

# A High-Level Modular Definition of the Semantics of C $\sharp$

Egon Börger<sup>a</sup>, Nicu G. Fruja<sup>b</sup>, Vincenzo Gervasi<sup>a</sup>,  
Robert F. Stärk<sup>b,\*</sup>

<sup>a</sup>*Dipartimento di Informatica, Università di Pisa,  
Via F. Buonarroti 2, I-56127 Pisa, Italy*

<sup>b</sup>*Computer Science Department, ETH Zürich, CH-8092 Zürich, Switzerland*

---

## Abstract

We propose a structured mathematical definition of the semantics of C $\sharp$  programs to provide a platform-independent interpreter view of the language for the C $\sharp$  programmer, which can also be used for a precise analysis of the ECMA standard of the language and as a reference model for teaching. The definition takes care to reflect directly and faithfully – as much as possible without becoming inconsistent or incomplete – the descriptions in the C $\sharp$  standard to become comparable with the corresponding models for Java in [1] and to provide for implementors the possibility to check their basic design decisions against an accurate high-level model. The model sheds light on some of the dark corners of C $\sharp$  and on some critical differences between the ECMA standard and the implementations of the language.

*Key words:* Semantics of programming languages, Abstract State Machines, C $\sharp$ , Java, .NET

---

## 1 Introduction

In this paper the method developed in [1] for a rigorous definition and analysis of Java and its implementation on the Java Virtual Machine (JVM) is applied to formalize the semantics of the entire language C $\sharp$ . We provide a

---

\* Corresponding author.

*Email addresses:* boerger@di.unipi.it (Egon Börger), fruja@inf.ethz.ch (Nicu G. Fruja), gervasi@di.unipi.it (Vincenzo Gervasi), staerk@inf.ethz.ch (Robert F. Stärk).

succinct, purely mathematical (thus platform-independent) model, which reflects as much as possible the intuitions and design decisions underlying the language as described in the ECMA standard [2] and in [3] and can be used as accurate and complete reference model by C# programmers, by implementors of the language and by students learning it. In Sect. 8 we point to some challenging applications of the model for proving interesting theorems about C# and its implementations.

The model clarifies a certain number of semantically relevant issues which are not handled by the ECMA standard, wherefore we also consulted the Microsoft Press books [4–6] and the documentation in [7–11]. A series of bugs and gaps in the ECMA standard for C# and in its implementation in .NET and incoherences between the two were detected during our attempt to build for the language a consistent and complete yet abstract *ground model* (in the sense described in [12]). Some of them are mentioned in this paper to shed light on some dark corners of C#, for a complete discussion we refer the reader to the companion paper [13]. As a rule we adhere to an established scientific tradition for which one of the goals of defining the meaning of programs is to accurately specify the freedom the compiler writer has for the implementation. Nevertheless we also want our model to support the practice of programming. Therefore, whenever we see for a language construct an incoherence or a to-be-closed gap between on the one side the view offered by the ECMA standard, which should support the understanding also by programmers, and on the other side the view current compilers seem to have, we give in our model a pragmatic preference to abstractly defining what the programmer is allowed to expect from the execution of his code in the current implementations of C# [7–9]. In each case we explicitly discuss the discovered discrepancy so that the parameters of the design decision become clear. To support the experimentation with the model a project has been started to refine the model developed here to .NET-executable AsmL code [14], similarly to the AsmGofer refinement developed by Joachim Schmid [15,16] for the Java and JVM models in [1].

To provide the programmer with a transparent view of the intricate interaction of various language features which depend on the run-time environment, our model comes as an *abstract interpreter*, which provides a simple way to reflect those run-time-related features encountered upon executing a given C# program. To exploit the flexibility the use of Abstract State Machines (ASMs) offers in high-level system modeling and to obtain the faithfulness and simplicity of abstract models the ASM method allows one to achieve, the interpreter takes the form of an ASM.<sup>1</sup> This allows us in particular to specify the static and the dynamic parts of the semantics separately, due to the ASM

---

<sup>1</sup> See Sect. 8 for more information on our choice of ASMs among the many frameworks in the literature to deal with language semantics.

classification of abstract states into a static and a dynamic part. The *dynamic semantics* of the language is captured operationally by ASM rules which describe the run-time effect of program execution on the abstract state of the program, the *static semantics* comes as a mainly declarative description of the relevant syntactical and compile-time checked language features (like typing rules, rules for definite assignment and reachability, name resolution, method resolution for overloaded methods, etc.) and of pre-processing directives (like `#define`, `#undef`, `#if`, `#else`, `#endif`, etc.), which are mostly reflected in the attributed abstract syntax tree our model starts from.

To keep the size of the models small and to facilitate the understanding of clusters of language constructs in terms of local state transformations, similarly to the decomposition of Java and the JVM in [1] we structure the  $C\sharp$  programming language into *layered modules of orthogonal language features*, namely

- the imperative core, related to sequential control by while programs, built from statements and expressions over the simple types of  $C\sharp$ ,
- classes, realizing procedural abstraction with global (module) variables and class initialization,
- object-orientation with class instances, instance methods, inheritance,
- exception handling,
- delegates together with events (including here for convenience also properties, indexers, attributes),
- concurrency (threads),
- so-called unsafe code with pointers and pointer arithmetic.

This yields a sequence of sublanguages  $C\sharp_I$ ,  $C\sharp_C$ ,  $C\sharp_O$ ,  $C\sharp_E$ ,  $C\sharp_D$ ,  $C\sharp_T$ ,  $C\sharp_U$  which altogether describe the entire language  $C\sharp$ . Each language  $L$  in the sequence extends its predecessor and for each one we build a submachine  $EXEC\text{SHARP}_L$  which is a conservative (purely incremental) extension of its predecessor. The model  $EXEC\text{SHARP}$  for the entire language  $C\sharp$  is a composition of all submachines.

$$\begin{aligned} EXEC\text{SHARP} &\equiv \\ &EXEC\text{SHARP}_I \\ &EXEC\text{SHARP}_C \\ &EXEC\text{SHARP}_O \\ &EXEC\text{SHARP}_E \\ &EXEC\text{SHARP}_T \\ &EXEC\text{SHARP}_D \\ &EXEC\text{SHARP}_U \end{aligned}$$

This approach supports a systematic piecemeal introduction of the numerous language constructs in teaching  $C\sharp$  (or similar programming languages).

To keep the definition of the models succinct, we avoid tedious and routine rep-

etitions concerning language constructs which can be reduced in well-known ways to the core constructs in our models. Whenever instead of a direct formalization of a construct we use a syntactical translation to constructs dealt with in the core model, we have to justify that the translation is correct with respect to the semantics of the construct as intended by the standard. The ASM model we define provides a basis to rigorously formulate and mathematically prove the intended equivalence.<sup>2</sup> Since such a justification follows well-known patterns, it is skipped in this paper, but to remind the reader of the problem we usually mention it.

The handling of truly concurrent threads, not limited to interleaving or similar simple scheduling techniques, is closely related to the underlying memory model. Since the description of this memory model goes much beyond this paper, the submodel  $C\sharp_{\mathcal{T}}$  and its further analysis is postponed to a separate paper [17].

By and large one can correctly understand an ASM as pseudo-code operating over abstract data (structures in the sense of logic). Therefore we skip a detailed definition of ASMs, which is available in textbook form in Chapter 2 of the *AsmBook* [18]. Since our paper is not a tutorial or manual on  $C\sharp$ , we restrict our explanations of language constructs to features a reader will appreciate who is already knowledgeable about the basic concepts of object-oriented programming. In a technical report [19] also the remaining details which are skipped in this paper are spelt out completely, together with further explanations and examples.

The paper is structured by the modularization we propose for the language description. The basic framework of our model is introduced in Sect. 2 together with the interpreter for the imperative kernel  $C\sharp_{\mathcal{I}}$  of the language. Successively one more section is added for each model refinement to capture the related language extension. In general each section has a first part where the static assumptions of the model are formulated, followed by a second part which contains the dynamics expressed by the ASM transition rules operating on the corresponding state components. In general at each layer the interpreter consists of two submachines, one defining expression evaluation and one defining statement execution.

---

<sup>2</sup> One has to define an extension of the core model by a direct formalization of the construct in question and then to prove that this model is equivalent to the core model modulo the syntactical translation of the construct.

$$\begin{aligned}
Exp & ::= Lit \mid Vexp \mid Uop \ Exp \mid Exp \ Bop \ Exp \mid Exp \ '?' \ Exp \ ':' \ Exp \\
& \quad \mid ( \ Type \ ') \ ExpSexp \mid ( \ Exp \ ') \mid \text{'checked'} \ ( \ Exp \ ') \\
& \quad \mid \text{'unchecked'} \ ( \ Exp \ ') \\
Vexp & ::= Loc \\
Sexp & ::= Vexp \ '=' \ Exp \mid Vexp \ Aop \ Exp \mid Vexp \ '++' \mid Vexp \ '--' \\
Uop & ::= '+' \mid '-' \mid '!' \mid '~' \\
Bop & ::= '*' \mid '/' \mid '%' \mid '+' \mid '-' \mid '<<' \mid '>>' \mid '<' \mid '>' \mid '<=' \mid '>=' \mid '==' \\
& \quad \mid '!=' \mid '&' \mid '^' \mid '|' \\
Aop & ::= '*=' \mid '/=' \mid '%=' \mid '+=' \mid '-=' \mid '<<=' \mid '>>=' \mid '&=' \mid '^=' \mid '|=' \\
Stm & ::= ';' \mid Sexp \ ';' \mid \text{'break'} \ ';' \mid \text{'continue'} \ ';' \mid \text{'goto'} \ Lab \ ';' \\
& \quad \mid \text{'if'} \ ( \ Exp \ ') \ Stm \ \text{'else'} \ Stm \\
& \quad \mid \text{'while'} \ ( \ Exp \ ') \ Stm \mid \text{'do'} \ Stm \ \text{'while'} \ ( \ Exp \ ') \\
& \quad \mid \text{'for'} \ ( \ [Sexp] \ ';' \ [Exp] \ ';' \ [Sexp] \ ') \ Stm \\
& \quad \mid \text{'switch'} \ ( \ Exp \ ') \ '{ \ {Case \ {Case} \ Bstm \ {Bstm}} \ }' \\
& \quad \mid \text{'goto'} \ \text{'case'} \ Cexp \ ';' \mid \text{'goto'} \ \text{'default'} \ ';' \\
& \quad \mid \text{'checked'} \ Block \mid \text{'unchecked'} \ Block \mid Block \\
Sexp & ::= Sexp \ '{ \ , \ Sexp \ } \\
Case & ::= \text{'case'} \ Cexp \ ':' \mid \text{'default'} \ ':' \\
Block & ::= '{ \ {Bstm} \ }' \\
Bstm & ::= Type \ Loc \ ';' \mid \text{'const'} \ Type \ Loc \ '=' \ Cexp \ ';' \mid Lab \ ':' \ Stm \mid Stm
\end{aligned}$$

Fig. 1. Grammar of expressions and statements in  $C\sharp_{\mathcal{I}}$ .

## 2 The imperative core $C\sharp_{\mathcal{I}}$

In this section we define the model for  $C\sharp_{\mathcal{I}}$ , which defines the basic machinery of the ASM model for  $C\sharp$ . It describes the semantics of the sequential imperative core of  $C\sharp$  with to be executed statements (appearing in method bodies) and to be evaluated expressions (appearing in statements) built using predefined operators over simple types. The computations of this interpreter are supposed to start with an arbitrary but fixed  $C\sharp$  program. We separate syntax and compile-time matters from run-time issues by assuming that the program is given as an attributed syntax tree (i.e. annotated abstract syntax tree resulting from parsing and elaboration), trying to achieve model simplicity also by assuming some useful syntactical simplifications which will be mentioned as we build the model. Before defining the transition rules for the dynamic semantics of  $C\sharp_{\mathcal{I}}$ , we formulate what has to be said about the static semantics.

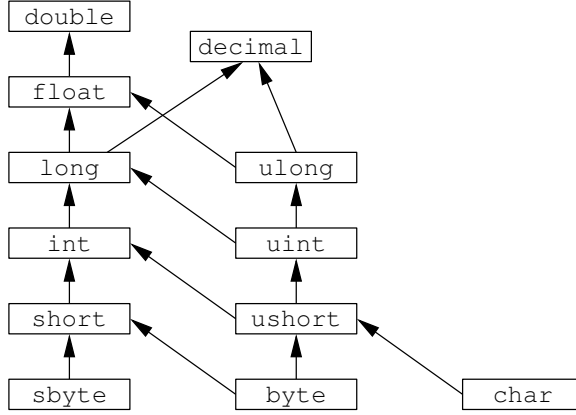


Fig. 2. The simple types of  $C\#_{\mathcal{I}}$ .

### 2.1 Static semantics of $C\#_{\mathcal{I}}$

We view the grammar in Fig. 1, which defines expressions and statements of the sublanguage  $C\#_{\mathcal{I}}$ , as defining also the corresponding ASM domains  $Exp$  and  $Stm$ . To avoid lengthy repetitions we include here already the distinctions between checked and unchecked expressions and blocks, though they are semantically irrelevant in the submodel  $C\#_{\mathcal{I}}$  and start to play a role only with  $C\#_{\mathcal{E}}$ . The set  $Vexp$  of variable expressions (lvalues) consists in this model of the local variables only and will be refined below.  $Sexp$  denotes the set of statement expressions than can be used either as (result yielding) expressions or as (result discarding) statements, such as an assignment to a variable expression using '=' or an assignment operator from the set  $Aop$  or '++' or '--'.  $Lit$  denotes the set of literals, similarly for  $Type$ ,  $Lab$  and the set  $Cexp$  of constant expressions whose value is known at compile time. When referring to the set of sequences of elements from a set  $Item$  we write  $Items$ , e.g.  $Sexp_s$  for the set of sequences of statement expressions. We usually write lower case letters  $e$  to denote elements of a set  $E$ , e.g.  $lit$  for elements of  $Lit$ .

The descriptions of implicit numeric conversions in [2, §13.1] and of binary numeric promotions in [2, §14.2.6] can be succinctly formulated as follows, using the type graph in Fig. 2 for the simple types of  $C\#$ , which are the types of  $C\#_{\mathcal{I}}$  (for a classification of the types of  $C\#$  see Fig. 4).

**Definition 1 (Implicit conversion [2, §13.1])** *We say that there exists an implicit numeric conversion from type  $A$  to  $B$  (written  $A \prec B$ ) iff there exists a finite, non-empty path of arrows from  $A$  to  $B$  in the type graph in Fig. 2. We write  $A \preceq B$  for  $A \prec B$  or  $A = B$ . A type  $C$  is called an upper bound of  $A$  and  $B$  iff  $A \preceq C$  and  $B \preceq C$ . A type  $C$  is the least upper bound of  $A$  and  $B$  iff*

- $C$  is an upper bound of  $A$  and  $B$  and

- $C \preceq D$  for each upper bound  $D$  of  $A$  and  $B$ .

We write  $\text{sup}(A, B)$  for the least upper bound of  $A$  and  $B$  if it exists.

We assume all the type constraints (on the operand and result values) and precedence conventions listed in [2] for the predefined (arithmetical, relational, bit and boolean logical) operators and the expression types. As usual each expression node  $exp$  in the attributed syntax tree has as attribute its compile-time type  $\text{type}(exp)$ .

About type conversions at compile-time we assume that type casts are inserted in the syntax tree if necessary. For example, if a binary numeric operator  $bop$  is applied to arguments in  $e_1 \ bop \ e_2$ , then the least upper bound  $T$  of the types of  $e_1$  and  $e_2$  must exist and the expression is transformed into  $(T)e_1 \ bop \ (T)e_2$ .

**Definition 2 (Binary numeric promotion [2, §14.2.6])**

The binary numeric promotion consists of applying the following rules:

- If the least upper bound of  $A$  and  $B$  exists, then
  - if  $\text{sup}(A, B) \preceq \text{int}$ , then  $A$  and  $B$  are converted to  $\text{int}$ ,
  - otherwise,  $A$  and  $B$  are converted to  $\text{sup}(A, B)$ .
- If the least upper bound of  $A$  and  $B$  does not exist, then a compile-time error occurs.

We also assume the syntactical constraints for statements listed in [2], e.g. the following ones for blocks (where the *scope of a local variable (local constant)* is defined as the block in which it is declared, the *scope of a label* is the block in which the label is declared, and a local variable is identified by its name *and* the position of its declaration, so that in particular local variables with the same name in disjoint blocks are considered as different):

- It is not allowed to refer to a local variable (local constant) in a textual position that precedes its declaration.
- It is not allowed to declare another local variable or local constant with the same name in the scope of a local variable (local constant).
- It is not allowed for two labels with the same name to have overlapping scopes.
- A `goto Lab` must be in the scope of a label with name *Lab*.
- Expressions in *constant declarations* are evaluated at compile-time.

To simplify the exposition of our model we assume some standard syntactical reductions as indicated in Table 1. The correctness of these replacements with respect to [2] can easily be checked on the basis of our semantics model for  $C\ddagger$ .

**Control-flow analysis.** During the static program analysis where the compiler has to verify that the given program is well-typed, predicates *reachable*

Table 1  
Standard syntactical reductions.

$exp_1 \ \&\& \ exp_2$	$exp_1 \ ? \ exp_2 \ : \ \mathbf{false}$
$exp_1 \    \ exp_2$	$exp_1 \ ? \ \mathbf{true} \ : \ exp_2$
$\mathbf{if} \ (exp) \ stm$	$\mathbf{if} \ (exp) \ stm \ \mathbf{else} \ ;$
$++vexp$	$vexp \ += \ 1$
$--vexp$	$vexp \ -= \ 1$
$\mathbf{int} \ x = 1, \ y, \ z = x * 2;$	$\mathbf{int} \ x; \ x = 1; \ \mathbf{int} \ y; \ \mathbf{int} \ z; \ z = x * 2;$
$\mathbf{for} \ (t \ loc = exp; \ tst; \ step) \ stm$	$\{ \ t \ loc; \ \mathbf{for} \ (loc = exp; \ tst; \ step) \ stm \}$

and *normal* with the following intended meaning are computed for statements, using the type information contained in the attributed syntax tree as the result of parsing and elaboration:

$$\begin{aligned}
 \mathit{reachable}(stm) &\iff stm \text{ can be reached} \\
 \mathit{normal}(stm) &\iff stm \text{ can terminate normally} \\
 &\iff \text{the end point of } stm \text{ can be reached}
 \end{aligned}$$

One of the language design goals was to guarantee the following two properties for programs to be accepted by the compiler:

- during the program execution, only *reachable* positions are reached,
- normal termination happens only in *normal* positions.

