

A High-Level Modular Definition of the Semantics of C#

Egon Börger, Georgian Fruja, Vincenzo Gervasi, Robert F. Stärk

October 10, 2003

Abstract

We propose a structured mathematical definition of the semantics of C# programs to provide a platform-independent view of the language for the C# programmer, which can also be used for a precise analysis of the ECMA [7] standard of the language. The definition takes care to reflect directly and faithfully – as much as possible without becoming inconsistent or incomplete – the descriptions in the C# standard to become comparable with the corresponding models for Java in [15] and to provide for implementors the possibility to check their basic design decisions against an accurate high-level model. In particular we will highlight some of the major differences between the ECMA standard and the implementation of the language in .NET.

1 Introduction

In this work the method developed in [15] for a rigorous definition and analysis of Java and its implementation on the Java Virtual Machine (JVM) is applied to formalize the semantics of C#. We provide a 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 [7] and in [10] and clarifies a certain number of semantically relevant issues which are not handled by that standard. We also consulted the Microsoft Press books [1, 11, 12]. A series of bugs and gaps in the ECMA standard for C# and in .NET and incoherences between the two were detected during our attempt to build a consistent and complete yet abstract model for the language. To support the experimentation with the model a project has been started to refine the model developed here to .NET-executable AsmL code [8], similarly to the AsmGofer refinement developed by Joachim Schmid [13, 14] for the Java and JVM models in [15].

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 State Machine (ASM) whose notion of run provides a transparent way to reflect those run-time-related features, which are encountered upon executing a given C# program. The use of ASMs¹ also allows us to specify the static and the dynamic parts of the semantics separately. 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 annotated 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 [15] we structure the C# 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#), classes (realizing procedural abstraction with class initialization and global (module) variables), object-orientation (with class instances, instance methods, inheritance), exception handling, delegates together with events (including for convenience here also properties, indexers, attributes), concurrency, unsafe code with

¹more precisely the classification of abstract states into a static and a dynamic part

pointer arithmetic. This yields a sequence of sublanguages $C\#_{\mathcal{I}}$, $C\#_c$, $C\#_o$, $C\#_{\varepsilon}$, $C\#_{\mathcal{D}}$, $C\#_{\mathcal{T}}$, $C\#_{\mathcal{U}}$ which altogether describe the entire language $C\#$. Each language in the sequence extends its predecessor and for each one we build a submachine which is a conservative (purely incremental) extension of its predecessor. The model for the entire language $C\#$ is a parallel composition of all submachines.

To keep the definition of the models succinct, we avoid tedious and routine repetitions concerning language constructs which can be reduced in well-known ways to the core constructs in our models. Since the handling of truly concurrent threads, not limited to interleaving or similar simple scheduling techniques, is closely related to the underlying memory model whose description goes much beyond this paper, the submodel $C\#_{\mathcal{T}}$ and its further analysis is postponed to a forthcoming separate paper [6].

Since by and large one can correctly understand an ASM as pseudo-code operating over abstract data, we skip a detailed definition of ASMs which is available in textbook form in Chapter 2 of the *AsmBook* [5]. The basic framework of the model is introduced together with the model for the imperative kernel $C\#_{\mathcal{I}}$ of the language which is then refined by successive extensions to the full model. Since our paper is not a tutorial or manual on $C\#$, we restrict our explanations here to features a reader will appreciate who is already knowledgeable about the basic concepts of object-oriented programming. In a technical report [3] also the remaining details which are skipped in this paper are spelt out completely, together with further explanations and examples.

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

In this section we define the model for $C\#_{\mathcal{I}}$, which defines the basic machinery of the ASM model for $C\#$ and describes the semantics of the sequential imperative core of $C\#$ 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 $C\#$ program. We separate syntax and compile-time matters from run-time issues by assuming that the program is given as an 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 $C\#_{\mathcal{I}}$, we formulate what has to be said about the static semantics.

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 $C\#_{\varepsilon}$. 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 on the top-level like 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. $Sexps$ 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 [7, §13.1] and of binary numeric promotions in [7, §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 2.2 (Implicit conversion [7, §13.1]) There is 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 simple 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 .

```

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

```

Figure 1: Grammar of expressions and statements in $C\#_{\mathcal{I}}$.

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 [7] for the predefined (arithmetical, relational, bit and boolean logical) operators and the expression types. We assume that each expression node exp in the abstract syntax tree is annotated with 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 $e_1 \text{ } bop \text{ } e_2$, then the least upper bound T of the types of e_1 and e_2 must exist and the expression is transformed into $(T) e_1 \text{ } bop \text{ } (T) e_2$.

Definition 2.3 (Binary numeric promotion [7, §14.2.6]) 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 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 [7], 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.

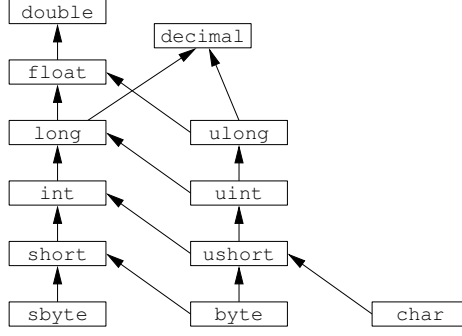


Figure 2: The simple types of C#.

