in the context of a heap and a frame. One can remark how a frame is pushed onto the stack in *pcall* and *fcall* rules, a new frame F_0 is actually replacing the current frame F in the subsumption of these rules and used to execute the body of the called subprogram. This semantics mixes up small-step and big-step reductions. Indeed, most of the rules are small-step except for *pcall* and *fcall* rules in which we link the computation to its result by making a sequence of reductions in the rule subsumption.

5.3 Data Structures Interpretation

The link between Petri nets and their LLVM implementation is formalised with a family of interpretation functions for all data structures. This allows to formalise the behavioural requirements on the interfaces presented in section 4.

An interpretation is a partial function which maps a pair formed by a heap and a pointer to a Petri net object: a marking, a set of markings, a multiset of tokens or a single token, depending on the interpreted object. Interpretations are denoted by $\llbracket H, p \rrbracket^\star$, where $H \in \mathbb{H}$, $p \in \mathbb{P} \cup \mathbb{D}$ and \star is an annotation describing the interpreted object (for instance we use $mset(t)$ instead of \star to interpret a multiset over a type t). Whenever p is a pointer, we assume that the interpretation depends only on data that is accessible from p , *i.e.*, $p \downarrow_H$. Moreover every interpretation function has to respect the following consistency requirement.

Requirement 1 (Consistency). *Let H, H' be two heaps, $\llbracket \cdot, \cdot \rrbracket^\star$ an interpretation function, and p, p' two pointers or values. If $(H, p) =_{st} (H', p')$ then $\llbracket H, p \rrbracket^\star = \llbracket H', p' \rrbracket^\star$.*

As presented in section 4, we use data structures and functions as basic blocks for constructing our algorithms, they are either predefined or produced by the compilation process. Each of these functions and data structures is specified (actually, axiomatized) by a formal interface. In particular, this helps to ensure independence and modularity between components both in a programmatic and formal way. Specifying an interface leads to define a set of primitives that respect given derivations and interpretations. For example, let H be a heap and F a frame such that, $F(x_{mset}) = p_{mset}$ is a pointer on a multiset structure storing

$$\begin{array}{c}
\frac{F(x) = p \quad H(p) = (t, v)}{(\text{load } x)_{H,F} \rightsquigarrow (v)_H} \text{load} \quad \frac{F(x) = p \quad p[i]_H \text{ is defined}}{(\text{gep } x, 0, i)_{H,F} \rightsquigarrow (p[i]_H)_H} \text{gep}_0 \\
\\
\frac{\overline{F(x_1) + F(x_2)} = v}{(\text{add } x_1, x_2)_{H,F} \rightsquigarrow (v)_H} \text{add} \quad \frac{(H', p) = \text{new}(H, t)}{(\text{alloc } t)_{H,F} \rightsquigarrow (p)_{H \oplus H'}} \text{alloc} \\
\\
\frac{\overline{F(x_1) \text{ op } F(x_2)} = v \quad \text{op} \in \{<, \leq, =, \neq, \geq, >\}}{(\text{icmp } \text{op}, x_1, x_2)_{H,F} \rightsquigarrow (v)_H} \text{icmp} \\
\\
\frac{1 \leq i \leq n \quad l_i \neq \perp}{(\text{phi } (x_1, l_1), \dots, (x_n, l_n))_{H, (l_i, l_c, \beta)} \rightsquigarrow (\beta(x_i))_H} \text{phi} \\
\\
\frac{f(a_1, \dots, a_n) \in \text{dom}(P) \quad F_0 = (\perp, f(a_1, \dots, a_n), \{a_1 \mapsto F(r_1), \dots, a_n \mapsto F(r_n)\}) \quad (P(f(a_1, \dots, a_n)))_{H, F_0} \rightsquigarrow^* (v)_{H'}}{(\text{fcall } f(r_1, \dots, r_n))_{H,F} \rightsquigarrow (v)_{H'}} \text{fcall} \\
\\
\frac{}{(\text{ret } r)_{H,F} \rightsquigarrow (F(r))_H} \text{ret}
\end{array}$$