These two properties are obtained by checking two sufficient conditions via so-called reachability rules, which can be inductively defined for  $C\sharp_{\mathcal{I}}$  in Table 2 (similarly for **do**, **for**, **switch**).<sup>3</sup> For constant boolean expressions in conditional and while statements we assume that they are replaced in the abstract syntax tree by **true** or **false**.

Unreachable statements indicate programming errors and therefore generate compile-time warnings. Function bodies that can terminate normally generated compile-time errors, since at run-time execution could fall off the bottom of the code array.

Another language design goal was to achieve the type safety of well-typed  $C\sharp$  programs, i.e. that a) variables at run-time contain values that are *compatible* with the declared types, and b) expressions are evaluated at run-time to values that are *compatible* with their compile-time types. Among the desired consequences of the type safety of a program one has that at run-time

<sup>3</sup> We include these rules here to place the corresponding natural language specification in [2] on a firm ground for a mathematical proof of the above two properties as part of a type safety proof for  $C\sharp$ .



Table 2

Reachability rules for  $C\sharp_{\mathcal{I}}$ .

$s$ is a function body	$\implies$	$reachable(s)$
$reachable(;)$	$\implies$	$normal(;)$
$reachable(e;)$	$\implies$	$normal(e;)$
$reachable(\{ \})$	$\implies$	$normal(\{ \})$
$reachable(\{s \dots \})$	$\implies$	$reachable(s)$
$normal(s_i)$ in $\{ \dots s_i s_{i+1} \dots \}$	$\implies$	$reachable(s_{i+1})$
$reachable(\text{goto } l;)$ in $\{ \dots l : s \dots \}$	$\implies$	$reachable(l : s)$
$normal(s)$	$\implies$	$normal(\{ \dots s \})$
$reachable(\text{if } (e) s_1 \text{ else } s_2) \wedge e \neq \text{false}$	$\implies$	$reachable(s_1)$
$reachable(\text{if } (e) s_1 \text{ else } s_2) \wedge e \neq \text{true}$	$\implies$	$reachable(s_2)$
$normal(s_1) \vee normal(s_2)$	$\implies$	$normal(\text{if } (e) s_1 \text{ else } s_2)$
$reachable(\text{while } (e) s) \wedge e \neq \text{false}$	$\implies$	$reachable(s)$
$reachable(\text{while } (e) s) \wedge e \neq \text{true}$	$\implies$	$normal(\text{while } (e) s)$
$reachable(\text{break};)$ in $s$	$\implies$	$normal(\text{while } (e) s)$

its variables will never contain *undefined* values, that there are no *dangling* references, that the program cannot *corrupt* the memory, and that the dynamic *method lookup* always succeeds. Using the notation explained in the next section such invariants can be made precise and be proven to hold under appropriate assumptions.<sup>4</sup>

To guarantee the type safety the compiler checks a sufficient condition computing predicates *before*, *after* (for occurrences of statements and expressions in a function body) and *true*, *false* (for the two possible evaluation results of boolean expressions), which implement the so-called definite assignment rules to assure that a variable is *definitely assigned* before its value is used. The situation is illustrated in Fig. 3. Unfortunately the picture does not reflect reality. Microsoft has decided that in verified IL (intermediate language) code local variables are initialized by the run-time system with zero values.<sup>5</sup> Hence, also

<sup>4</sup> For example the following invariants can be proved to hold at run-time: a)  $before(pos) \subseteq Defined$  where  $Defined = \{x \in Loc \mid mem(locals(x)) \neq Undefined\}$ , b)  $after(pos) \subseteq Defined$  if  $values(pos) = Norm$  or  $values(pos) \in Value$ . Specifically for boolean expressions holds  $true(pos) \subseteq Defined$  if  $values(pos) = True$ , the same for *false*. Such proofs can be carried out on the basis of the model developed in this paper, using the pattern developed in [1, Ch.8] for proving that Java is type safe. For a different approach see [20].

<sup>5</sup> Maybe to simplify the job of the JIT verifiers, as one of our referees suggested.

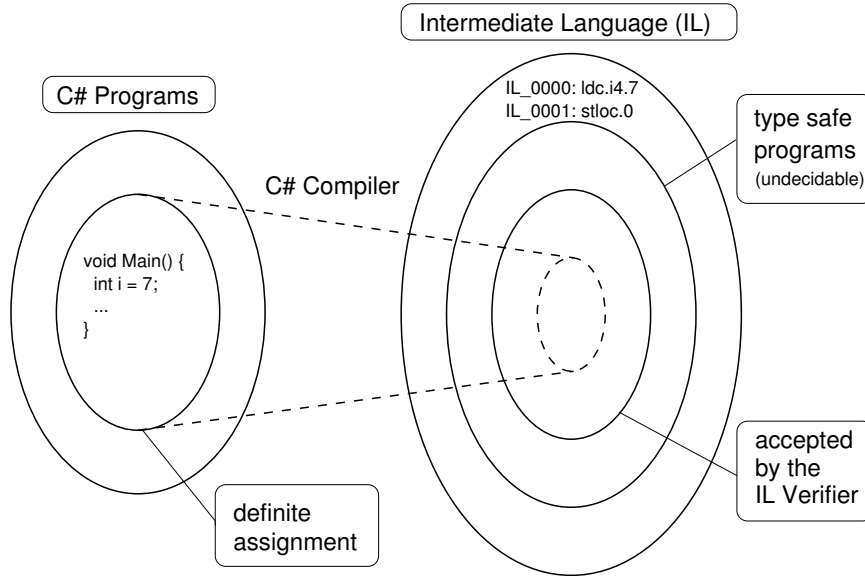


Fig. 3. Definite assignment and IL verification.

source code programs that do not fulfill the definite assignment constraints are accepted by the IL verifier.

A variable occurring in a position is called definitely assigned there, if on every execution path leading to that position (in the abstract syntax tree) a value is assigned to the variable. Thus the intended meaning of the above predicates is as follows, where by “elaboration” of an *item* we mean “execution”, if *item* is a statement, and “evaluation” if it is an expression.

$x \in \text{before}(\text{item})$ :  $x$  is definitely assigned *before* the elaboration of *item*

$x \in \text{after}(\text{item})$  :  $x$  is definitely assigned *after* normal elaboration of *item*

$x \in \text{true}(\text{exp})$  :  $x$  is definitely assigned *after* *exp* evaluates to *true*

$x \in \text{false}(\text{exp})$  :  $x$  is definitely assigned *after* *exp* evaluates to *false*

To provide a basis for a mathematical analysis, we turn the verbally stated definite assignment rules of [2, §12.3.3] into a precise set of equational constraints, where  $\text{vars}(\text{stm}) = \{x \mid \text{stm is in the scope of } x\}$ .

Table 3 contains the constraints for the statements. Table 4 contains the equations for specific boolean expressions, which are imposed for the eager (short-circuit) evaluation of boolean expressions. Note that there is no equation in Table 4 for *after* sets since by definition  $\text{after}(\text{exp}) = \text{true}(\text{exp}) \cap \text{false}(\text{exp})$ . If *exp* is a boolean expression which is not an instance of one of the expressions in Table 4, then the following are constraints for *exp*:  $\text{true}(\text{exp}) = \text{after}(\text{exp})$  and  $\text{false}(\text{exp}) = \text{after}(\text{exp})$ .

Table 3

Definite assignment for statements.

$s$ is a function body	$before(s) = \emptyset$
$;$	$after(;) = before(;)$
$exp;$	$before(exp) = before(exp;), after(exp;) = after(exp)$
<b>break;</b>	$after(\mathbf{break};) = vars(\mathbf{break};)$
<b>continue;</b>	$after(\mathbf{continue};) = vars(\mathbf{continue};)$
<b>goto</b> $l;$	$after(\mathbf{goto} \ l; ) = vars(\mathbf{goto} \ l; )$
$stm = \{s_1 \dots s_n\}$	$before(s_1) = before(stm), after(stm) = after(s_n),$ $before(s_{i+1}) = after(s_i) \cap goto(s_{i+1})$ where $goto(l:s) =$ $\bigcap \{before(\mathbf{goto} \ l; ) \mid \mathbf{goto} \ l; \text{ reachable in } stm\}$ and $goto(s) = vars(s)$ if $s$ is not a labeled statement
$stm = \mathbf{if} \ (e) \ s_1 \ \mathbf{else} \ s_2$	$before(e) = before(stm), before(s_1) = true(e)$ $before(s_2) = false(e), after(stm) = after(s_1) \cap after(s_2)$
$stm = \mathbf{while} \ (e) \ s$	$before(e) = before(stm), before(s) = true(e),$ $after(stm) = false(e) \cap break(s)$ where $break(s) = \bigcap \{before(\mathbf{break};) \mid \mathbf{break}; \text{ reachable in } s\}$

Table 4

Definite assignment for boolean expressions.

<b>true</b>	$true(\mathbf{true}) = before(\mathbf{true}), false(\mathbf{true}) = vars(\mathbf{true})$
<b>false</b>	$true(\mathbf{false}) = vars(\mathbf{false}), false(\mathbf{false}) = before(\mathbf{false})$
$exp = !e$	$before(e) = before(exp), true(exp) = false(e)$ $false(exp) = true(e)$
$exp = (e_1 \ \&\& \ e_2)$	$before(e_1) = before(exp), before(e_2) = true(e_1),$ $true(exp) = true(e_2), false(exp) = false(e_1) \cap false(e_2)$
$exp = (e_1 \    \ e_2)$	$before(e_1) = before(exp), before(e_2) = false(e_1),$ $true(exp) = true(e_1) \cap true(e_2), false(exp) = false(e_2)$
$exp = (e_0 \ ? \ e_1 \ : \ e_2)$	$before(e_0) = before(exp), before(e_1) = true(e_0)$ $before(e_2) = false(e_0), true(exp) = true(e_1) \cap true(e_2)$ $false(exp) = false(e_1) \cap false(e_2)$

Table 5 contains the equations for non-boolean expressions. In all other cases, if  $exp$  is an expression which has the *direct subexpressions*  $e_1, e_2, \dots, e_n$ , then the left-to-right evaluation scheme yields

- (1)  $before(e_1) = before(exp)$
- (2)  $before(e_{i+1}) = after(e_i)$  for  $i \in [1 .. n - 1]$
- (3)  $after(exp) = after(e_n)$

Table 5

Definite assignment for arbitrary expressions.

$loc$	$loc \in before(loc), after(loc) = before(loc)$
$lit$	$after(lit) = before(lit)$
$exp = (loc = e)$	$before(e) = before(exp), after(exp) = after(e) \cup \{loc\}$
$exp = (loc op = e)$	$loc \in before(exp), before(e) = before(exp)$ $after(exp) = after(e)$
$exp = (e_0 ? e_1 : e_2)$	$before(e_0) = before(exp), before(e_1) = true(e_0)$ $before(e_2) = false(e_0), after(exp) = after(e_1) \cap after(e_2)$

Due to the goto statement the above constraints do not specify in a unique way the sets of variables that have to be considered as definitely assigned. Consider the following block (from [21]):

```
{ int i = 1 ; L : goto L ; }
```

Then the constraints of the definite assignment analysis are satisfied for both  $before(L : \text{goto } L;) = \emptyset$  and  $before(L : \text{goto } L;) = \{i\}$ . Hence during the analysis the greatest sets of variables that satisfy the constraints for *before* and *after* have to be computed (cf. [21]). For blocks without goto statements, however, it can be proved from the above axioms that the *before* set determines the *after* set in a unique way.

## 2.2 Dynamic semantics for $C\sharp_I$

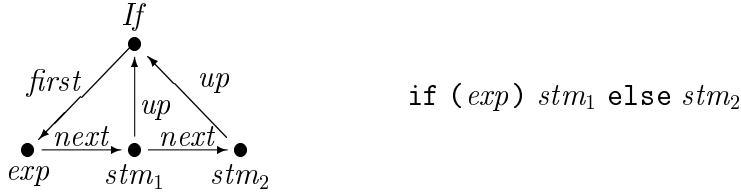
The dynamic semantics for  $C\sharp_I$  describes the effect of statement execution and expression evaluation upon the program state, so that the transition rule for  $C\sharp_I$  (the same for its extensions) has the form

EXEC $C\sharp_I \equiv$   
EXEC $C\sharp_I$ EXP $I$   
EXEC $C\sharp_I$ STM $I$

The first subrule defines one execution step in the evaluation of expressions; the second subrule defines one step in the execution of statements.

To make the further model refinements possible via purely incremental extensions, our definition proceeds by walking through the attributed syntax tree and computing at each node the effect of the program construct attached to the node. We formalize the walk by a cursor  $\blacktriangleright$ , whose position in the tree – represented by a dynamic function  $pos: Pos$  – is updated using static tree functions, leading from a node in the tree down to its *first* child, from there to the *next* brother or *up* to the parent node (if any), as illustrated by the following

self-explanatory example. The moves of  $pos$  contain implicitly the control-flow graph of  $C\sharp$ . A function  $label: Pos \rightarrow Label$  decorates nodes with the information which identifies the grammar rule associated to the node. For the sake of notational succinctness we use concrete syntax of  $C\sharp$  to describe the labels, thus avoiding the explicit introduction of auxiliary non-terminals the reader probably does not want to see. In the following example the four possible cursor positions are reachable from the root by following the tree functions  $first$ ,  $next$  and  $up$ . The  $label$  of the root node is the auxiliary non-terminal  $If$ , identifying the grammar rule which produces in one step  $if (exp) stm_1 \text{ else } stm_2$ .



For updating the values of local variables in the memory we use two dynamic functions  $locals: Loc \rightarrow Adr$  and  $mem: Adr \rightarrow SimpleValue \cup \{Undef\}$ , which assign to local variables memory addresses and store the values there. Since in  $C\sharp_{\mathcal{I}}$  the values are of simple types, the equation  $Value = SimpleValue \cup Adr$  holds, which will be refined in the extended models to include also references and structs. The uniquely identified local variables are modeled by stipulating  $Loc = Identifier \times Pos$ , where  $Pos$  is the set of positions in the abstract syntax tree.

The indirection through memory addresses is not really needed in  $C\sharp_{\mathcal{I}}$ . In  $C\sharp_{\mathcal{I}}$  one could assign values directly to local variables without storing them in an abstract memory. The addresses, however, are needed later for call-by-reference with  $ref$  and  $out$  parameters (one of the major differences between  $C\sharp$  and Java from the modelling point of view).

Statements can terminate normally or abruptly, where in  $C\sharp_{\mathcal{I}}$  the reasons of abruption are from the set  $Abr = Break \mid Continue \mid Goto(Lab)$ , to be refined for the extended models. We use a dynamic function  $values: Pos \rightarrow Result$  to store intermediate evaluation results from the set

$$Result = Value \cup Abr \cup \{Undef, Norm\}.$$

For the initial state we assume

- $mem(i) = Undef$  for every  $i \in Adr$
- $pos = \text{root position of the attributed syntax tree}$
- $locals(x) \in Adr$  for every variable  $x$ <sup>6</sup>

<sup>6</sup> This amounts to assuming that the compiler chooses an address for each variable.

As intermediate *values* at a position  $p$  the cursor is at or is passing to, the computation may yield directly a simple value; at *AddressPositions* as defined below it may yield an address; but it may also yield a *memValue* which has to be retrieved indirectly via the given address (where for  $C\sharp_{\mathcal{I}}$  the memory value of a given type  $t$  at a given address  $adr$  is defined by  $memValue(adr, t) = mem(adr)$ ; the parameter  $t$  will become relevant only in the refinement of *memValue* in  $C\sharp_{\mathcal{O}}$  and  $C\sharp_{\mathcal{U}}$ ). This is described by the following two macros:

$$\begin{aligned} \text{YIELD}(val, p) &\equiv \\ &\quad values(p) := val \\ &\quad pos := p \end{aligned}$$

$$\begin{aligned} \text{YIELDINDIRECT}(adr, p) &\equiv \\ &\quad \text{if } AddressPos(p) \text{ then } \text{YIELD}(adr, p) \\ &\quad \text{else } \text{YIELD}(memValue(adr, type(p)), p) \end{aligned}$$

We will use the macros in the two forms  $\text{YIELD}(val) \equiv \text{YIELD}(val, pos)$  and  $\text{YIELDUP}(val) \equiv \text{YIELD}(val, up(pos))$ . Similarly we have two forms also for the second macro:  $\text{YIELDINDIRECT}(adr)$  and  $\text{YIELDUPINDIRECT}(adr)$ .

Being in a context where an address and not a value is required can be defined as follows:

$$\begin{aligned} AddressPos(\alpha) &\iff FirstChild(\alpha) \wedge \\ &\quad (label(up(\alpha)) \in \{++, --\} \vee label(up(\alpha)) \in Aop) \\ \text{where } FirstChild(\alpha) &\iff first(up(\alpha)) = \alpha \end{aligned}$$

To further reduce any notational overhead not needed by the human reader, in spelling out the rules below we identify positions with the occurrences of the syntactical constructs nodes are decorated with. This explains updates like  $pos := exp$  or  $pos := stm$ , which are used as shorthand for updating  $pos$  to the node labeled with the corresponding occurrence of  $exp$  respectively  $stm$ .<sup>7</sup> Furthermore, for a succinct formulation we use a macro  $context(pos)$  to describe the context of the currently to be handled expression or statement or intermediate result, which has to be matched against the syntactically possible cases (in the textual order of the rule) to select the appropriate computation step. If the elaboration of the subtree at the position  $pos$  has not yet started, then  $context(pos)$  is the construct encoded by the labels of  $pos$  and of its children. Otherwise, if  $pos$  carries already its result in *values*,  $context(pos)$  is the pseudo-construct encoded by the labels of the parent node of  $pos$  and of its children after replacing the already evaluated constructs by their *values* in the corresponding node. This explains notations like  $uop \blacktriangleright val$

<sup>7</sup> An identification of this kind, which is common in mathematics, has clearly to be resolved in an executable version of the model. See for example the formulation of the ASM model for Java in [1].

to describe the *context* of  $pos$ , where  $pos$  is marked with the cursor ( $\blacktriangleright$ ), resulting from the successful evaluation of the argument  $exp$  of the construct  $uop\ exp$  (encoded by  $up(pos)$  and its child  $pos$ ), just before  $uop$  is applied to  $val$  to  $\text{YIELDUP}(Apply(uop, val))$ .