To simplify the exposition of our model we assume some standard syntactical reductions as indicated in the following table:

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

During the static program analysis where the compiler has to verify that the given program is well-typed, predicates *reachable* and *normal* with the following intended meaning are computed for statements, using the type information contained in the annotated syntax tree as result of parsing and elaboration:

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

For the correctness of our model we have to guarantee the following two program properties: a) during the program execution, only *reachable* positions are reached, b) normal termination happens only in *normal* positions. These two properties, which are undecidable, are guaranteed by checking two sufficient conditions via so-called reachability rules, which can be inductively defined for C# as follows (similarly for `do`, `for`, `switch`). For constant boolean expressions in conditional and while statements we assume that they are replaced in the abstract syntax tree by `true` or `false`.

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

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 of the bottom of the code array.

We also have to reflect in the model the type safety of well-typed C# programs, i.e. that a) variables at run-time contain values that are *compatible* with the declared types, and b) expressions

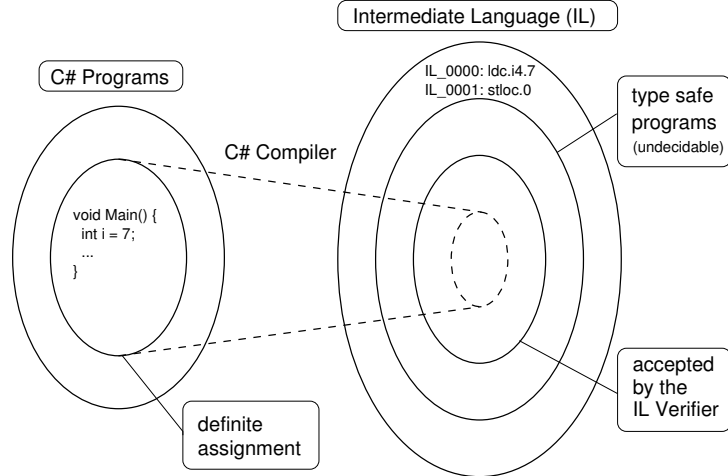


Figure 3: Definite assignment and IL verification.

are evaluated at run-time to values that are *compatible* with their compile-time types. The desired consequence of the type safety of a program is that at run-time 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 assumptions².

To guarantee the type safety, which is an undecidable program property, 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 the reality. For unknown reasons Microsoft has decided that in verified IL (intermediate language) code local variables are initialized by the run-time system with zero values. Hence, also 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.

$$\begin{aligned}
 x \in \text{before}(\text{item}) &\iff x \text{ is definitely assigned } \textit{before} \text{ the elaboration of } \textit{item} \\
 x \in \text{after}(\text{item}) &\iff x \text{ is definitely assigned } \textit{after} \text{ normal elaboration of } \textit{item} \\
 x \in \text{true}(\text{exp}) &\iff x \text{ is definitely assigned } \textit{after} \text{ exp evaluates to } \textit{true} \\
 x \in \text{false}(\text{exp}) &\iff x \text{ is definitely assigned } \textit{after} \text{ exp evaluates to } \textit{false}
 \end{aligned}$$

We make the definite assignment rules of [7, §12.3.3] precise by the following constraints, where $\text{vars}(\text{stm}) = \{x \mid \text{stm} \text{ is in the scope of } x\}$.

- If s is a function body, then $\text{before}(s) = \emptyset$
- $\text{after}(\text{break};) = \text{vars}(\text{break};)$
- $\text{after}(\text{;}) = \text{before}(\text{;})$
- $\text{after}(\text{continue};) = \text{vars}(\text{continue};)$
- $\text{before}(\text{exp}) = \text{before}(\text{exp};)$
- $\text{after}(\text{goto } l;) = \text{vars}(\text{goto } l;)$
- $\text{after}(\text{exp};) = \text{after}(\text{exp})$

For blocks $\text{stm} = \{s_1 \dots s_n\}$ the constraints are as follows:

²For example the following invariants can be proved to hold at run-time: a) $\text{before}(\text{pos}) \subseteq \text{Defined}$ where $\text{Defined} = \{x \in \text{Loc} \mid \text{mem}(\text{locals}(x)) \neq \text{Undef}\}$, b) $\text{after}(\text{pos}) \subseteq \text{Defined}$ if $\text{values}(\text{pos}) = \text{Norm}$ or $\text{values}(\text{pos}) \in \text{Value}$. Specifically for boolean expressions holds $\text{true}(\text{pos}) \subseteq \text{Defined}$ if $\text{values}(\text{pos}) = \text{True}$, the same for *false*. Such proofs can be carried out using the pattern developed in [15, Ch.8] for proving that Java is type safe.