Fig. 6. Rules for expressions

$$\begin{array}{c}
\frac{(\text{cmd})_{H,F} \rightsquigarrow (\text{seq}')_{H',F'}}{(\text{cmd}; \text{seq})_{H,F} \rightsquigarrow (\text{seq}'; \text{seq})_{H',F'}} \text{seq} \\
\\
\frac{}{(\text{skip}; \text{seq})_{H,F} \rightsquigarrow (\text{seq})_{H,F}} \text{skip} \\
\\
\frac{}{(\text{br } l)_{H, (l_p, l_c, \beta)} \rightsquigarrow (P(l))_{H, (l_c, l, \beta)}} \text{branch}_1 \\
\\
\frac{(\beta(r) = \text{true} \wedge l = l_1) \vee (\beta(r) = \text{false} \wedge l = l_2)}{(\text{br } r, l_1, l_2)_{H, (l_p, l_c, \beta)} \rightsquigarrow (P(l))_{H, (l_c, l, \beta)}} \text{branch}_2 \\
\\
\frac{f(a_1, \dots, a_n) \in \text{dom}(P) \quad F_0 = (\perp, f(a_1, \dots, a_n), \{a_1 \mapsto F(r_1), \dots, a_n \mapsto F(r_n)\}) \quad (P(f(a_1, \dots, a_n)))_{H, F_0} \rightsquigarrow^* (\text{ret})_{H',F'}}{(\text{pcall } f(r_1, \dots, r_n))_{H,F} \rightsquigarrow (\text{skip})_{H',F}} \text{pcall} \\
\\
\frac{(\text{expr})_{H,F} \rightsquigarrow (v)_{H'}}{(x = \text{expr})_{H,F} \rightsquigarrow (\text{skip})_{H', F \oplus \{x \mapsto v\}}} \text{assign} \\
\\
\frac{F(r_p) = p \quad H(p) = (t, d) \quad H' = \{p \mapsto (t, F(r_{\text{new}}))\}}{(\text{store } r_{\text{new}}, r_p)_{H,F} \rightsquigarrow (\text{skip})_{H \oplus H', F}} \text{store}
\end{array}$$

Fig. 7. Rules for sequences and commands

elements of type t . Under these conditions, procedure add_{mset} is specified by:

$$(pcall\ add_{mset}(x_{mset}, x))_{H,F} \rightsquigarrow (skip)_{H \oplus H',F} \quad (1)$$

$$dom(H) \cap dom(H') \subseteq p_{mset} \downarrow_H \quad (2)$$

$$\llbracket H \oplus H', p_{mset} \rrbracket^{mset(t)} = \llbracket H, p_{mset} \rrbracket^{mset(t)} + \{\llbracket H, F(x) \rrbracket^t\} \quad (3)$$

Condition (1) describes by a reduction the result of the call, condition (2) restrict the updates to be localised in the heap and condition (3) interprets the computation in terms of Petri nets objects. Similarly, any implementation of the marking structures has to respect the two following requirements.

Requirement 2 (Soundness). *Let H, H' be two heaps, F, F' two frames, $p_{mrq} \in dom(H)$ a pointer to a marking structure, and p_s a pointer to a place nested in p_{mrq} (i.e., $p_s \in p_{mrq} \downarrow_H$). If $\llbracket H, p_{mrq} \rrbracket^{mrq}(s) = \llbracket H, p_s \rrbracket^s, (seq)_{H,F} \rightsquigarrow (seq')_{H \oplus H',F'}$ and $p_{mrq} \notin dom(H')$ then*

$$\llbracket H \oplus H', p_{mrq} \rrbracket^{mrq}(s) = \llbracket H \oplus H', p_s \rrbracket^s$$

Requirement 3 (Separation). *Let p_{mrq} be a pointer on a structure marking in a heap H . If p_s and $p_{s'}$ are pointers to distinct places in this structure then we have $p_s \downarrow_H \cap p_{s'} \downarrow_H = \emptyset$.*

The *soundness* property ensures that any update of a place through a pointer returned by get_s is actually made on the marking (not on a copy). The *separation* property ensures that places do not share memory so that updating a place does not have side effects on other places.

6 Correctness and Termination Results

We present now the two main results proving the correctness of functions $fire_t$ (theorem 1) and $succ_t$ (theorem 2). Both these results are shown in a minimal context, i.e., a heap that just contains the required pointers. An auxiliary theorem (not presented here) allows to generalise both results to any context that includes the minimal one.