It thus remains to define the two submachines for expression evaluation and statement execution. This is done in a modular fashion, grouping behaviorally similar instructions into one parameterized instruction.<sup>8</sup>

**Expression evaluation rules.** We are now ready to define the machine  $\text{EXECSHARPEXP}_I$  in a compositional way, namely proceeding expression-wise: for each syntactical form of expressions there is a set of rules covering each intermediate phase of their evaluation. The machine passes control from unevaluated expressions to the appropriate subexpressions until an atom (a literal or a local variable) is reached. It can continue its computation only as long as no operator exception occurs, as a consequence it does not distinguish between checked and unchecked expression evaluation – the extension by rules to handle exceptions is defined in the model extension  $C\#_{\mathcal{E}}$ . The expressions for numeric casts will be refined in  $C\#_{\mathcal{O}}$  and in  $C\#_{\mathcal{E}}$ . The macro  $\text{WRITEMEM}(adr, t, val)$  denotes here  $mem(adr) := val$ ; it will be refined in the model for  $C\#_{\mathcal{O}}$ .

```

EXECSHARPEXPI ≡ match context(pos)
  lit → YIELD(ValueOfLiteral(lit))
  loc → YIELDINDIRECT(locals(loc))

  uop exp → pos := exp
  uop  $\blacktriangleright$  val → if  $\neg$ UopException(uop, val) then YIELDUP(Apply(uop, val))

  exp1 bop exp2 → pos := exp1
   $\blacktriangleright$  val bop exp → pos := exp
  val1 bop  $\blacktriangleright$  val2 → if  $\neg$ BopException(bop, val1, val2) then
    YIELDUP(Apply(bop, val1, val2))

  exp0 ? exp1 : exp2 → pos := exp0
   $\blacktriangleright$  val ? exp1 : exp2 → if val then pos := exp1 else pos := exp2
  True ?  $\blacktriangleright$  val : exp → YIELDUP(val)
  False ? exp :  $\blacktriangleright$  val → YIELDUP(val)

  loc = exp → pos := exp
  loc =  $\blacktriangleright$  val → {WRITEMEM(locals(loc), type(loc), val), YIELDUP(val)}

  (type) exp → pos := exp
  (type)  $\blacktriangleright$  val → if type(pos) ∈ NumericType ∧ type ∈ NumericType then
    if  $\neg$ UopException(type, val) then
      YIELDUP(Convert(type, val))

```

<sup>8</sup> The specializations can be regained instruction-wise by mere parameter expansion, a form of refinement that is easily proved to be correct.

```

vexp op= exp → pos := vexp
▶ adr op= exp → pos := exp
adr op= ▶ val → let t = type(up(pos)) and v = memValue(adr, t) in
    if  $\neg$ BopException(op, v, val) then
        let w = Apply(op, v, val) and r = Convert(t, w) in
            {WRITEMEM(adr, t, r), YIELDUP(r)}

vexp op → pos := vexp // for postfix operators op ∈ {++, --}
▶ adr op → let old = memValue(adr, type(pos)) in
    if  $\neg$ UopException(op, old) then
        WRITEMEM(adr, type(up(pos)))
        Apply(op, old)
        YIELDUP(old)

checked(exp) → pos := exp
checked(▶val) → YIELDUP(val)
unchecked(exp) → pos := exp
unchecked(▶val) → YIELDUP(val)

```

Being in a checked context is used to define whether operators throw an overflow exception (in which case a rule will be added in the model for  $C\#_{\mathcal{E}}$ ). The general rule is that operators for the type `decimal` always throw overflow exceptions whereas operators for integral types only throw overflow exceptions in a checked context except for the division by zero. By default every position is unchecked, unless explicitly declared otherwise.

$$\begin{aligned}
 \text{Checked}(\alpha) &\iff \text{label}(\alpha) = \text{Checked} \vee \\
 &\quad (\text{label}(\alpha) \neq \text{Unchecked} \wedge \text{up}(\alpha) \neq \text{Undef} \wedge \text{Checked}(\text{up}(\alpha))) \\
 \text{UopException}(uop, val) &\iff \text{Checked}(pos) \wedge \text{Overflow}(uop, val) \\
 \text{BopException}(bop, val_1, val_2) &\iff \\
 &\quad \text{DivisonByZero}(bop, val_2) \vee \text{DecimalOverflow}(bop, val_1, val_2) \vee \\
 &\quad (\text{Checked}(pos) \wedge \text{Overflow}(bop, val_1, val_2))
 \end{aligned}$$

**Statement execution rules.** The machine  $\text{EXEC}\text{SHARP}\text{STM}_I$  is defined below statement-wise. It transfers control from structured statements to the appropriate substatements, until the current statement has been computed normally or abruptly the computation. Abruptions trigger the control to propagate through all the enclosing statements up to the target labeled statement. The concept of propagation is defined for  $C\#_{\mathcal{I}}$  in such a way that in the refined models it is easily extended to abruptions due to return from procedures or to exceptions.<sup>9</sup> In case of a new execution of the body of a while statement, the previously computed intermediate results have to be cleared.<sup>10</sup> For

<sup>9</sup> For  $C\#_{\mathcal{I}}$  alone it would be simpler to transfer control directly by updating *pos* to the value of a corresponding static function.

<sup>10</sup> `CLEARVALUES` is needed in the present rule formulation, which is close to an



the sake of brevity we skip the analogous transition rules for statements `do`, `for`, `switch`, `goto case`, `goto default`. Since we formulate the model for the human reader, we also use the `..`-notation, for example in the rules for abruption or for sequences of block statements. This avoids having to fuss with an explicit formulation of the context, typically determined by a walk through a list. This simplification, which is tailored for the human reader, can easily be resolved for an executable model version without increasing the number of rules.

```

EXECCSHARPSTMI ≡ match context(pos)
  ; → YIELD(Norm)
  exp; → pos := exp
  ▶ val; → YIELDUP(Norm)
  break; → YIELD(Break)
  continue; → YIELD(Continue)
  goto lab; → YIELD(Goto(lab))
  if (exp) stm1 else stm2 → pos := exp
  if (▶ val) stm1 else stm2 → if val then pos := stm1 else pos := stm2
  if (True) ▶ Norm else stm → YIELDUP(Norm)
  if (False) stm else ▶ Norm → YIELDUP(Norm)
  while (exp) stm → pos := exp
  while (▶ val) stm → if val then pos := stm
  else YIELDUP(Norm)
  while (True) ▶ Norm → {pos := up(pos), CLEARVALUES(up(pos))}
  while (True) ▶ Break → YIELDUP(Norm)
  while (True) ▶ Continue → {pos := up(pos), CLEARVALUES(up(pos))}
  while (True) ▶ abr → YIELDUP(abr)
  type loc; → YIELD(Norm)
  lab: stm → pos := stm
  lab: ▶ Norm → YIELDUP(Norm)
  checked block → pos := block
  checked ▶ Norm → YIELDUP(Norm)
  unchecked block → pos := block
  unchecked ▶ Norm → YIELDUP(Norm)
  ... ▶ abr ... → if up(pos) ≠ Undef ∧ PropagatesAbr(up(pos)) then
  YIELDUP(abr)

```

---

executable format. In a more abstract SOS-style, as used for the Java model in [1], it wouldn't be necessary because there the intermediate *values* can be written into a dynamic function for the still to be executed rest program.

{ }	→ YIELD( <i>Norm</i> )
{ <i>stm ...</i> }	→ <i>pos</i> := <i>stm</i>
{ ... ▶ <i>Norm</i> }	→ YIELDUP( <i>Norm</i> )
{ ... ▶ <i>Norm stm ...</i> }	→ <i>pos</i> := <i>stm</i>
{ ... ▶ <i>Goto</i> ( <i>l</i> ) ... }	→ <b>let</b> $\alpha = \text{GotoTarget}(\text{first}(\text{up}(\text{pos})), l)$ <b>if</b> $\alpha \neq \text{Undef}$ <b>then</b> { <i>pos</i> := $\alpha$ , CLEARVALUES( <i>up</i> ( <i>pos</i> )) } <b>else</b> YIELDUP( <i>Goto</i> ( <i>l</i> ))
{ ... ▶ <i>abr ...</i> }	→ YIELDUP( <i>abr</i> )

In  $C\#_{\mathcal{I}}$  abruptions are propagated upwards except at the following statements:

$\text{PropagatesAbr}(\alpha) \iff \text{label}(\alpha) \notin \{\text{Block}, \text{While}, \text{Do}, \text{For}, \text{Switch}\}$

To compute the target of a label in a list of block statements we define:

$\text{GotoTarget}(\alpha, l) =$   
  **if**  $\text{label}(\alpha) = \text{Lab}(l)$  **then**  $\alpha$   
  **elseif**  $\text{next}(\alpha) = \text{Undef}$  **then**  $\text{Undef}$   
  **else**  $\text{GotoTarget}(\text{next}(\alpha), l)$

The auxiliary macro CLEARVALUES( $\alpha$ ) to clear all values in the subtree at position  $\alpha$  can be defined by recursion as follows, proceeding from top to bottom and from left to right:<sup>11</sup>

CLEARVALUES( $\alpha$ )  $\equiv$   
   $\text{values}(\alpha) := \text{Undef}$   
  **if**  $\text{first}(\alpha) \neq \text{Undef}$  **then** CLEARVALUESSEQ( $\text{first}(\alpha)$ )

CLEARVALUESSEQ( $\alpha$ )  $\equiv$   
  CLEARVALUES( $\alpha$ )  
  **if**  $\text{next}(\alpha) \neq \text{Undef}$  **then** CLEARVALUESSEQ( $\text{next}(\alpha)$ )

### 3 $C\#_{\mathcal{C}}$ : refining $C\#_{\mathcal{I}}$ by static class features

In this section we refine the imperative core  $C\#_{\mathcal{I}}$  to  $C\#_{\mathcal{C}}$  by adding classes (modules) concentrating upon their static features (static fields, methods, constructors), including their initialization and the parameter mechanism that provides value, **ref** and **out** parameters. For such a refinement we a) extend the ASM

<sup>11</sup> Intuitively it should be clear that the execution of this recursive ASM provides simultaneously – in one step – the set of all updates of all its recursive calls, as is needed here for the clearing purpose; see [22] for a precise definition.

universes and functions, or introduce new ones, to reflect the grammar extensions for expressions and statements, b) add the appropriate constraints needed for the static analysis of the new items (type constraints, definite assignment rules), c) extend some of the macros, e.g. *PropagatesAbr*( $\alpha$ ), to make them work also for the newly occurring cases, d) add rules which define the semantics of the new instructions that operate over the new domains.

### 3.1 Static semantics of $C\sharp_C$

In  $C\sharp_C$  a program is a set of compilation units, each coming with “using directives” and declarations of names spaces (including a global namespace) and types (for classes and interfaces<sup>12</sup>) in the global namespace. For simplicity of exposition we disregard “using” directives and nested namespaces by assuming everywhere the adoption of (equivalent) fully qualified names. The precise syntax of classes and their static members, the rules for the accessibility of types and members via the access modifiers (public, internal, protected, private) and illustrating examples are spelt out in [19]. We define here the extension of the grammars for *Vexp*, *Sexp*, *Stm* and thereby of the corresponding ASM domains, which reflects the introduction of sets of *Classes* with static *Fields* and static *Methods* in  $C\sharp_C$ . The new set *Arg* of arguments appearing here reflects that besides value parameters also **ref** and **out** parameters can be used.

$$\begin{aligned}
 Vexp &::= \dots \mid Field \mid Class \text{ ‘.’ } Field \\
 Sexp &::= \dots \mid Meth ( [Args] ) \mid Class \text{ ‘.’ } Meth ( [Args] ) \\
 Arg &::= Exp \mid \text{ ‘ref’ } Vexp \mid \text{ ‘out’ } Vexp \\
 Args &::= Arg \{ \text{ ‘,’ } Arg \} \\
 Stm &::= \dots \mid \text{ ‘return’ } Exp \text{ ‘;’ } \mid \text{ ‘return’ } \text{ ‘;’ }
 \end{aligned}$$

The type constraints for the new expressions and the return statement are spelt out in [19]. The difference between **ref** and **out** parameters at function calls and in function bodies is reflected by including as *AddressPositions* all nodes whose parent node is labeled by **ref** or **out** and by adding the following definite assignment constraints:

- **ref** arguments must be definitely assigned *before* the function call.
- **out** arguments are definitely assigned *after* the function call.
- **ref** parameters and value parameters of a function are definitely assigned at the beginning of the function body.
- **out** parameters must be definitely assigned when the function returns.

<sup>12</sup>Note that struct and enum types and delegates are introduced by further refinement steps below.

Therefore the definite assignment constraints for expressions are extended by the following constraints for general argument expressions in function calls and for **ref** and **out** argument expressions:

- For  $exp = M(args)$ :
  - $before(args) = before(exp)$
  - $RefParams(args) \subseteq after(args)$
  - $after(exp) = after(args) \cup OutParams(args)$
- For  $exp = (\mathbf{ref} \ e)$  or  $exp = (\mathbf{out} \ e)$ :
  - $before(e) = before(exp)$
  - $after(exp) = after(e)$

The definite assignment constraints for statements are extended for function bodies and return statements as follows:

- If  $s$  is the body of  $M$ , then
  - $before(s) = ValueParams(M) \cup RefParams(M)$
- If  $stm = \mathbf{return};$  is in  $M$ , then
  - $OutParams(M) \subseteq before(stm)$
  - $after(stm) = vars(stm)$
- If  $stm = \mathbf{return} \ e;$  is in  $M$ , then
  - $before(e) = before(stm)$
  - $OutParams(M) \subseteq after(e)$
  - $after(stm) = vars(stm)$

The presence of to-be-initialized classes and of method calls is reflected by the introduction of new universes to denote methods, the initialization status of a type (which will be refined below by exceptions) and the sequence of still active method calls (frame stack):

$$Meth = Type \times Msig$$

$$TypeState = Linked \mid InProgress \mid Initialized$$

$$Frame = Meth \times Pos \times Locals \times Values$$

$$\text{where } Values = (Pos \rightarrow Result) \text{ and } Locals = (Loc \rightarrow Adr)$$

A method signature  $Msig$  consists of the name of a method plus the sequence of types of the arguments of the method. A method is uniquely determined by the type in which it is declared and its signature. The reasons for abruptions are extended by method return:

$$Abr = \dots \mid Return \mid Return(Value)$$

### 3.2 Dynamic semantics of $C\sharp_c$

To dynamically handle the addresses of static fields (global or class variables), the initialization state of types, the current method, the execution stack and the (initially) to be initialized type we use the following new dynamic functions:

$$\begin{aligned} \text{globals} &: \text{Type} \times \text{Field} \rightarrow \text{Adr} & \text{frames} &: \text{List}(\text{Frame}) \\ \text{typeState} &: \text{Type} \rightarrow \text{TypeState} & \text{meth} &: \text{Meth} \end{aligned}$$

We extend the stipulations for the initial state as follows:

- $\text{typeState}(c) = \text{Linked}$  for each class  $c$
- $\text{meth} = \text{EntryPoint::Main}()$  [*EntryPoint* is the main class]
- $\text{pos} = \text{body}(\text{meth})$  [The root position of the body]
- $\text{locals} = \text{values} = \emptyset$  and  $\text{frames} = []$

The submachine  $\text{EXEC}\text{SHARP}_C$  extends the machine  $\text{EXEC}\text{SHARP}_I$  for  $C\sharp_I$  by additional rules for the evaluation of the new expressions and for the execution of return statements. In the same way the further refinements in the sections below consist in the parallel addition of appropriate submachines.

$$\begin{aligned} \text{EXEC}\text{SHARP}_C &\equiv \\ &\text{EXEC}\text{SHARP}_I \\ &\text{EXEC}\text{SHARP}\text{EXP}_C \\ &\text{EXEC}\text{SHARP}\text{STM}_C \end{aligned}$$

**Expression evaluation rules.** The rules for class field evaluation in the submachine  $\text{EXEC}\text{SHARP}\text{EXP}_C$  are analogous to those for the evaluation of local variables in  $\text{EXEC}\text{SHARP}\text{EXP}_I$ , except for using *globals* instead of *locals* and for the additional clause for class initialization. The rules for method calls use the macro `INVOKESTATIC` defined below and reflect that the arguments are evaluated from left to right.<sup>13</sup>

$$\begin{aligned} \text{EXEC}\text{SHARP}\text{EXP}_C &\equiv \mathbf{match} \text{context}(\text{pos}) \\ c.f &\rightarrow \mathbf{if} \text{Initialized}(c) \mathbf{then} \text{YIELDINDIRECT}(\text{globals}(c::f)) \\ &\quad \mathbf{else} \text{INITIALIZE}(c) \end{aligned}$$

<sup>13</sup> These are the rules to be modified in case one wants to specify another evaluation order for expressions, involving the use of the ASM **choose** construct if some non-deterministic choice has to be formulated. For a discussion of such model variations we refer to [23] where an ASM model is developed which can be instantiated to capture the different expression evaluation strategies in Ada95, C, C++, Java,  $C\sharp$  and Fortran.

```

c.f = exp → pos := exp
c.f = ► val → if Initialized(c) then
    WRITEMEM(globals(c::f), type(c::f), val)
    YIELDUP(val)
    else INITIALIZE(c)

c.m(args) → pos := (args)
c.m►(vals) → INVOKESTATIC(c::m, vals)

ref vexp → pos := vexp
ref ►adr → YIELDUP(adr)

out vexp → pos := vexp
out ►adr → YIELDUP(adr)

() → YIELD([])
(arg, ...) → pos := arg
(val1, ..., ►valn) → YIELDUP([val1, ..., valn])
(... ►val, arg ...) → pos := arg

```

The macro `INVOKESTATIC` invokes the method if the initialization of its class is not triggered, otherwise it initializes the class. The initialization of a class (or struct, see Sect. 4) is not triggered if the class is already initialized.<sup>14</sup> For methods which are not declared external, `INVOKEMETHOD` updates the frame stack and the current frame in the expected way, allocating via `INITLOCALS` for every local variable or value parameter a new address and copying every value argument there. Since we will also have to deal with external methods – whose declaration includes an `extern` modifier and which may be implemented using a language other than  $C\sharp$  – we provide here for their invocation a sub-machine `INVOKEEXTERN`, to be defined separately depending on the class of external (e.g. library) methods.<sup>15</sup> The predicate *StaticCtor* recognizes static class constructors; their implicit call interrupts the member access at *pos*, to later return to the evaluation of *pos* instead of *up(pos)*. We separate the current frame—consisting of *meth*, *pos*, *locals* and *values*—from the stack of such frames to notationally smoothen the transition from  $C\sharp_{\mathcal{I}}$  to  $C\sharp_c$ .

```

INVOKESTATIC(c::m, vals) ≡
    if not TriggerInit(c) then INVOKEMETHOD(c::m, vals)
    else INITIALIZE(c)
    where TriggerInit(c) =  $\neg$ Initialized(c)

```

<sup>14</sup> As analyzed in [13], it is also not triggered if the class is marked with the implementation flag *beforefieldinit*, indicating that the reference of the static method does not trigger the class (or struct) initialization. If one wants to model this flag the definition has to be refined to  $\textit{TriggerInit}(c) = \neg\textit{Initialized}(c) \wedge \neg\textit{beforefieldinit}(c)$  and furthermore in Sect. 5.)

<sup>15</sup> For an illustration of this use of external methods see below the model for delegates.

```

INVOKEMETHOD( $c::m, vals$ )  $\equiv$ 
  if extern  $\in modifiers(c::m)$  then INVOKEEXTERN( $c::m, vals$ )
  else let  $p = \text{if } StaticCtor(c::m) \text{ then } pos \text{ else } up(pos)$  in
     $frames := push(frames, (meth, p, locals, values))$ 
     $meth := c::m$ 
     $pos := body(c::m)$ 
     $values := \emptyset$ 
  INITLOCALS( $c::m, vals$ )

```

The following definition for the initialization of local variables reflects that C $\sharp$  permits to pass function call parameters by value, as Java does, but also by **reference**. Also **out** parameters are allowed, treated as **ref** parameters except that they need not be definitely assigned until the function call returns.<sup>16</sup>

In the following definition, all (also simultaneous) applications of the external function *new* during the computation of the ASM are supposed to provide pairwise different fresh elements from the underlying domain *Adr*. See [24] and [18, 2.4.4] for a justification of this assumption. See also the end of Sect. 4 where we provide an abstract specification of the needed memory allocation to addresses of references and objects of struct type and to their instance fields.  $paramIndex(c::m, x)$  yields the index of the formal parameter  $x$  in the signature of  $c::m$ .

```

INITLOCALS( $c::m, vals$ )  $\equiv$ 
  forall  $x \in LocalVars(c::m)$  do // addresses for local variables
     $locals(x) := new(Adr, type(x))$ 
  forall  $x \in ValueParams(c::m)$  do // copy value arguments
    let  $adr = new(Adr, type(x))$  in
       $locals(x) := adr$ 
      WRITEMEM( $adr, type(x), vals(paramIndex(c::m, x))$ )
  forall  $x \in RefParams(c::m) \cup OutParams(c::m)$  do // ref, out arguments
     $locals(x) := vals(paramIndex(c::m, x))$ 

```

**Statement execution rules.** The rules for method return in the submachine EXEC $\sharp$ STM $_C$  trigger an abruptio upon returning from a method, resulting (via the propagation of this abruptio to the beginning of the method body where it occurred) in the execution of the machine EXITMETHOD. The rule to YIELDUP(*Norm*) does not capture falling off the method body, but yields up the result of the normal execution of the invocation of a method with void return type in an expression statement.

EXEC $\sharp$ STM $_C \equiv \text{match } context(pos)$

<sup>16</sup> To reflect different parameter passing mechanisms as encountered in other programming languages, due to the modular character of our model essentially only the above submachine INITLOCALS would have to be appropriately modified.

```

return exp; → pos := exp
return ▶ val; → YIELDUP(Return(val))
return; → YIELD(Return)
Return → if pos = body(meth) ∧ ¬Empty(frames) then
    EXITMETHOD(Norm)
Return(val) → if pos = body(meth) ∧ ¬Empty(frames) then
    EXITMETHOD(val)
▶ Norm; → YIELDUP(Norm)

```

The machine EXITMETHOD restores the frame of the invoker and passes the result value (if any). Upon normal return from a static constructor it also updates the *typeState* of the relevant class as *Initialized*. We also add a rule FREELOCALS to free the memory used for local variables and value parameters, using an abstract notion FREEMEMORY of how addresses of local variables and value parameters are actually de-allocated.<sup>17</sup>

```

EXITMETHOD(result) ≡
let (oldMeth, oldPos, oldLocals, oldValues) = top(frames) in
  meth := oldMeth
  pos := oldPos
  locals := oldLocals
  frames := pop(frames)
  if StaticCtor(meth) ∧ result = Norm then
    typeState(type(meth)) := Initialized
    values := oldValues
  else
    values := oldValues ⊕ {oldPos ↦ result}
  FREELOCALS

```

```

FREELOCALS ≡
forall x ∈ LocalVars(meth) ∪ ValueParams(meth) do
  FREEMEMORY(locals(x), type(x))

```

Following [2, §17.11,17.4.5.1,10.11,10.4.5.1] a type *c* is considered as initialized if its static constructor has been invoked (see the update of *typeState*(*c*) to *InProgress* in INITIALIZE below) or has terminated normally (see the update of *typeState*(*c*) to *Initialized* in EXITMETHOD above). We therefore define:

$$\textit{Initialized}(c) \iff \textit{typeState}(c) = \textit{Initialized} \vee \textit{typeState}(c) = \textit{InProgress}$$

To initialize a class its static constructor is invoked (*.cctor* = class constructor). This macro will be further refined in  $C\sharp_{\mathcal{E}}$  to account for exceptions during

---

<sup>17</sup> Under the assumption of a potentially infinite supply of addresses, which is often made when describing the semantics of a programming language, one can dispense with FREELOCALS.



an initialization.

```
INITIALIZE(c) ≡
  if typeState(c) = Linked then
    typeState(c) := InProgress
    forall f ∈ staticFields(c) do
      let t = type(c::f) in WRITEMEM(globals(c::f), t, defaultValue(t))
    INVOKEMETHOD(c::.ctor, [])
```

Note that in C $\sharp$  the initialization of a class does not trigger the initialization of its direct base class, differing on this point from Java where the rule for calling static constructors (see [1, Fig.4.5]) triggers the initialization of the superclass in case the superclass is not yet initialized.

With respect to the execution of initializers of static class fields the ECMA standard [2, §17.4.5.1] says that the static field initializers of a class correspond to a sequence of assignments that are executed in the textual order in which they appear in the class declaration. If a static constructor exists in the class, execution of the static field initializers occurs immediately prior to executing that static constructor. Otherwise, the static field initializers are executed at an *implementation-dependent* time prior to the first use of a static field of that class. Our model expresses the decision taken by Microsoft’s current C $\sharp$  compiler, which in the second case creates a static constructor. If one wants to reflect also the non-determinism suggested by the ECMA formulation, one can formalize the implementation-dependent external control by a monitored function *typeToBeInitialized* (which by the way can also be used for the classes and structs classified by an implementation flag as *beforefieldinit* type). The C $\sharp$  interpreter then takes the following form:<sup>18</sup>

```
if typeToBeInitialized ≠ Undef then
  INITIALIZE(typeToBeInitialized)
else EXECCSHARP
```

#### 4 Refinement C $\sharp_0$ of C $\sharp_c$ by object related features

In this section we refine the static class features of C $\sharp_c$  by adding objects (for class instances, comprising arrays and structs) together with *instance* fields, methods and constructors<sup>19</sup> as well as inheritance (including overriding

<sup>18</sup>This is discussed in detail in [13]. The reader finds there also a detection of further class initialization features that are missing in the ECMA specification, related to the definition of when a static class constructor has to be executed and to the initialization of structs.

<sup>19</sup>Destructors or finalizers which relate to garbage collection are not modeled here.

and overloading of methods). Accordingly we extend the ASM universes and functions of  $C\sharp_C$  to reflect the new expressions and statements together with the appropriate constraints and new rules, using appropriate refinements of some of the macros to define the semantics of the new instructions of  $C\sharp_O$ . For the detailed definition of the syntax of (members of) classes, interfaces, structs, etc., and of the constraints for the new modifiers (**abstract**, **sealed**, **readonly**, **volatile**, **virtual**, **override**) together with illustrating examples, we refer the reader to [19].

#### 4.1 Static semantics of $C\sharp_O$

The first extension concerns the sets  $Exp$ ,  $Vexp$ ,  $Sexp$  where the new reference and array types appear.  $Rank$  serves to denote the dimension of array types;  $NonArrayType$  stands for value types, classes and interfaces and will be extended in  $C\sharp_D$  to comprise also delegates. Value types represent a feature that distinguishes  $C\sharp$  from Java. In  $C\sharp_I$  we have cast expressions  $(t)exp$  where the type  $t$  and the type of  $exp$  are numeric types. Here, we extend the grammar to  $(t)exp$  where  $t$  and the type of  $exp$  can be any type. A  $RefExp$  is an expression of a reference type and an  $ArrayExp$  is an expression of an array type.

$$\begin{aligned}
 Exp & ::= \dots \mid \text{'null'} \mid \text{'this'} \mid \text{'typeof'} \text{'(' } RetType \text{' )'} \mid Exp \text{'is'} Type \\
 & \quad \mid Exp \text{'as'} RefType \mid \text{'(' } Type \text{' )'} Exp \mid Exp \text{'.'} Field \\
 & \quad \mid \text{'new'} NonArrayType \text{'[' } Exps \text{' ]'} \{ Rank \} [ArrayInitializer] \\
 Vexp & ::= \dots \mid Vexp \text{'.'} Field \mid RefExp \text{'.'} Field \mid \text{'base'} \text{'.'} Field \\
 & \quad \mid ArrayExp \text{'[' } Exps \text{' ]'} \\
 Sexp & ::= \dots \mid \text{'new'} Type ( [Args] ) \mid Exp \text{'.'} Meth ( [Args] ) \\
 & \quad \mid \text{'base'} \text{'.'} Meth ( [Args] ) \\
 Exps & ::= Exp \{ \text{' , ' } Exp \} \\
 Rank & ::= \text{'[' } \{ \text{' , ' } \} \text{' ]'}
 \end{aligned}$$

A **this** in an instance constructor or instance method of a struct is considered to be a  $Vexp$ . When a **this** occurs in a class it is not a  $Vexp$ .

The extended type classification where simple types become aliases for struct types is re-assumed by Fig. 4. We refer the reader to [19] for the detailed list of new type constraints. Also the constraints for overriding and overloading of methods and the resolution of overloaded methods at compile-time are spelt out there.

The subtype relation (i.e. the standard implicit conversion) is based on the inheritance relation – defined as a finite tree with root **object** – together with the “implements” relation between classes and interfaces. It is defined as

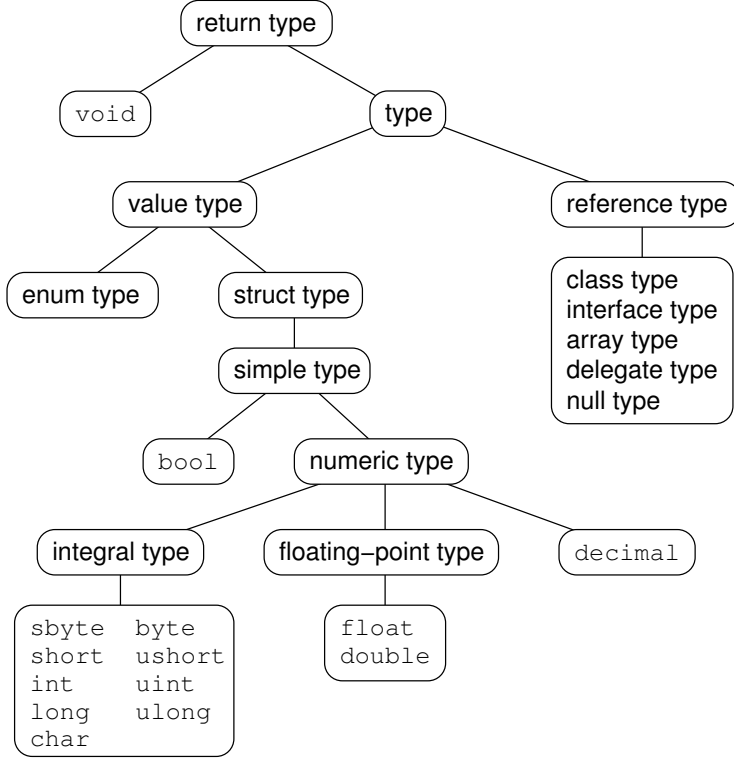


Fig. 4. The classification of types of C#.

follows (and should not be confused with the classification of types in Fig. 4):

- $T$  any type  $\implies T \preceq \mathbf{object}$  and  $T \preceq T$
- class  $S$  derived from  $T \implies S \preceq T$
- class, interface or struct  $S$  implements interface  $T \implies S \preceq T$
- $T$  array type  $\implies T \preceq \mathbf{System.Array}$
- $T$  delegate type  $\implies T \preceq \mathbf{System.Delegate}$
- $T$  value type  $\implies T \preceq \mathbf{System.ValueType}$
- $T$  array or delegate type  $\implies T \preceq \mathbf{System.ICloneable}$
- $T$  reference type  $\implies \Lambda \preceq T$  [ $\Lambda$  is the null type]
- $S$  and  $T$  reference types,  $S \preceq T \implies S[R_1] \cdots [R_k] \preceq T[R_1] \cdots [R_k]$

Note that types of one category in Fig. 4 can be subtypes of another (disjoint) category. For example, if a struct type  $S$  implements the interface  $I$ , then (the value type)  $S$  is considered to be a subtype of (the reference type)  $I$ .

We list here the additional definite assignment rules for local variables of struct type:

- If  $p$  is a local variable of a struct type  $S$ , then  $p.f$  is considered as a local variable for each instance field  $f$  of  $S$ .
- A local variable  $p$  of struct type  $S$  is definitely assigned  $\iff p.f$  is definitely assigned for each instance field  $f$  of  $S$ .

We assume that as a result of field and method resolution the attributed syntax tree has exact information. Each field access has the form  $T::f$  where  $f$  is a field declared in the type  $T$ . Each method call has the form  $T::m(args)$  where  $m$  is the signature of a method declared in type  $T$ . Moreover, at compile-time certain expressions are reduced to basic expressions as follows.

For the base access of fields and methods we have:

- $\mathbf{base}.f$  in class  $C$  is replaced by  $\mathbf{this}.B::f$ , where  $B$  is the first base class of  $C$  where the field  $f$  is declared.
- $\mathbf{base}.m(args)$  in class  $C$  is replaced by  $\mathbf{this}.B::M(args)$ , where  $B::M$  is the method signature of the method selected by the compiler (the set of applicable methods is constructed starting in the *base class* of  $C$ ). This selection algorithm is described in [19], formalizing the conditions stated in [2, §14.4.2/3].

For instance field access and class instance creation we have:

- If  $f$  is a field, then  $f$  is replaced by  $\mathbf{this}.T::f$ , where  $f$  is declared in  $T$ .
- Let  $T$  be a class type. Then the instance creation  $\mathbf{new } T::M(args)$  is replaced by  $\mathbf{new } T.T::M(args)$ .

Hence we split an instance creation expression into a creation part and an invocation of an instance constructor. To make the splitting correctly reflect the intended meaning of  $\mathbf{new } T::M(args)$ , we assume in our model that class instance constructors return the value of  $\mathbf{this}$ .

For instance constructors of structs one has to reflect that in addition they need an address for  $\mathbf{this}$ . Also for constructors of structs we assume that they return the value of  $\mathbf{this}$ . Let  $S$  be a struct type. Then the following transformations are applied:

- An assignment  $vexp = \mathbf{new } S::M(args)$  is replaced by  $vexp.S::M(args)$ . This reflects that such a  $\mathbf{new}$  triggers no object creation or memory allocation since structs get their memory allocated at declaration time.
- Other occurrences of  $\mathbf{new } S::M(args)$  are replaced by  $x.S::M(args)$ , where  $x$  is a new temporary local variable of type  $S$ .

For automatic boxing we have:

- An assignment  $vexp = exp$  is replaced by  $vexp = (T)exp$  if  $type(exp)$  is a value type,  $T = type(vexp)$  and  $T$  is a reference type. In this case we must have  $type(exp) \preceq T$ .
- An argument  $arg$  is replaced by  $(T)arg$  if  $type(arg)$  is a value type, the selected method expects an argument of type  $T$  and  $T$  is a reference type. In this case we must have  $type(arg) \preceq T$ .

## 4.2 Dynamic semantics for $C\sharp_{\mathcal{O}}$

We are now ready to describe the extension of the dynamic state for the model of  $C\sharp_{\mathcal{O}}$ . The domain of values is extended to contain also references (assuming  $Ref \cap Adr = \emptyset$  and  $null \in Ref$ ) and struct values:  $Value = SimpleValue \cup Adr \cup Ref \cup Struct$ . The set  $Struct$  of struct values can be defined as the set of mappings from  $StructType::Field$  to  $Value$ . The value of an instance field of a value of struct type  $T$  can then be extracted by applying the map to the field name, i.e.  $structField(val, T, f)$ .