Theorem 1. *Let M be a marking, H a heap and p_{mrq} a pointer on a marking structure such that $dom(H) = p_{mrq} \downarrow_H$ and $\llbracket H, p_{mrq} \rrbracket^{mrq} = M$. Let $\beta \stackrel{\text{df}}{=} \{x_1 \mapsto v_1, \dots, x_n \mapsto v_n\}$ be a LLVM mode for transition t , which implies that each v_i is a value or a pointer encoding a token in place s_i : $\llbracket H, v_i \rrbracket^{\ell(s_i)} \in M(s_i)$. Let F be a frame such that $\beta_F \stackrel{\text{df}}{=} \beta \oplus \{x_{mrq} \mapsto p_{mrq}\}$. If*

$$M[t, \beta]M' \quad \text{and} \quad (fcall\ fire_t(x_{mrq}, x_1, \dots, x_n))_{H,F} \rightsquigarrow (p'_{mrq})_{H \oplus H'}$$

then

$$\llbracket H \oplus H', p'_{mrq} \rrbracket^{mrq} = M' \quad \text{and} \quad dom(H) \cap dom(H') = \emptyset$$

Proof (sketch). This theorem is the direct application of two lemmas showing the correctness of tokens consumption and production respectively. Both are proved by induction on the number of places in the marking structure. \square

Corollary 1. *Under the same hypothesis, the call to fire_t terminates.* \square

Theorem 2. *Let F be a frame and p_{mrq} a pointer in a heap H such that $p_{mrq} \downarrow_H = \text{dom}(H)$, $\llbracket H, p_{mrq} \rrbracket^{mrq} = M$, $\beta_F(x_{mrq}) = p_{mrq}$ and $\beta_F(x_{next}) = p_{next}$. If*

$$(fcall\ succ_t(x_{mrq}, x_{next}))_{H,F} \rightsquigarrow^* (p_{next})_{H \oplus H'} \quad \text{and} \quad \llbracket H, p_{next} \rrbracket^{set} = E$$

where E is a set of markings, then

$$\text{dom}(H) \cap \text{dom}(H') = \emptyset$$

$$\llbracket H \oplus H', p_{next} \rrbracket^{set} = E \cup \{M' \mid \exists \beta, M[t, \beta]M'\}$$

Proof (sketch). The first result is a consequence of the formal interfaces of the called functions. The second result is proved as two inclusions:

\supseteq . This is the consequence of two lemmas:

- by applying reduction rules, we show that if execution goes through a block $header_{t,k}$ then it will necessarily reach a block $footer_{t,k+1}$, and every block annotated by an index greater than k must be executed also;
- consequently, all combinations of tokens for the input places are actually enumerated, which implies that all potential modes of t are considered.

\subseteq . To prove that only actual successor markings are added, we first remark that if a marking is added into the set, then this happens in block $body_{t,0}$. So it is enough to prove that this block is executed only if the binding is actually a mode for t , which can be proved using the reduction rules backward to show that the guard necessarily evaluated to *true*. \square

Corollary 2. *Under the same hypothesis, the call to succ_t terminates.* \square

7 Conclusion

We have shown how a Petri net can be compiled targeting a fragment of the LLVM language. This compilation produces code that provides the primitives to compute the state space of the compiled Petri net model. Then we have defined a formal semantics for the fragment of the LLVM language we use. To produce a readable and usable system of inference rules, we have defined a memory model based on explicit heaps and stacks. Finally we have proved the correctness of the code generated by our compiler. The full proofs provided in [5,4] are quite long because they are very much detailed to improve our confidence into their correctness and to ease there later validation using a proof assistant. But notice

also that they are at the same time quite easy to follow. It is worth noting also that our proofs are modular thanks to clearly defined interfaces with appropriate axiomatisation. As a consequence, we should avoid issues when parts of the generated code are replaced, for instance, to use a more efficient data structure, or alternative state space exploration approaches (like in [8]).

The fragment of LLVM we have considered is rather limited with respect to the number of instructions. However, it is at the same time quite representative of the full language. Indeed, it includes the necessary to handle the stack and the heap which are conceptually the most complicated parts of the language. Extending our fragment to include all the LLVM computational instructions (like arithmetic) would be an easy but tedious work. Adding the instructions to manipulate the stack (like *unwind* for exception handling) looks quite straightforward. The most complicated is probably adding full support for pointers, which would require to refine our heap model (in particular, pointer arithmetic would have to be defined).