Two dynamic functions keep track of the  $runTimeType: Ref \rightarrow Type$  of references and of the type object  $typeObj: RetType \rightarrow Ref$  of a given type, where  $RetType ::= Type \mid \text{'void'}$ . The memory function is extended to store also references:  $mem: Adr \rightarrow SimpleValue \cup Ref \cup \{Undef\}$ . For boxing we need a dynamic function  $valueAdr: Ref \rightarrow Adr$  to record the address of a value in a box. If  $runTimeType(ref)$  is a *value type*  $t$ , then  $valueAdr(ref)$  is the address of the struct value of type  $t$  stored in the box. The static function  $instanceFields: Type \rightarrow Powerset(Type::Field)$  yields the set of instance fields of any given type  $t$ ; if  $t$  is a class type, it includes the fields declared in base classes of  $t$ . We use the common programming notation  $Type::Field$  instead of the set-theoretic product set notation. We abstract from the implementation-dependent layout of structs and objects and use a function  $fieldAdr: (Adr \cup Ref) \times Type::Field \rightarrow Adr$  to record addresses of fields. This function satisfies the following properties:

- If  $t$  is a *struct type*, then  $fieldAdr(adr, t::f)$  is the address of field  $f$  of a value of type  $t$  stored in  $mem$  at address  $adr$ .
- A value of struct type  $t$  at address  $adr$  occupies the following addresses in  $mem$ :

$$\{fieldAdr(adr, f) \mid f \in instanceFields(t)\}$$

- If  $runTimeType(ref)$  is a *class type*, then  $fieldAdr(ref, t::f)$  is the address of field  $t::f$  of the object referenced by  $ref$ .
- An object of class  $c$  is represented by a reference  $ref$  with the property  $runTimeType(ref) = c$  and occupies the following addresses in  $mem$ :

$$\{fieldAdr(ref, f) \mid f \in instanceFields(c)\}$$

A function  $elemAdr: Ref \times \mathbb{N}^* \rightarrow Adr$  records addresses of array elements. The **this** reference is treated as first parameter and is passed by value. Therefore  $paramIndex(c::m, \mathbf{this}) = 0$  and **this** is an element of  $ValueParams(c::m)$ .

For the refinement of the  $EXEC\sharp_{\mathcal{O}}$  transition rules it suffices to add the new rule  $EXEC\sharp_{\mathcal{O}}EXP_{\mathcal{O}}$  for evaluating the new expressions, since  $C\sharp_{\mathcal{O}}$

introduces no new statements.

$$\begin{aligned} \text{EXEC SHARP}_O &\equiv \\ &\text{EXEC SHARP}_C \\ &\text{EXEC SHARP EXP}_O \end{aligned}$$

For better readability we organize the numerous  $\text{EXEC SHARP EXP}_O$  rules for each of the new expressions into parallel submachines each of which collects the rules for expressions which belong to the same category (for type testing and casting, for fields, for arrays). As analyzed in [13], the invocation of an instance constructor of a class may trigger the class initialization.

$$\begin{aligned} \text{EXEC SHARP EXP}_O &\equiv \mathbf{match} \text{ context}(pos) \\ &\mathbf{null} \rightarrow \text{YIELD}(null) \\ &\mathbf{this} \rightarrow \text{YIELD INDIRECT}(locals(\mathbf{this})) \\ &\text{TEST CAST EXP}_O \\ &\text{FIELD EXP}_O \\ &\mathbf{new } c \rightarrow \mathbf{let } ref = \mathbf{new}(Ref, c) \mathbf{in} \\ &\quad \text{runTimeType}(ref) := c \\ &\quad \mathbf{forall } f \in \text{instanceFields}(c) \mathbf{do} \\ &\quad \quad \mathbf{let } adr = \text{fieldAdr}(ref, f) \mathbf{and } t = \text{type}(f) \mathbf{in} \\ &\quad \quad \quad \text{WRITE MEM}(adr, t, \text{defaultValue}(t)) \\ &\quad \quad \text{YIELD}(ref) \\ &exp.T::M(args) \rightarrow pos := exp \\ &\blacktriangleright val.T::M(args) \rightarrow \mathbf{if } \text{StructValueInvocation}(up(pos)) \mathbf{then} \\ &\quad // \text{create home for struct value} \\ &\quad \mathbf{let } adr = \mathbf{new}(Adr, \text{type}(pos)) \mathbf{in} \\ &\quad \quad \text{WRITE MEM}(adr, \text{type}(pos), val) \\ &\quad \quad \text{values}(pos) := adr \\ &\quad pos := (args) \\ &val.T::M\blacktriangleright(vals) \rightarrow \mathbf{if } \text{InstanceCtor}(M) \wedge \text{TriggerInit}(T) \mathbf{then} \\ &\quad \text{INITIALIZE}(T) \\ &\quad \mathbf{elseif } val \neq null \mathbf{then} \\ &\quad \quad \text{INVOKE INSTANCE}(T::M, val, vals) \\ &\text{ARRAY EXP}_O \end{aligned}$$

A struct value invocation is a method invocation on a struct value which is not stored in a variable. For such struct values the above rule creates a temporary storage area (called “home”) to be passed in the invocation as value of **this**.

$$\begin{aligned} \text{StructValueInvocation}(exp.T::M(args)) &\iff \\ &\text{type}(exp) \in \text{StructType} \wedge exp \notin Vexp \end{aligned}$$

The rules for casting in  $\text{TEST CAST EXP}_O$  use the new macro  $\text{YIELD UP BOX}$  defined below. Note that in expressions ‘ $exp \text{ is } t$ ’ and  $(t)exp$  the type  $t$  can

be any type, whereas in ‘*exp as t*’ the type *t* must be a reference type. The type of ‘*exp is t*’ is `bool`, while the type of  $(t) \text{ exp}$  and ‘*exp as t*’ is *t*.

```

TESTCASTEXPo ≡
  typeof(t) → YIELD(typeObj(t))
  exp is t → pos := exp
  ▶ val is t → if type(pos) ∈ ValueType then
    YIELDUP(type(pos) ≼ t) // compile-time property
  else
    YIELDUP(val ≠ null ∧ runTimeType(val) ≼ t)
  exp as t → pos := exp
  ▶ val as t → if type(pos) ∈ ValueType then
    YIELDUPBOX(type(pos), val) // box a copy of the value
  elseif (val ≠ null ∧ runTimeType(val) ≼ t) then
    YIELDUP(val) // pass reference through
  else YIELDUP(null) // convert to null reference
  (t) exp → pos := exp
  (t) ▶ val → if type(pos) ∈ ValueType then
    // compile-time identity
    if t = type(pos) then YIELDUP(val)
    // box value
    if t ∈ RefType then YIELDUPBOX(type(pos), val)
  if type(pos) ∈ RefType then
    if t ∈ RefType ∧ (val = null ∨ runTimeType(val) ≼ t) then
      YIELDUP(val) // pass reference through
    if t ∈ ValueType ∧ val ≠ null ∧ t = runTimeType(val) then
      // un-box a copy of the value
      YIELDUP(memValue(valueAdr(val), t))

```

The rules for instance field access and assignment in `FIELDEXPo` are analogous to those for class variables, adding the evaluation of the instance, using *fieldAdr* instead of *globals*, and instead of `WRITEMEM` the macro `SETFIELD` defined below. The second rule for instance field access has to distinguish two cases, depending upon the statically known instance type. Since this type information is already known at the time of static analysis, it could be resolved by introducing two separate constructs for field access, as one of our referees observed pointing also out that in fact the CLI has a single, overloaded instruction for field access with overloading to be resolved by the JIT. However from the modeling point of view, having two separate constructs for field access would lead to essentially the same two rules we have formulated here, except for having as rule guard a matching condition for the two constructs instead of the type test. We use  $\text{type}(\text{exp}.t::f) = \text{type}(t::f)$ .

```

FIELDEXPo ≡

```

```

exp . t::f → pos := exp
▶ val . t::f → if type(pos) ∈ ValueType ∧ val ∉ Adr then
    YIELDUP(structField(val, type(pos), t::f))
    elseif val ≠ null then
        YIELDUPINDIRECT(fieldAdr(val, t::f))

exp1 . t::f = exp2 → pos := exp1
▶ val . t::f = exp → pos := exp
val1 . t::f = ▶ val2 → if val1 ≠ null then
    SETFIELD(val1, t::f, val2)
    YIELDUP(val2)

```

$C\#_O$  supports single dimensional as well as multi-dimensional arrays. Array types are read from right to left. For example, `int[][] , ]` is the type of single-dimensional arrays of two-dimensional arrays with elements of type `int`. By  $dim(n)$  we denote a sequence of  $n - 1$  commas, hence  $T[dim(n)]$  is the type of  $n$ -dimensional arrays with elements of type  $T$ . The length of the  $i$ th dimension of an  $n$ -dimensional array represented by a reference  $ref$  is stored as the value of  $dimLength(ref, i)$ . Note that the rules for using array indexing expressions as `rvalue` respectively as `lvalue` appear together as subgroups of `ARRAYEXPO`, separated by pattern matching.

```

ARRAYEXPO ≡
new T [exp1 , ... , expn] [R1] ⋯ [Rk] → pos := exp1
new T [l1 , ... , ▶ li , expi+1 , ... , expn] [R1] ⋯ [Rk] → pos := expi+1
new T [l1 , ... , ▶ ln] [R1] ⋯ [Rk] →
if ∀ i ∈ [1 .. n] (0 ≤ li) then
    let S = T [R1] ⋯ [Rk] in
        let ref = new(Ref, [l1 , ... , ln], S) in
            runTimeType(ref) := T [dim(n)] [R1] ⋯ [Rk]
            forall i ∈ [1 .. n] do dimLength(ref, i - 1) := li
            forall α ∈ [0 .. l1 - 1] × ⋯ × [0 .. ln - 1] do
                WRITEMEM(elemAdr(ref, α), S, defaultValue(S))
            YIELDUP(ref)

exp0 [exp1 , ... , expn] → pos := exp0
▶ ref [exp1 , ... , expn] → pos := exp1
ref [i1 , ... , ▶ ik , expk+1 , ... , expn] → pos := expk+1
ref [i1 , ... , ▶ in] →
if ref ≠ null ∧ ∀ k ∈ [1 .. n] (0 ≤ ik < dimLength(ref, k - 1)) ∧
    (RefOrOutArg(up(pos)) ∧ type(up(pos)) ∈ RefType →
        elementType(runTimeType(ref)) = type(up(pos)))
then
    YIELDUPINDIRECT(elemAdr(ref, (i1 , ... , in)))

exp0 [exp1 , ... , expn] = exp → pos := exp0
▶ ref [exp1 , ... , expn] = exp → pos := exp1

```



```

ref [i1, ..., ▶ ik, expk+1, ..., expn] = exp → pos := expk+1
ref [i1, ..., ▶ in] = exp → pos := exp
ref [i1, ..., in] = ▶ val →
  let T = elementType(runTimeType(ref)) in
  if ref ≠ null ∧ ∀k ∈ [1..n] (0 ≤ ik < dimLength(ref, k - 1)) ∧
    (type(pos) ∈ RefType → runTimeType(val) ≼ T)
  then
    WRITEMEM(elemAdr(ref, (i1, ..., in)), T, val)
    YIELDUP(val)

```

**Refinement of macros.** Invocation of instance methods splits into virtual and non-virtual calls. The function *lookup* yields the class where the given method specification is defined in the class hierarchy, depending on the runtime type of the given reference.

```

INVOKEINSTANCE(T::M, val, vals) ≡
  if callKind(T::M) = Virtual then // indirect call, val ∈ Ref
    let S = lookup(runTimeType(val), T::M) in
    let this = if S ∈ StructType then valueAdr(val) else val in
    INVOKEMETHOD(S::M, [this] · vals)
  if callKind(T::M) = NonVirtual then // direct call, val ∈ Adr ∪ Ref
    INVOKEMETHOD(T::M, [val] · vals)

```

In  $C\sharp_{\mathcal{O}}$  the notion of reading from the memory is refined by extending the simple equation  $memValue(adr, t) = mem(adr)$  of  $C\sharp_{\mathcal{I}}$  to fit also reference and struct types. This is done by the following simultaneous recursive definition of *memValue* and *getField* along the given struct type.

```

memValue(adr, t) =
  if t ∈ SimpleType ∪ RefType then mem(adr)
  elseif t ∈ StructType then {f ↦ getField(adr, f) | f ∈ instanceFields(t)}

getField(adr, t::f) = memValue(fieldAdr(adr, t::f), type(t::f))

```

Writing to memory is refined recursively together with **SETFIELD** along the given struct type:

```

WRITEMEM(adr, t, val) ≡
  if t ∈ SimpleType ∪ RefType then mem(adr) := val
  elseif t ∈ StructType then
    forall f ∈ instanceFields(t) do SETFIELD(adr, f, val(f))

```

```

SETFIELD(adr, t::f, val) ≡ WRITEMEM(fieldAdr(adr, t::f), type(t::f), val)

```

The notion of *AddressPos* from  $C\sharp_{\mathcal{C}}$  is refined to include also lvalue nodes of *StructType*, so that address positions are of the following form: **ref** □, **out** □,

$\square ++, \square --, \square op = exp, \square .f, \square .m(args).$

$AddressPos(\alpha) \iff FirstChild(\alpha) \wedge$   
 $label(up(\alpha)) \in \{\mathbf{ref}, \mathbf{out}, ++, --\} \vee label(up(\alpha)) \in Aop \vee$   
 $(label(up(\alpha)) = '.' \wedge \alpha \in Vexp \wedge type(\alpha) \in StructType)$

YIELDUPBOX creates a box for a given value of a given type and returns its reference. The run-time type of a reference to a boxed value of struct type  $t$  is defined to be (the value type)  $t$  of the value. There is no need to introduce special reference types for boxed values. If  $type(exp)$  is a value type that implements the interface  $I$ , then  $type(exp) \preceq I$  and the value can be boxed using ' $(I) exp$ ' or ' $exp \text{ as } i$ '.

$YIELDUPBOX(t, val) \equiv \mathbf{let} \text{ } ref = new(Ref) \mathbf{and} \text{ } adr = new(Adr, t) \mathbf{in}$   
 $runTimeType(ref) := t$   
 $valueAdr(ref) := adr$   
 $WRITEMEM(adr, t, val)$   
 $YIELDUP(ref)$

The struct value is copied in both cases, when it is boxed and when it is un-boxed.

**ASM function new.** We now justify in the context of the basic parallel execution mechanism of ASM rules the sequentiality which is used in the following macros:

$\mathbf{let} \text{ } adr = new(Adr, T) \mathbf{in} P$   
 $\mathbf{let} \text{ } ref = new(Ref, T) \mathbf{in} P$   
 $\mathbf{let} \text{ } ref = new(Ref, [l_1, \dots, l_n], T) \mathbf{in} P$

In the context of the machine EXEC SHARP this comes up to specify an abstract memory management. In fact  $\mathbf{let} \text{ } adr = new(Adr, T) \mathbf{in} P$  stands for the sequential execution of a new address allocation (which uses the ASM construct **import** to provide a completely fresh element) followed by  $P$ :

$\mathbf{let} \text{ } adr = new(Adr, T) \mathbf{in} P \equiv (\mathbf{import} \text{ } adr \mathbf{do} \text{ } ALLOCADR(adr, T)) \mathbf{seq} P$

where the operator **seq** for sequential execution of two ASMs  $M, N$  is to be understood as defined for turbo ASMs in [25] (alternatively see [18, Ch.4.1]), namely as binding into one overall ASM step the two steps of first executing  $M$  in the given state and then  $N$  in the resulting state. Similarly  $\mathbf{let} \text{ } ref = new(Ref, T) \mathbf{in} P$  stands for the sequential execution of address allocation for all instance fields of a given type followed by  $P$ :

$\mathbf{let} \text{ } ref = new(Ref, T) \mathbf{in} P \equiv$   
 $\mathbf{import} \text{ } ref \mathbf{do}$

```

    Ref(ref) := True
    ALLOCFIELDS(ref, instanceFields(T))
  seq P

```

Similarly we define the address allocation for elements of an  $n$ -dimensional array:

```

let ref = new(Ref, [l1, ..., ln], T) in P ≡
  import ref do
    Ref(ref) := True
    forall α ∈ [0 .. l1 - 1] × ... × [0 .. ln - 1] do
      import adr do
        elemAdr(ref, α) := adr
        ALLOCADR(adr, T)
  seq P

```

The two macros for allocation of addresses and fields can be recursively defined as follows, relying again upon the definition of recursive turbo ASMs in [22] (or see alternatively [18, Ch.4.1.2]):

```

ALLOCADR(adr, T) ≡
  Adr(adr) := True
  if T ∈ StructType then ALLOCFIELDS(adr, instanceFields(T))

```

```

ALLOCFIELDS(x, fs) ≡
  forall f ∈ fs import adr do
    fieldAdr(x, f) := adr
    ALLOCADR(adr, type(f))

```

## 5 Refinement $C_{\#}^{\varepsilon}$ of $C_{\#}^{\circ}$ by exception handling

In this section we extend  $C_{\#}^{\circ}$  with the exception handling mechanism of  $C_{\#}^{\#}$ , which separates normal program code from exception handling code. To this purpose exceptions are represented as objects of predefined system exception classes or of user-defined application exception classes. Once created (thrown), these objects trigger an abruptio of the normal program execution to catch the exception – in case it is compatible with one of the exception classes appearing in the program in an enclosing try-catch-finally statement. Optional finally statements are guaranteed to be executed independently of whether the try statement completes normally or is abruptly.