In this paper, we have considered an “optimistic” approach in that we assume that the LLVM code provided in the model annotations is correct and terminates. Moreover, we did not make any assumption about the finiteness of the state space or the boundedness of integer values that are assumed not to overflow. In practice, these are however important issues. Approaches based on abstract interpretation of assembly code like [18] may be helpful to prove such properties on the compiled model before to start the state space exploration, ensuring that it will run safely (and allowing to avoid implementing checks in the generated code).

Future works will address a generalisation of the presented approach to compile a wider variety of coloured Petri nets, in particular nets embedding annotation languages easier to use for the modeller than LLVM. Moreover, we would like to refine requirement 3 to allow for a logical separation instead of a physical separation as it is currently defined. This would enable us for implementing memory sharing and thus saving a lot of memory during a state space exploration. We are also interested in particular in exploiting remarkable structures of Petri net models that allow to optimise the code generated by the compiler. Such optimisations also need to be formally proved and preliminary results about this can be found in [4]. A longer term goal is to prove the whole compilation chain to obtain the core (*i.e.*, state space exploration) of a certified explicit model-checker for coloured Petri nets. A complementary aspect is to evaluate the performance of the state space generation, which is of course another important motivation when working on a model-checker. As shown in [6], the current implementation is efficient and can outperform state-of-the-art tools. So, certification is not an objective that contradicts efficiency.

References

1. Clarke, E., Emerson, A., Sifakis, J.: Model checking: Algorithmic verification and debugging. ACM Turing Award (2007)
2. ADT Coq/INRIA. The Coq proof assistant, <http://coq.inria.fr>

3. Evangelista, S.: Méthodes et outils de vérification pour les réseaux de Petri de haut niveau. PhD thesis, CNAM, Paris, France (2006)
4. Fronc, L.: Analyse efficace des réseaux de Petri par des techniques de compilation. Master's thesis, MPRI, university of Paris 7 (2010), http://www.ibisc.fr/~lfronc/pub/LF_2010.pdf
5. Fronc, L., Pommereau, F.: Proving a Petri net model-checker implementation. Technical report, IBISC, University of Évry (2010), <http://goo.gl/WMzhp>
6. Fronc, L., Pommereau, F.: Optimizing the compilation of Petri nets models. In: Proc. of SUMo 2011. CEUR, vol. 726 (2011), <http://ceur-ws.org/Vol-726>
7. Holzmann, G.J., et al.: Spin, formal verification, <http://spinroot.com>
8. Holzmann, G.J., Peled, D., Yannakakis, M.: On nested depth-first search. In: Proc. of the 2nd Spin Workshop. AMS (1996)
9. Jensen, K., Kristensen, L.M.: Coloured Petri Nets: Modelling and Validation of Concurrent Systems. Springer, Heidelberg (2009); ISBN: 978-3-642-00283-0
10. Lattner, C.: LLVM language reference manual, <http://llvm.org/docs/LangRef.html>
11. Lattner, C., et al.: The LLVM compiler infrastructure., <http://llvm.org>
12. Lattner, C., et al.: LLVM related publications., <http://llvm.org/pubs>
13. Lattner, C., et al.: LLVM users, <http://llvm.org/Users.html>
14. Pajault, C., Evangelista, S.: Helena: a high level net analyzer, <http://helena.cnam.fr>
15. Paulson, L., Nipkow, T., Wenzel, M.: The Isabelle proof assistant, <http://www.cl.cam.ac.uk/research/hvg/Isabelle>
16. Pommereau, F.: Quickly prototyping Petri nets tools with SNAKES. Petri net newsletter (2008)
17. Reinke, C.: Haskell-coloured Petri nets. In: Koopman, P., Clack, C. (eds.) IFL 1999. LNCS, vol. 1868, pp. 181–198. Springer, Heidelberg (2000)
18. Xavier Rival. Traces Abstraction in Static Analysis and Program Transformation Abstraction de Traces en Analyse Statique et Transformations de Programmes. PhD thesis, Computer Science Department, École Normale Supérieure (2005)
19. Verma, K.N., Goubault-Larrecq, J., Prasad, S., Arun-Kumar, S.: Reflecting bDDs in coq. In: Kleinberg, R.D., Sato, M. (eds.) ASIAN 2000. LNCS, vol. 1961, pp. 162–181. Springer, Heidelberg (2000)