## 5.1 Static semantics of $C\sharp_{\mathcal{E}}$

For the refinement of  $\text{EXEC}\text{SHARP}_O$  by exceptions, as in the previous section it suffices to add the new rules for exception handling and to extend the static semantics. The set of statements is extended by `throw` and `try-catch-finally` statements satisfying the following constraints:

$$\begin{aligned} \text{Stm} & ::= \dots \mid \text{'throw' } \text{Exp } \text{' ;' } \mid \text{'throw' } \text{' ;' } \\ & \quad \mid \text{'try' } \text{Block } \{ \text{Catch} \} [ \text{'catch' } \text{Block} ] [ \text{'finally' } \text{Block} ] \\ \text{Catch} & ::= \text{'catch' } \text{' (' } \text{ClassType } [ \text{Loc} ] \text{' )' } \text{Block} \end{aligned}$$

- every `try-catch-finally` statement contains at least one *catch clause*, *general catch clause* (i.e. of form `catch block`), or *finally block*
- no `return` statements are allowed in finally blocks
- a `break`, `continue`, or `goto` statement is not allowed to jump out of a finally block
- a `throw` statement without expression is only allowed in catch blocks
- the exception classes in a *Catch* clause appear there in a non-decreasing type order, more precisely, for every try-catch statement of the form

$$\text{try block catch } (E_1 x_1) \text{ block}_1 \dots \text{catch } (E_n x_n) \text{ block}_n$$

the following holds  $i < j \implies E_j \not\preceq E_i$  (and  $E_i \preceq \text{System.Exception}$ ).

In our model the sets of abruptions and type states have to be extended by exceptions. Due to the presence of `throw` statements without expression, a stack of references is needed to record exceptions which are to be re-thrown.

$$\text{Abr} = \dots \mid \text{Exc}(\text{Ref}), \text{TypeState} = \dots \mid \text{Exc}(\text{Ref}), \text{excStack}: \text{List}(\text{Ref})$$

To simplify the exposition we assume that general catch clauses `'catch block'` are replaced at compile-time by `'catch (Object x) block'` with a new variable  $x$ . We also reduce `try-catch-finally` statements to `try-catch` and `try-finally` statements as follows. Both reductions can easily be shown to correctly express the ECMA specification.

$$\begin{array}{ccc} \begin{array}{l} \text{try TryBlock} \\ \text{catch } (E_1 x_1) \text{ CatchBlock}_1 \\ \vdots \\ \text{catch } (E_n x_n) \text{ CatchBlock}_n \\ \text{finally FinallyBlock} \end{array} & \implies & \begin{array}{l} \text{try } \{ \\ \quad \text{try TryBlock} \\ \quad \text{catch } (E_1 x_1) \text{ CatchBlock}_1 \\ \quad \vdots \\ \quad \text{catch } (E_n x_n) \text{ CatchBlock}_n \\ \} \text{finally FinallyBlock} \end{array} \end{array}$$

If a static constructor throws an exception, and no catch clauses exists to catch it, then this exception is wrapped into a `TypeInitializationException` by translating `static T() { BlockStatements }` into

```
static T() {
  try { BlockStatements }
  catch (Exception e) {
    throw new TypeInitializationException(T, e);
  }
}
```

The reachability rules and the definite assignment constraints for a try-catch  $stm \equiv \text{try } tryBlock \text{ catch } (E_1 x_1) catchBlock_1 \dots \text{catch } (E_n x_n) catchBlock_n$  are:

- If  $reachable(stm)$ , then  $reachable(tryBlock)$  and  $reachable(catchBlock_i)$  for every  $i \in [1..n]$ .
- If  $normal(tryBlock)$  or  $normal(catchBlock)$  for at least one  $i \in [1..n]$ , then  $normal(stm)$ .
- $before(tryBlock) = before(stm)$
- $before(catchBlock_i) = before(stm) \cup \{x_i\}$  for every  $i \in [1..n]$
- $after(stm) = after(tryBlock) \cap \bigcap_{i=1}^n after(catchBlock_i)$

For a statement  $stm$  of the form `try tryBlock finally finallyBlock` the rules and constraints are:

- If  $reachable(stm)$ , then  $reachable(tryBlock)$  and  $reachable(finallyBlock)$ .
- If  $normal(tryBlock)$  and  $normal(finallyBlock)$ , then  $normal(stm)$ .
- $before(tryBlock) = before(stm)$
- $before(finallyBlock) = before(stm)$
- $after(stm) = after(tryBlock) \cup after(finallyBlock)$

## 5.2 Dynamic semantics for $C\sharp_\varepsilon$

The transition rules for  $EXEC\text{SHARP}_E$  are defined by adding two submachines to  $EXEC\text{SHARP}_O$ . The first one provides the rules for handling the exceptions which may occur during the evaluation of expressions, the second one describes the meaning of the new throw and try-catch-finally statements.

$$\text{EXEC}\text{SHARP}_E \equiv$$

$$\text{EXEC}\text{SHARP}_O$$

$$\text{EXEC}\text{SHARP}\text{EXP}_E$$

## EXEC SHARP STM<sub>E</sub>

**Expression evaluation rules.** EXEC SHARP EXP<sub>E</sub> contains rules for each of the numerous forms of run-time exceptions defined in the subclasses of `System.Exception`. We give here seven characteristic examples and group them for the ease of presentation into parallel submachines by the form of expression they are related to, namely for arithmetical exceptions and for those related to cast expressions, reference expressions or array expressions. The notion of FAILUP we use is supposed to execute the code `throw new E()` at the parent position, so that we define the macro by invoking an internal method `ThrowE` with that effect for each of the exception classes  $E$  used as parameter of FAILUP.

```

EXEC SHARP EXPE ≡ match context(pos)
  uop ▶ val → if Checked(pos) ∧ Overflow(uop, val) then
    FAILUP(OverflowException)
  val1 bop ▶ val2 → if DivisionByZero(bop, val2) then
    FAILUP(DivideByZeroException)
    elseif DecimalOverflow(bop, val1, val2) ∨
      (Checked(pos) ∧ Overflow(bop, val1, val2))
    then FAILUP(OverflowException)

CAST EXCEPTIONS
NULL REF EXCEPTIONS
ARRAY EXCEPTIONS

```

```

FAILUP(E) ≡ INVOKEMETHOD(ExcSupport::ThrowE, [ ])

```

```

CAST EXCEPTIONS ≡ match context(pos)
  (t) ▶ val →
    if type(pos) ∈ RefType then
      if t ∈ RefType ∧ val ≠ Null ∧ runTimeType(val) ≠ t then
        FAILUP(InvalidCastException)
      if t ∈ ValueType then // attempt to unbox
        if val = Null then FAILUP(NullReferenceException)
        elseif t ≠ runTimeType(val) then
          FAILUP(InvalidCastException)
      if type(pos) ∈ NumericType ∧ t ∈ NumericType ∧ Checked(pos) ∧
        Overflow(t, val)
    then FAILUP(OverflowException)

```

The semantics of assignments as defined by the ECMA standard and formalized by the rule NULL REF EXCEPTIONS is violated by a compiler optimization in [7] related to the timing of the *Null* check for certain expressions, see the

analysis in [13].

```

NULLREFERENCEEXCEPTIONS  $\equiv$  match context(pos)
   $\blacktriangleright$  ref.t::f  $\rightarrow$  if ref = Null then
    FAILUP(NullReferenceException)
  ref.t::f =  $\blacktriangleright$  val  $\rightarrow$  if ref = Null then
    FAILUP(NullReferenceException)
  ref.T::M( $\blacktriangleright$  vals)  $\rightarrow$  if ref = Null then
    FAILUP(NullReferenceException)

```

If the address of an array element is passed as a *ref* or *out* argument to a method, then the run-time element type of the array must be *equal* to the parameter type that the method expects. If an object is assigned to an array element, then the type of the object must be a *subtype* of run-time element type of the array (array covariance problem). In both cases, if the condition is not satisfied, an `ArrayTypeMismatchException` is thrown.

```

ARRAYEXCEPTIONS  $\equiv$  match context(pos)
  new T[l1, ...,  $\blacktriangleright$ ln][R1]  $\cdots$  [Rk]  $\rightarrow$ 
    if  $\exists i \in [1..n]$  (ik < 0) then FAILUP(OverflowException)
  ref[i1, ...,  $\blacktriangleright$ in]  $\rightarrow$ 
    if ref = Null then FAILUP(NullReferenceException)
    elseif  $\exists k \in [1..n]$  (ik < 0  $\vee$  dimLength(ref, k - 1)  $\leq$  ik) then
      FAILUP(IndexOutOfRangeException)
    elseif RefOrOutArg(up(pos))  $\wedge$  type(up(pos))  $\in$  RefType  $\wedge$ 
      elementType(runTimeType(ref))  $\neq$  type(up(pos))
    then FAILUP(ArrayTypeMismatchException)
  ref[i1, ..., in] =  $\blacktriangleright$  val  $\rightarrow$ 
    if ref = Null then FAILUP(NullReferenceException)
    elseif  $\exists k \in [1..n]$  (ik < 0  $\vee$  dimLength(ref, k - 1)  $\leq$  ik) then
      FAILUP(IndexOutOfRangeException)
    elseif type(pos)  $\in$  RefType  $\wedge$ 
      runTimeType(val)  $\not\leq$  elementType(runTimeType(ref)) then
      FAILUP(ArrayTypeMismatchException)

```

**Statement execution rules.** The statement execution submachine splits naturally into submachines for `throw`, `try-catch`, `try-finally` statements and a rule for the propagation of an exception (from the root position of a method body) to the method caller. To support a correct understanding of the exception messages that are printed to the console we include into the rule for `throw` statements the initialization of stack traces. The initialization of stack traces in Java and  $C\sharp$  is different. In Java, the stack trace is initialized with the complete trace up to the main function once and for all when the exception object is created. In  $C\sharp$  the stack trace is initialized with the empty trace each time when the exception object is thrown with `throw` and then augmented whenever the exception is propagated to a parent frame. The semantics of

the parameterless **throw**; instruction is explained in terms of the exception Stack *excStack*. When an exception is caught, it is pushed on top of the exception stack (which as explained above is needed to record exceptions which are to be re-thrown). Whenever a catch block terminates (normally or abruptly) the topmost element of the exception stack is deleted. No special rules are needed for general catch clauses ‘**catch** *block*’ in try-catch statements, due to their compile-time transformation mentioned above. The completeness of the try-finally rules is due to the constraints listed above, which restrict the possibilities for exiting a finally block to normal completion or triggering an exception.

```

EXEC SHARPSTME ≡ match context(pos)
  throw exp; → pos := exp
  throw ▶ ref; → if ref = Null then FAILUP(NullReferenceException)
                    else INITSTACKTRACE(ref, meth)
                      YIELDUP(Exc(ref))
  throw; → YIELD(Exc(top(excStack)))
  try block catch (E x) stm ... → pos := block
  try ▶ Norm catch (E x) stm ... → YIELDUP(Norm)
  try ▶ Exc(ref) catch (E1 x1) stm1 ... catch (En xn) stmn →
    if ∃i ∈ [1 .. n] runTimeType(ref) ≼ Ei then
      let j = min{i ∈ [1 .. n] | runTimeType(ref) ≼ Ei} in
        pos := stmj
        excStack := push(ref, excStack)
        WRITEMEM(locals(xj), object, ref)
      else YIELDUP(Exc(ref))
  try ▶ abr catch (E1 x1) stm1 ... catch (En xn) stmn → YIELDUP(abr)
  try Exc(ref) ... catch (...) ▶ res ... →
    {excStack := pop(excStack), YIELDUP(res)}
  try tryBlock finally finallyBlock → pos := tryBlock
  try ▶ res finally finallyBlock → pos := finallyBlock
  try res finally ▶ Norm → YIELDUP(res)
  try res finally ▶ Exc(ref) → YIELDUP(Exc(ref))
  Exc(ref) → if pos = body(meth) ∧ ¬Empty(frames) then
    if StaticCtor(meth) then typeState(type(meth)) := Exc(ref)
    else APPENDSTACKTRACE(ref, meth(top(frames)))
  EXITMETHOD(Exc(ref))

```

In case an exception happened in the static constructor of a type, its type state is set to that exception to prevent its re-initialization and instead to re-throw the old exception object. The refinement of the macro INITIALIZE defined in C<sub>#c</sub> re-throws the exception object of a type which had an exception in the



static constructor, thus preventing its re-initialization.<sup>20</sup>

```
INITIALIZE(c) ≡
  ...
  if typeState(c) = Exc(ref) then YIELD(Exc(ref))
```

## 6 Refinement $C\#_D$ of $C\#_\mathcal{E}$ by delegates

In this section we extend  $C\#_\mathcal{E}$  by features which distinguish  $C\#$  from other languages, e.g. Java. We start with delegates and then add further constructs whose semantics can be defined mainly by reducing them via syntactical translations to the language model developed so far: properties, indexers, overloaded operators, enumerators with the **foreach** statement, the **using** statement, events and attributes. We use the model developed so far as ground model (in the sense of [12]) for  $C\#$ , thus providing a basis to justify the correctness (with respect to the ECMA standard [2]) of the “semantics of syntactic sugar” introduced in this section to define the semantics for delegates, properties, etc.

### 6.1 Delegates

Delegate types in  $C\#$  are reference types that encapsulate a static or instance method with a specific signature, with the intention of having delegates playing the role of type-safe function pointers. A delegate type  $D$  is declared as follows:

```
delegate  $T$   $D(S_1\ x_1, \dots, S_n\ x_n)$ ;
```

It represents the type of methods that take  $n$  arguments of type  $S_1, \dots, S_n$  and have return type  $T$ . Delegate types appear as subtypes of `System.Delegate` and provide in particular the *callback* functionality and asynchronous event handling. More precisely, the characteristic ability of delegates is to call a list of multiple methods sequentially. This feature is realized by means of an *invocationList*:  $Ref \rightarrow Delegate^* \cup \{Undef\}$  with which each delegate instance is equipped upon its creation. Each such list is a per instance immutable, non-empty, ordered list of static methods or pairs of target objects and instance methods. Upon invocation of a delegate instance with arguments *args*, the

<sup>20</sup>For modeling the implementation flag *beforefieldinit* mentioned above this implies, as observed in [13], to refine also the predicate *TriggerInit*, used for invoking static or instance methods, namely to guarantee for a class in exception state its initialization even if the class is marked *beforefieldinit*:  $TriggerInit(c) = (\neg Initialized(c) \wedge \neg beforefieldinit(c)) \vee typeState(c) = Exc(ref)$ .

methods of its invocation list are called one after the other with these arguments *args*, returning to the caller of the delegate either the *return value* of the last list element or the first *exception* a list element has thrown during its execution, preventing the remaining list elements from being invoked.

Therefore we introduce a new universe  $Delegate = Meth \cup (Ref \times Meth)$ . To express the creation and use of new delegate expressions the sets *Exp*, *Sexp* are extended by additional grammar rules as follows, using a new set *Dexp* of delegate expressions:

$$\begin{aligned} Sexp &::= \dots \mid Exp ([Args] ) \\ Exp &::= \dots \mid \text{'new' } DelegateType \text{'(' } Dexp \text{'(')} \\ Dexp &::= Meth \mid Type \text{'.'} Meth \mid Exp \text{'.'} Meth \mid Exp \end{aligned}$$

A method  $T::M$  is called *compatible* with the delegate type  $D$  iff  $T::M$  and  $D$  have the same return type and the same number of parameters with the same parameter types (including `ref`, `out`, `params` modifiers). The type constraints on the new expressions are spelt out in [19].

We use the model  $EXEC\text{SHARP}STM_I$ , which includes a description of the `for` statement of  $C\sharp_I$ , to express the sequentiality of the execution of delegate invocation list elements. In fact the above delegate declaration can be translated for  $T \neq \text{void}$  in the following class:

```
sealed class D : System.Delegate {
  public T Invoke(S1 x1, ..., Sn xn) {
    T result;
    for (int i = 0; i < this._length() ; i++)
      result = this._invoke(i, x1, ..., xn);
    return result;
  }
  private extern int _length();
  private extern T _invoke(int i, S1 x1, ..., Sn xn);
}
```

A delegate invocation expression  $exp(args)$  can be syntactically translated into a normal method call  $exp.D::Invoke(args)$  where  $D$  is the type of  $exp$ .<sup>21</sup> It then suffices to refine the ASM rule `INVOKEEXTERN` defined in the model  $EXEC\text{SHARP}EXP_C$  to describe the meaning of the method  $D::_invoke$ , which is to invoke the  $i$ th element of the invocation list on the given arguments, and

<sup>21</sup> In [2, §10.4.7] the members of a delegate are defined to be the members inherited from the class `System.Delegate`. However neither .NET nor Rotor nor Mono do respect this stipulation since they add further methods to those inherited. One such example is the method `_invoke` we use here.

analogously of `_length`.

```

INVOKEEXTERN( $T::M$ ,  $vals$ )  $\equiv$ 
  if  $T \in DelegateType$  then
    if  $name(M) = \_length$  then
      DELEGATELENGTH( $vals(0)$ )
    if  $name(M) = \_invoke$  then
      INVOKEDELEGATE( $vals(0)$ ,  $vals(1)$ ,  $drop(vals, 2)$ )

```

```

DELEGATELENGTH( $ref$ )  $\equiv$ 
  YIELDUP( $length(invocationList(ref))$ )

```

```

INVOKEDELEGATE( $ref$ ,  $i$ ,  $vals$ )  $\equiv$ 
  match  $invocationList(ref)(i)$ 
     $T::M$   $\rightarrow$  INVOKESTATIC( $T::M$ ,  $vals$ )
    ( $target$ ,  $T::M$ )  $\rightarrow$  INVOKEINSTANCE( $T::M$ ,  $target$ ,  $vals$ )

```

Since there are no new statements appearing in  $C\sharp_D$ , the addition of the rule EXECSHARP<sub>D</sub> to EXECSHARP<sub>E</sub> consists in the following ASM subrule EXECSHARPEXP<sub>D</sub>, which defines the meaning of delegate instance creation. For a detailed analysis of the discrepancy we exhibit here between the ECMA standard and the .NET implementation see [13].

```

EXECSHARPEXPD  $\equiv$  match  $context(pos)$ 
  new  $D(T::M) \rightarrow$ 
    let  $d = new(Ref, D)$  in
       $runTimeType(d) := D$ 
       $invocationList(d) := [T::M]$ 
      YIELD( $d$ )
  new  $D(exp.T::M) \rightarrow pos := exp$ 
  new  $D(\blacktriangleright ref.T::M) \rightarrow$ 
    if  $ref = Null$  then FAILUP(NullReferenceException)
    else let  $d = new(Ref, D)$  in
       $runTimeType(d) := D$ 
       $invocationList(d) := [(ref, T::M)]$ 
      YIELDUP( $d$ )
  new  $D(exp) \rightarrow pos := exp$ 
  new  $D(\blacktriangleright ref) \rightarrow$ 
    if  $ref = Null$  then FAILUP(NullReferenceException)
    else let  $d = new(Ref, D)$  in
       $runTimeType(d) := D$ 
       $invocationList(d) := invocationList(ref)$  // ECMA §14.5.10.3
      // Microsoft .NET Framework:
      //  $invocationList(d) := [(ref, D::Invoke)]$ 
      YIELDUP( $d$ )

```

To be complete, one should add some rules which reflect the special character of delegate invocation lists. As usual for lists, two invocation lists are *equal* (`==`) iff they have the same length and the elements of the lists are pairwise equal; they can be *combined* (concatenated with `+`) and sublists determined by a particular prefix and suffix condition can be *removed* from them (with `-`). To describe this specialization of list operations in our model it suffices to refine the macro `INVOKEEXTERN` by new rules for these operators `operator+`, `operator-`, `operator==`.

```

INVOKEEXTERN(T::M, vals) ≡
  ...
  if T ∈ DelegateType then
    if name(M) = operator+ then
      DELEGATECOMBINE(T, vals(0), vals(1))
    if name(M) = operator- then
      DELEGATEREMOVE(T, vals(0), vals(1))
    if name(M) = operator== then
      DELEGATEEQUAL(vals(0), vals(1))

```

Since invocation lists are considered to be immutable, combination and removal return *new* delegate instances unless one of the arguments is `null`. The `null` reference represents a delegate instance with an empty invocation list.

```

DELEGATECOMBINE(D, r1, r2) ≡
  if r1 = Null then YIELDUP(r2)
  elseif r2 = Null then YIELDUP(r1)
  else let d = new(Ref, D) in
    runTimeType(d) := D
    invocationList(d) := invocationList(r1) · invocationList(r2)
    YIELDUP(d)

```

```

DELEGATEREMOVE(D, r1, r2) ≡
  if r1 = Null then YIELDUP(Null)
  elseif r2 = Null then YIELDUP(r1)
  else let l1 = invocationList(r1) and l2 = invocationList(r2) in
    if l1 = l2 then YIELDUP(Null)
    elseif Subword(l2, l1) then let d = new(Ref, D) in
      runTimeType(d) := D
      invocationList(d) := prefix(l2, l1) · suffix(l2, l1)
      YIELDUP(d)
    else YIELDUP(r1)

```

The notions of *prefix* and *suffix* are defined here in terms of the *last* occurrence of a subword: *prefix*(*u*, *v*) is the part of *v* before the last occurrence of *u* in *v*

and  $\text{suffix}(u, v)$  the part of  $v$  after the last occurrence of  $u$  in  $v$ .

```

DELEGATEEQUAL( $r_1, r_2$ )  $\equiv$ 
  if  $r_1 = \text{Null} \vee r_2 = \text{Null}$  then YIELDUP( $r_1 = r_2$ )
  else let  $l_1 = \text{invocationList}(r_1)$  and  $l_2 = \text{invocationList}(r_2)$  in
    YIELDUP( $\text{length}(l_1) = \text{length}(l_2) \wedge \forall i < \text{length}(l_1) (l_1(i) = l_2(i))$ )

```

## 6.2 Properties, events and further features in $C_{\#D}^{\sharp}$

In this section we add further language features of  $C_{\#D}^{\sharp}$  whose semantics can be easily defined in terms of the model developed so far, essentially by simple syntactical reductions which one can easily justify to formalize correctly the explanations in [2].

**Properties.** Collections of a read and/or a write method for attributes of a class or struct are called *properties* in  $C_{\#D}^{\sharp}$  and declared in the following form (we skip the modifiers):

*Type Identifier* ‘{’ [‘get’ *Block*] [‘set’ *Block*] ‘}’

By definition a *read-write* property has a **get** and a **set** accessor, a *read-only* property has only a **get** accessor, a *write-only* property has only a **set** accessor. The identifier of a property  $P$  of type  $T$  can be used like a field identifier,<sup>22</sup> except that it cannot be passed as **ref** or **out** argument. Furthermore it is required that the body of a **get** accessor is the body of a method with return type  $T$ , that a **set** accessor has a value parameter named **value** of type  $T$  and that its body is the body of a **void** method. Using the signatures  $T \text{ get\_}P()$ ; and **void set\\_}P( $T \text{ value}$ )**;, which are reserved for get and set accessors, the intended semantics of properties is reduced to the semantics of methods, using the following syntactical reductions:

$$\begin{array}{l}
 T \text{ } P \{ \\
 \quad \text{get } \{ \text{getAccessor} \} \\
 \quad \text{set } \{ \text{setAccessor} \} \\
 \} \\
 \hline
 \end{array}
 \Longrightarrow
 \begin{array}{l}
 T \text{ get\_}P() \{ \text{getAccessor} \} \\
 \text{void set\_}P(T \text{ value}) \{ \\
 \quad \text{setAccessor} \\
 \} \\
 \hline
 \end{array}$$

$$\text{exp}.P \Longrightarrow \text{exp}.get\_P() \quad \text{exp}_1.P = \text{exp}_2; \Longrightarrow \text{exp}_1.set\_P(\text{exp}_2);$$

The above translation comprises also expressions of the form  $\text{exp}_1.P \text{ op} = \text{exp}_2$ , because we can assume that these compound assignments are compiled to

<sup>22</sup> Without knowing whether it is accessed directly or whether an accessor method is being called.

$\langle x = exp_1, y = x.get\_P() \text{ op } exp_2, x.set\_P(y), y \rangle$  with fresh local variables  $x, y$ , using as auxiliary operator the comma operator familiar from C/C++. This necessitates auxiliary rules for going through sequences of expressions of the following form:

$$\begin{aligned} \langle exp, \dots \rangle &\rightarrow pos := exp \\ \langle val_1, \dots, \blacktriangleright val_n \rangle &\rightarrow YIELDUP(val_n) \\ \langle \dots \blacktriangleright val, exp \dots \rangle &\rightarrow pos := exp \end{aligned}$$

**Indexers.** Indexers can be used like array elements except that they cannot contain `ref` or `out` parameters and their elements cannot be passed as `ref` or `out` arguments. They are declared in a class or struct type as follows (we skip the modifiers):

*Type* `'this' '[' [Params] ']' '{' ['get' Block] ['set' Block] '}'`

Analogously to the constraints for properties, for an indexer of type  $T$  with parameters  $p$ , the body of a `get` accessor is the body of a method with parameters  $p$  and return type  $T$ , the body of a `set` accessor is the body of a `void` method with parameters  $p$  and an implicit value parameter named `value` of type  $T$ . A base class indexer can be accessed by `base[exps]`. Using the signatures  $T$  `get_Item(params)` and `void set_Item(params, T value)`, which are reserved for `get` and `set` accessors, the intended semantics of indexers is reduced to the semantics of arrays and methods via the following compile-time translation (and corresponding operator expression translation as explained for properties):

$$\begin{array}{ccc} T \text{ this}[params] \{ & & T \text{ get\_Item}(params) \{ \text{getAccessor} \} \\ \quad \text{get} \{ \text{getAccessor} \} & \implies & \text{void set\_Item}(params, T \text{ value}) \{ \\ \quad \text{set} \{ \text{setAccessor} \} & & \quad \text{setAccessor} \\ \} & & \} \end{array}$$

**Events.** Events can be declared in C $\sharp$  like fields as follows: in the form `'event' DelegateType Identifier ';' (we omit the modifiers)`, or like properties, in the form

`'event' DelegateType Identifier '{' 'add' Block 'remove' Block '}'`.

Outside the type that contains the declaration, an event  $X$  can only be used as the left-hand operand of `+=` and `-=` in expressions  $X += exp$  and  $X -= exp$  of type `void`; within the type that contains the declaration, field-like events can be used like fields of delegate types. The accessors of property-like events have to be bodies of `void` methods with an implicit parameter `value` of *DelegateType*.

The semantics of events in  $C\sharp$  follows the *Publish/Subscribe* pattern. A class publishes an event it can raise, so that any number of classes can subscribe to that event. When the event is actually raised, each subscriber is notified that the event has occurred, namely by calling a delegate whose invocation list is executed with the sender object and the event data as its arguments. This idea is realized as follows. The *event sender* class that raises an event named  $X$  has the member event  $X\_EventHandler$   $X$ ; where the delegate type  $X\_EventHandler$  for the event is declared as follows (with two arguments, the first one for the publisher and the second one for the event information object, which must be derived from the class `EventArgs`):

```
delegate void X_EventHandler(object sender, X_EventArgs e);
```

To consume the event, the *event receiver* declares an event-handling method `Receive_X` with the same signature as the event delegate:

```
void Receive_X(object sender, X_EventArgs e) { ... }
```

To *register* the event handler, the event receiver has to add the `Receive_X` method to the event  $X$  of the event sender object:

```
X += new X_EventHandler(this.Receive_X);
```

The event sender *raises* the event by invoking the invocation list of  $X$  with the sender object and the event data, e.g.

```
void On_X(X_EventArgs e) { if (X != null) X(this, e); }
```

It suffices to assign a meaning to the signatures `void add_X(D value)` and `void remove_X(D value)`, which are reserved for every event  $X$  of delegate type  $D$ . This is done by the following translation of field-like events, anticipating the `lock` statement of  $C\sharp_{\mathcal{T}}$  which is explained in [17].<sup>23</sup>

```
class C {
  private D __X;
  void add_X(D value) {
    lock (this) { __X = __X + value; }
  }
  void remove_X(D value) {
```

<sup>23</sup> If one prefers to be independent of the thread model  $C\sharp_{\mathcal{T}}$ , one can consider lock statements `lock (exp) stm` translated for single-thread execution by `{ Object o = exp; stm }` (with a fresh variable  $o$ ), which is then refined in  $C\sharp_{\mathcal{T}}$  for the multiple thread execution model.

```

    lock (this) { __X = __X - value; }
  }
}

```

**Further constructs.** Similar syntactical reductions to those given above can be used to define the semantics of overloaded standard mathematical operators and user-defined conversions, of enumeration related statements ‘`foreach (T x in exp) stm`’, of using statements ‘`using ( resource ) stm`’, of parameter arrays and of attributes.

## 7 Refinement $C\sharp_u$ by pointers in unsafe code

In this section we add the features  $C\sharp$  offers for using pointers (coming with address-of and dereferencing operators ‘`&`’, ‘`*`’, ‘`->`’ together with pointer arithmetic) to directly work on memory addresses, bypassing the type checking by the compiler – hence the name ‘`unsafe`’ code blocks. Java has no such unsafe extension. The extension includes a mechanism called *pinning* of objects to prevent the runtime during the execution of a ‘`fixed`’ statement to manage via the garbage collector memory one wants to address directly. Code for which (de-) allocation is not controlled by the runtime is called *unmanaged*. As an alternative to pinning, data of unmanaged type can also be ‘`stackalloc`’ated, instead of using the heap.

The refinement consists, besides some new rules, mainly in a definition of the *memory* function in terms of byte sequences. This is a typical data refinement, using an encoding of simple types and a corresponding refinement of the function *structField*.

### 7.1 Signature refinement for $C\sharp_u$

We refine *Type* by adding pointer types to value and reference types.

$$\begin{aligned}
 \textit{Type} & ::= \textit{ValueType} \mid \textit{RefType} \mid \textit{PointerType} \\
 \textit{PointerType} & ::= \textit{UnmanagedType} \textit{'*'} \mid \textit{'void'} \textit{'*'}
 \end{aligned}$$

where *unmanaged* types are types which are not managed and *managed* types are recursively defined as a) reference types or b) struct types that contain a field of a managed type or a pointer to a managed type. The subtype relation is extended to pointer types such that  $\Lambda \preceq T^* \preceq \textit{void}^*$ . *Exp* and *Vexp* are extended by address-of and dereferencing expressions and expressions to denote the values of a new function indicating the ‘`sizeof`’ unmanaged types. *Stm* is extended to reflect ‘`unsafe`’ code blocks, ‘`fixed`’ statements and



Table 6

Type constraints for unsafe expressions.

Expression	Constraints	Expression type
<code>sizeof(<i>t</i>)</code>	<i>t</i> unmanaged type	<code>int</code>
<code>*<i>e</i></code>	$type(e) = T*$ , $T \neq \text{void}$	<i>T</i>
<code>&amp;<i>v</i></code>	<i>v</i> a fixed variable	$T*$ , where $T = type(v)$
<code><i>e</i> -&gt; <i>m</i></code>	$type(e) = T*$ , $T \neq \text{void}$	$type(T::m)$
<code><i>e</i>[<i>i</i>]</code>	$type(e) = T*$ , $T \neq \text{void}$ , $type(i)$ integral	<i>T</i>

‘`stackalloc`’ation of arrays. ‘`unsafe`’ can also appear as modifier for classes, structs, interfaces, delegates as well as for fields, methods, properties, indexers, operators, events, constructors, destructors.

$$Exp ::= \dots \mid \text{'\&'} Vexp \mid Exp \text{'->'} Meth ([Args]) \mid Exp \text{'->'} Field$$

$$\mid \text{'sizeof'} \text{'(' } UnmanagedType \text{'\>'}$$

$$Vexp ::= \dots \mid \text{'*'} Exp$$

$$Stm ::= \dots \mid \text{'unsafe'} Block \mid \text{'fixed'} \text{'(' } PointerType Loc = Exp \text{'\>' } Stm$$

$$Bstm ::= \dots \mid PointerType Loc \text{'='}$$

$$\text{'stackalloc'} UnmanagedType \text{'[' } Exp \text{'\>'} \text{'\>'}$$

In the following expressions, the basic arithmetical operators are used for pointer increment and decrement, pointer addition and subtraction, pointer comparison, and pointer conversion (where *p* and *q* are of a pointer type, *i* is of integer type):

- `++p`, `--p`, `p++`, `p--`, `p + i`, `i + p`, `p - i`, `p - q`, `p == q`, `p != q`, `p < q`, `p <= q`,  
`p > q`, `p >= q`
- `(T*)i`, `(T*)p`, `(int)p`, `(uint)p`, `(long)p`, `(ulong)p`

On the types of the new expressions the constraints in Table 6 are imposed. We assume the dereferencing and member access operator `e -> m` to be translated to `(*e).m`, similarly `e[i]` to `*(e + i)`.

For statements the following type constraints in Table 7 are assumed. A variable is called *moveable* (by the garbage collector) iff it is not a fixed variable. Fixed variables are (by recursive definition): local variables, value parameters, `*exp` for `exp` of pointer type, and instance field expressions `v.f` if *v* is a fixed variable of struct type *T* and *f* is an instance field of *T*.

The local variable *p* in the fixed statement is called a *pinned* local variable. A pinned local variable is a read-only variable. It is not allowed to assign a new value to it in the body of the fixed statement.

Table 7

Type constraints for unsafe statements.

Statement	Constraints
$T* p = \text{stackalloc } T[exp];$	$type(exp) = \text{int}, T$ unmanaged
$\text{fixed}(\text{char} * p = exp) \text{ stm}$	$type(exp) = \text{string}, p$ read-only in $stm$
$\text{fixed}(T * p = exp) \text{ stm}$	$type(exp) = T[R], T$ unmanaged, $p$ read-only in $stm$
$\text{fixed}(T * p = \&vexp) \text{ stm}$	$type(vexp) = T, T$ unmanaged, $vexp$ a moveable variable, $p$ read-only in $stm$

The principal refinement in the ASM extension  $\text{EXECCSHARP}_U$  for  $\text{C}\sharp_U$  is that of the *memory* together with its operators, where the set of *SimpleValues* is replaced by *Bytes* (8-bit strings), using non-negative integers as memory addresses ( $\text{Adr} = \mathbb{N}$ ):

$$\text{mem}: \text{Adr} \rightarrow \text{Byte} \cup \text{Ref} \cup \{\text{Undef}\}$$

The partial functions to *encode* (resp. *decode*) values of a given simple type  $T$  by byte sequences, of a length (number of bytes) depending on  $\text{sizeof}(T)$ , satisfy for values  $val$  the equations

$$\text{decode}(T, \text{encode}(val)) = val \quad \text{and} \quad \text{length}(\text{encode}(val)) = \text{sizeof}(T).$$

For every pointer type  $T*$  holds  $\text{sizeof}(T*) = \text{sizeof}(\text{void}^*)$ .

A function  $\text{fieldOffset}: \text{UnmanagedType} \times \text{Field} \rightarrow \mathbb{N}$  is used to describe the layout of unmanaged structs. It has to satisfy the following constraint for every unmanaged struct type  $T$  and instance field  $f$  of  $T$  (overlapping fields are allowed in  $\text{C}\sharp_U$ ):

$$\text{fieldOffset}(T, f) + \text{sizeof}(\text{type}(f)) \leq \text{sizeof}(T)$$

We assume that if  $adr$  is an address allocated using  $\text{new}(\text{Adr}, T)$  for struct type  $T$ , then for every instance field  $f$  of  $T$  the equation  $\text{fieldAdr}(adr, f) = adr + \text{fieldOffset}(T, f)$  holds.

To determine the layout of arrays with unmanaged element type we stipulate the following refinement of the function  $\text{elemAdr}$  which reflects that array elements are stored such that the indices of the right most dimension are increased first, then the next left dimension, and so on. For  $\text{runTimeType}(\text{ref}) = T[\text{dim}(n)]$ , where  $T$  is an unmanaged type and  $l_i = \text{dimLength}(\text{ref}, i - 1)$

for  $i \in [1 \dots n]$ , we assume the following:

$$\begin{aligned} elemAdr(ref, [i_1, i_2, \dots, i_n]) = \\ elemAdr(ref, [0, \dots, 0]) + (\dots (i_1 \cdot l_2 + i_2) \cdot l_3 + \dots + i_n) \cdot sizeOf(T) \end{aligned}$$

## 7.2 Transition rule refinement for unsafe code

Besides the rules below which define the semantics of the new expressions and statements, to be added to  $EXEC\text{SHARP}_D$ , we have to data refine the notions of reading from and writing to memory for values of unmanaged type.

$$\begin{aligned} memValue(adr, t) = \\ \mathbf{if} \ t \in RefType \ \mathbf{then} \ mem(adr) \\ \mathbf{elseif} \ t \in UnmanagedType \ \mathbf{then} \\ \quad [mem(adr + i) \mid i \in [0 \dots sizeOf(t) - 1]] \\ \mathbf{elseif} \ t \in StructType \ \mathbf{then} \\ \quad \{f \mapsto getField(adr, f) \mid f \in instanceFields(t)\} \end{aligned}$$

$$getField(adr, t::f) = memValue(fieldAdr(adr, t::f), type(t::f))$$

$$\begin{aligned} WRITE\text{MEM}(adr, t, val) \equiv \\ \mathbf{if} \ t \in RefType \ \mathbf{then} \ mem(adr) := val \\ \mathbf{elseif} \ t \in UnmanagedType \ \mathbf{then} \\ \quad \mathbf{forall} \ i \in [0 \dots sizeOf(t) - 1] \ \mathbf{do} \ mem(adr + i) := val(i) \\ \mathbf{elseif} \ t \in StructType \ \mathbf{then} \\ \quad \mathbf{forall} \ f \in instanceFields(t) \ \mathbf{do} \ SET\text{FIELD}(adr, f, val(f)) \end{aligned}$$

$$\text{SET}\text{FIELD}(adr, t::f, val) \equiv \text{WRITE}\text{MEM}(fieldAdr(adr, t::f), type(t::f), val)$$

Values of unmanaged struct types are directly represented as sequences of bytes. Hence, the function *structField* has to be refined to extract a subsequence in case of unmanaged struct types:

$$\begin{aligned} structField(val, T, f) = \\ \mathbf{if} \ T \in ManagedType \ \mathbf{then} \ val(f) \\ \mathbf{else} \ \mathbf{let} \ n = fieldOffset(T, f) \ \mathbf{in} \ [val(i) \mid n \leq i < n + sizeOf(type(f))] \end{aligned}$$

In the rules for  $EXEC\text{SHARP}\text{EXP}_U$  we have  $\&$   $\square$  as additional address position. We follow the implementation in Rotor and .NET in formulating the *Null* check to prevent writing to null addresses; the ECMA standard describes this check as optional.

$$\begin{aligned} EXEC\text{SHARP}\text{EXP}_U \equiv \mathbf{match} \ context(pos) \\ \quad sizeof(T) \rightarrow \text{YIELD}(sizeof(T)) \end{aligned}$$

```

&exp → pos := exp
&►adr → YIELDUP(adr)

*exp → pos := exp
*►adr → if adr = Null then // null pointer check optional
    FAILUP(NullReferenceException)
    else YIELDUPINDIRECT(adr)

*exp1 = exp2 → pos := exp1
*►adr = exp2 → pos := exp2
*adr = ►val → if adr = Null then // null pointer check optional
    FAILUP(NullReferenceException)
    else
        WRITEMEM(adr, type(pos), val)
        YIELDUP(val)

```

The rules for pointer arithmetic can be summarized as follows:

```

Apply(+ (T*, int), adr, i) = adr + i · sizeOf(T)
Apply(+ (int, T*), i, adr) = adr + i · sizeOf(T)
Apply(- (T*, T*), adr1, adr2) = (adr1 - adr2) / sizeOf(T)
Convert(T*, adr) = adr = Convert(S, adr)
    for S ∈ {int, uint, long, ulong}
Convert(T*, i) = i

```

In the execution of the `stackalloc` statement we assume that `new(adr, T, i)` allocates  $i$  consecutive chunks of addresses of size `sizeOf(T)` which are later de-allocated on method exit in `FREELOCALS`.

```

EXECCSHARPSTMU ≡ match context(pos)
    unsafe block → pos := block
    unsafe ►Norm → YIELDUP(Norm)

T* loc = stackalloc T[exp]; → pos := exp
T* loc = stackalloc T[►i]; → let adr = new(Adr, T, i) in
    WRITEMEM(locals(loc), T*, adr)
    YIELDUP(Norm)

```

The run-time execution of fixed statements can be explained by syntactical transformations.

Statement	Run-time execution
<code>fixed (char* p = exp) stm</code>	{ char* p; p = Cstring(exp); stm }
<code>fixed (T* p = exp) stm</code>	{ T* p; p = &exp[0]; stm }
<code>fixed (T* p = &amp;vexp) stm</code>	{ T* p; p = &vexp; stm }

In the first case, it is assumed that  $Cstring(s)$  is an internal function that returns the address of the first element of a C-style null-terminated character array representation of the string  $s$ . How it is related to the original representation of the string is not specified in [2].

## 8 Related Work and Conclusion

One of our referees would like to see a critical assessment of the ASM method we used for this work and a comparison to alternative approaches. Some justification of the kind from the perspective of semantic methods for programming languages has been given in [26, Sect.4], containing concrete illustrations of and references to the numerous and earlier competing proposals. This was at a time when ASMs were applied for the first time to (successfully) specify an industrial language standard, namely the ISO Prolog standard [27]. A decade later, a broader comparison of the then well-developed ASM method with respect to other system design and analysis frameworks has been provided in [28,29]. However, a systematic, comprehensive and authoritative evaluation of the multitude of system design and analysis proposals in the literature remains a highly desirable and challenging task to be accomplished, even if limited to the use of the major so-called formal methods for the development and investigation of programming languages and their implementations. From the perspective of practical system design and analysis some comparative studies of this kind have been published, see e.g. [30–32] (the interested reader may also consult the corresponding ASM-based work in [33–35]). For work centered around Java and the JVM the reader finds in [36] a collection of formal-method-approaches to language specification and analysis; [37] contains an excellent, detailed and at the time complete review of the huge literature on the subject (including an evaluation of the ASM-based Java/JVM investigations), with a focus on safety issues and their impact on smart cards. We cannot perform here a similar analysis for work on C $\sharp$  or other major programming languages. This explains why the references in this paper stick to C $\sharp$  documentation from ECMA and Microsoft and to some ASM work we have built upon directly.

For the work presented in this paper we set ourselves a more modest though not completely trivial major goal, namely to test whether the method developed in [1] for the definition and a proven to be correct implementation of a real-life programming language like Java scales naturally to the somewhat richer and more complex C $\sharp$ . It is up to the reader to judge whether this ASM reuse case study for a real-life complex model succeeded. For the formalization of other programming languages something can also be learnt directly from the formalization of the semantics of C $\sharp$  worked out here. For example, how to “divide and conquer” the static and the dynamic semantics of a language, how to

separate the description of conceptually independent programming constructs by dividing them into sublanguages, how to unify and streamline the formalization of similar constructs by appropriate parameterizations (which means abstractions), how to model and evaluate variations of specific features (e.g. expression evaluation, parameter passing mechanism, class initialization, etc.) by varying macros, rules and/or domains together with their operations, how to extend within a single framework the model for a language core by a form of bootstrapping (including in particular syntactical translations) to a model for the entire language, etc.

There are several by-products of the work presented here. Through the ASM-model-oriented analysis of the ECMA standard for C# we found several bugs and gaps in the formulation of the standard and in its .NET implementation as well as some incoherences between the two, as documented in detail in [13] in terms of our ASM model for C#. Another by-product of the high-level modular interpreter defined here is the support it provides to teachers of C#, in particular if they want to shed light on certain subtle language features which are not clarified by the ECMA documents. In the forthcoming paper [38] we are going to work out a concrete comparison of the two models for C# and for Java, which will allow us to formulate in a precise technical manner where and in which respect the two languages differ among each other and from other programming languages – methodologically, semantically and pragmatically. As a specific part of this reuse-case-study the second author is investigating how the main new features of C# 2.0 can be modeled by appropriate extensions of the ASM model developed here for C#, in particular generic types (parametric polymorphism), anonymous methods and iterators. Last but not least, with our C# model and its extension to threads in [17] we have laid the ground for a mathematical analysis and possibly mechanical verification of interesting properties of the language and its implementation, like type safety,<sup>24</sup> compiler correctness and completeness, correctness of (a mathematical model to be developed for) garbage collection, security, etc. For the correctness of the definite assignment analysis performed by a C# compiler, we may refer the interested reader to [21]. We hope somebody will feel challenged to use our model for precisely formulating and proving such theorems for C# and to build a corresponding model for Microsoft’s Common Language Runtime together with a compilation scheme from C# to IL code, applying to our model the powerful ASM refinement technique [40] along the lines shown in the ASM-based Java/JVM study in [1].

The questions asked by our referees lead us to mention another practical and

---

<sup>24</sup>For a fragment of Microsoft’s Intermediate Language, which is executed by Microsoft’s Common Language Runtime, a type safety proof has been given in [20], based upon Syme’s method [39] for writing functional specifications which can be subject to theorem proving in HOL.

industrially viable use that can be made of a modeling and analysis activity as the one reported in this paper, except if the extreme time pressure usually imposed on developers to produce executable code from incomplete verbal specifications (mostly formulated in natural language) prevents them from at least once trying out a more reliable option. Here is a concrete example what could have been done. On September 27, 2000, the penultimate day of his sabbatical stay with Microsoft Research in Redmond, in a seminar talk to representatives of the C# development team, the first author suggested to use the method, at the time formulated and presented in terms of Java/JVM for publication in what became the Jbook [1], for the following five fundamental activities in relation to the at-the-time ongoing development of what became known as the C# language with the underlying CLR virtual machine:

- defining an ASM model as *executable specification of critical language constructs* or layers (if not of the entire language) and of the mapping to IL code,
- *generating test cases* for the implementing code from the ASM model,
- using the ASM model as *oracle for test evaluations* and for comparing model test runs with code test runs,
- using the ASM model as *internal documentation* for future language extensions and for relating other .NET languages to C#, in particular those which are equipped already with an ASM model of their semantics,
- using the ASM model as *basis for writing innovative handbooks* for users and for maintenance professionals, where the innovative character derives from being a) accurate yet simple and easy to understand, b) complete and detailed yet succinct.

## Acknowledgements

We gratefully acknowledge partial support of this work by a Microsoft grant within the ROTOR project during the year 2002–2003. We thank Bruno Quarta for attracting us to the C# modeling work, even if *post festam*, as part of the ROTOR project. We also thank two anonymous referees for valuable criticism which helped us improve the exposition.

## References

- [1] R. F. Stärk, J. Schmid, E. Börger, Java and the Java Virtual Machine—Definition, Verification, Validation, Springer-Verlag, 2001.
- [2] C# Language Specification, Standard ECMA–334, <http://www.ecma-international.org> (2001).

- [3] A. Hejlsberg, S. Wiltamuth, P. Golde, C# Language Specification, Addison-Wesley, 2003.
- [4] T. Archer, A. Whitechapel, Inside C#, Microsoft Press, 2002.
- [5] J. Prosise, Programming Microsoft .NET, Microsoft Press, 2002.
- [6] J. Richter, Applied Microsoft .NET Framework Programming, Microsoft Press, 2002.
- [7] Visual Studio .NET 2003, <http://msdn.microsoft.com/vstudio/>.
- [8] SSCLI (Rotor) web site, <http://www.sscli.net>.
- [9] Mono compiler for C#, <http://www.go-mono.com/c-sharp.html>.
- [10] Common Language Infrastructure, Standard ECMA-335, <http://www.ecma-international.org> (2003).
- [11] D. Stutz, T. Neward, G. Shilling, Shared Source CLI Essentials, O'Reilly, 2003.
- [12] E. Börger, The ASM ground model method as a foundation of requirements engineering, in: N. Dershowitz (Ed.), Manna-Symposium, Vol. 2772 of Lecture Notes in Computer Science, Springer-Verlag, 2003.
- [13] N. G. Fruja, Specification and implementation problems for C#, in: B. Thalheim, W. Zimmermann (Eds.), Abstract State Machines 2004, Lecture Notes in Computer Science, Springer-Verlag, 2004.
- [14] Foundations of Software Engineering Group, Microsoft Research, AsmL, Web pages at <http://research.microsoft.com/foundations/AsmL/> (2001).
- [15] J. Schmid, Executing ASM specifications with AsmGofer, Web pages at <http://www.tydo.de/AsmGofer>.
- [16] J. Schmid, Refinement and implementation techniques for Abstract State Machines, Ph.D. thesis, University of Ulm, Germany (2002).
- [17] R. F. Stärk, E. Börger, An ASM specification of C# threads and the .NET memory model, in: B. Thalheim, W. Zimmermann (Eds.), Abstract State Machines 2004, Lecture Notes in Computer Science, Springer-Verlag, 2004.
- [18] E. Börger, R. F. Stärk, Abstract State Machines. A Method for High-Level System Design and Analysis, Springer-Verlag, 2003.
- [19] E. Börger, N. G. Fruja, V. Gervasi, R. Stärk, A complete formal definition of the semantics of C#, Technical report, In preparation (2004).
- [20] A. D. Gordon, D. Syme, Typing a multi-language intermediate code, in: Proc. 28th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, London, 2001, pp. 248–260.
- [21] N. G. Fruja, The correctness of the definite assignment analysis in C#, Technical report 435, Computer Science Department, ETH Zürich (2004).



- [22] E. Börger, T. Bolognesi, Remarks on turbo ASMs for computing functional equations and recursion schemes, in: E. Börger, A. Gargantini, E. Riccobene (Eds.), *Abstract State Machines 2003 – Advances in Theory and Applications*, Vol. 2589 of *Lecture Notes in Computer Science*, Springer-Verlag, 2003, pp. 218–228.
- [23] W. Zimmermann, A. Dold, A framework for modeling the semantics of expression evaluation with Abstract State Machines, in: E. Börger, A. Gargantini, E. Riccobene (Eds.), *Abstract State Machines 2003–Advances in Theory and Applications*, Vol. 2589 of *Lecture Notes in Computer Science*, Springer-Verlag, 2003, pp. 391–406.
- [24] N. G. Fruja, R. F. Stärk, The hidden computation steps of turbo Abstract State Machines, in: E. Börger, A. Gargantini, E. Riccobene (Eds.), *Abstract State Machines 2003 – Advances in Theory and Applications*, Vol. 2589 of *Lecture Notes in Computer Science*, Springer-Verlag, 2003, pp. 244–262.
- [25] E. Börger, J. Schmid, Composition and submachine concepts for sequential ASMs, in: P. Clote, H. Schwichtenberg (Eds.), *Computer Science Logic (Proceedings of CSL 2000)*, Vol. 1862 of *Lecture Notes in Computer Science*, Springer-Verlag, 2000, pp. 41–60.
- [26] E. Börger, A logical operational semantics for full Prolog. Part I: Selection core and control, in: E. Börger, H. Kleine Büning, M. M. Richter, W. Schönfeld (Eds.), *CSL’89. 3rd Workshop on Computer Science Logic*, Vol. 440 of *Lecture Notes in Computer Science*, Springer-Verlag, 1990, pp. 36–64.
- [27] E. Börger, K. Dässler, Prolog: DIN papers for discussion, ISO/IEC JTC1 SC22 WG17 Prolog Standardization Document 58, National Physical Laboratory, Middlesex, England (1990).
- [28] E. Börger, High-level system design and analysis using Abstract State Machines, in: D. Hutter, W. Stephan, P. Traverso, M. Ullmann (Eds.), *Current Trends in Applied Formal Methods (FM-Trends 98)*, Vol. 1641 of *Lecture Notes in Computer Science*, Springer-Verlag, 1999, pp. 1–43.
- [29] E. Börger, Abstract State Machines: A unifying view of models of computation and of system design frameworks, *Annals of Pure and Applied Logic* To appear.
- [30] C. Lewerentz, T. Lindner, Formal Development of Reactive Systems. Case Study “Production Cell”, Vol. 891 of *Lecture Notes in Computer Science*, Springer-Verlag, 1995.
- [31] J.-R. Abrial, E. Börger, H. Langmaack, The steam boiler case study: Competition of formal program specification and development methods, in: J.-R. Abrial, E. Börger, H. Langmaack (Eds.), *Formal Methods for Industrial Applications. Specifying and Programming the Steam-Boiler Control*, Vol. 1165 of *Lecture Notes in Computer Science*, Springer-Verlag, 1996, pp. 1–12.
- [32] M. Broy, S. Merz, K. Spies, Formal Systems Specification – The RPC-Memory Specification Case Study, Vol. 1169 of *Lecture Notes in Computer Science*, Springer-Verlag, 1996.

- [33] E. Börger, L. Mearelli, Integrating ASMs into the software development life cycle, *J. Universal Computer Science* 3 (5) (1997) 603–665.
- [34] C. Beierle, E. Börger, I. Durdanović, U. Glässer, E. Riccobene, Refining abstract machine specifications of the steam boiler control to well documented executable code, in: J.-R. Abrial, E. Börger, H. Langmaack (Eds.), *Formal Methods for Industrial Applications. Specifying and Programming the Steam-Boiler Control*, no. 1165 in *Lecture Notes in Computer Science*, Springer-Verlag, 1996, pp. 62–78.
- [35] J. Huggins, Broy-Lampert Specification Problem: A Gurevich Abstract State Machine Solution, Technical Report CSE-TR-320-96, EECS Dept., University of Michigan (1996).
- [36] J. Alves-Foss, *Formal Syntax and Semantics of Java*, Vol. 1523 of *Lecture Notes in Computer Science*, Springer-Verlag, 1998.
- [37] P. Hartel, L. Moreau, Formalizing the safety of Java, the Java Virtual Machine and Java Card, *ACM Computing Surveys* 33 (4) (2001) 517–558.
- [38] E. Börger, R. F. Stärk, Exploiting abstraction for specification reuse. The Java/C# case study, in: M. B. et al. (Ed.), *Proc. FMCO'03*, *Lecture Notes in Computer Science*, Springer-Verlag, 2004.
- [39] D. Syme, *Declarative theorem proving for operational semantics*, Ph.D. thesis, University of Cambridge (1998).
- [40] E. Börger, The ASM refinement method, *Formal Aspects of Computing* 15 (2003) 237–257.