- $before(s_1) = before(stm)$
- $after(stm) = after(s_n)$
- $before(s_{i+1}) = after(s_i) \cap goto(s_{i+1})$ where
 $goto(l:s) = \bigcap \{before(goto\ l;)\mid goto\ l\text{ is reachable in }stm\}$ and $goto(s) = vars(s)$ if s is not a labeled statement

For $stm = \text{if } (e) s_1 \text{ else } s_2$ one has to require:

- $before(e) = before(stm)$
- $before(s_1) = true(e)$
- $before(s_2) = false(e)$
- $after(stm) = after(s_1) \cap after(s_2)$

Constraints for while statements $stm = \text{while } (e) s$:

- $before(e) = before(stm)$
- $before(s) = true(e)$
- $after(stm) = false(e) \cap break(s)$ where
 $break(s) = \bigcap \{before(\text{break};)\mid \text{break};\text{ reachable in }s\}$

For boolean expressions we have the following general constraints.

If $exp \in \{\text{true}, \text{false}, !e, e_1 \ \&\& \ e_2, e_1 \ || \ e_2, e_0 ? e_1 : e_2\}$, then

- $after(exp) = true(exp) \cap false(exp)$

otherwise

- $true(exp) = after(exp)$
- $false(exp) = after(exp)$

In addition for specific boolean expressions, the following specific constraints are imposed for the eager ('short-circuit') evaluation of Boolean expressions.

- $true(\text{true}) = before(\text{true})$
- $false(\text{true}) = vars(\text{true})$
- $true(\text{false}) = vars(\text{false})$
- $false(\text{false}) = before(\text{false})$

For negations $exp = !e$:

- $before(e) = before(exp)$
- $true(exp) = false(e)$
- $false(exp) = true(e)$

For conjunctions $exp = (e_1 \ \&\& \ e_2)$:

- $before(e_1) = before(exp)$
- $before(e_2) = true(e_1)$
- $true(exp) = true(e_2)$
- $false(exp) = false(e_1) \cap false(e_2)$

For disjunctions $exp = (e_1 \ || \ e_2)$:

- $before(e_1) = before(exp)$
- $before(e_2) = false(e_1)$
- $true(exp) = true(e_1) \cap true(e_2)$
- $false(exp) = false(e_2)$

If $exp = (e_0 ? e_1 : e_2)$, then

- $before(e_0) = before(exp)$
- $before(e_1) = true(e_0)$
- $before(e_2) = false(e_0)$
- $true(exp) = true(e_1) \cap true(e_2)$
- $false(exp) = false(e_1) \cap false(e_2)$

For general expressions the constraints for definite assignment are as follows.

- $loc \in before(loc)$
- $after(loc) = before(loc)$
- $after(lit) = before(lit)$

For simple assignments $exp = (loc = e)$ we have

- $before(e) = before(exp)$
- $after(exp) = after(e) \cup \{loc\}$

For compound assignments $exp = (loc \ op = e)$ we have

- $loc \in before(exp)$
- $before(e) = before(exp)$
- $after(exp) = after(e)$

If $exp = (e_0 ? e_1 : e_2)$, then

- $before(e_0) = before(exp)$
- $before(e_1) = true(e_0)$
- $before(e_2) = false(e_0)$
- $after(exp) = after(e_1) \cap after(e_2)$

In all other cases, if exp is an expression with *direct subexpressions* e_1, e_2, \dots, e_n , then the left-to-right evaluation scheme yields

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

Due to the goto statement the above constraints specify the sets of variables that have to be considered as definitely assigned not in a unique way. For blocks without goto statements, however, it can be shown that the *before* set determines the *after* set in a unique way.

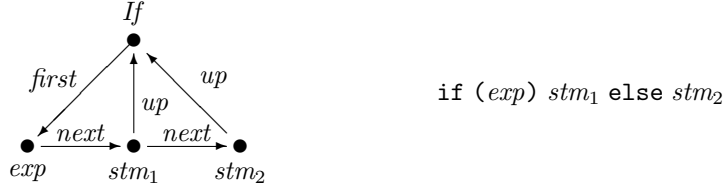
2.4 Transition rules for $C\#_{\mathcal{I}}$

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

$$\begin{aligned} \text{EXEC}\text{C}\text{\#}\text{SHARP}_{\mathcal{I}} &\equiv \\ &\text{EXEC}\text{C}\text{\#}\text{SHARP}\text{EXP}_{\mathcal{I}} \\ &\text{EXEC}\text{C}\text{\#}\text{SHARP}\text{STM}_{\mathcal{I}} \end{aligned}$$

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 abstract 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. A function $label: Pos \rightarrow Label$ decorates nodes with the information which identifies the grammar rule associated to the node. In the example the *label* of the root node is *If*.



As a side effect the values of local variables in the memory are updated using two dynamic function $locals: Loc \rightarrow Adr$ and $mem: Adr \rightarrow SimpleValue \cup \{Undef\}$ to assign to local variables memory addresses where the values are stored. Since in $C\#_{\mathcal{I}}$ 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 $C\#_{\mathcal{I}}$. In $C\#_{\mathcal{I}}$ 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 $C\#$ and Java from the modelling point of view).

Statements can terminate normally or abruptly, where in $C\#_{\mathcal{I}}$ the reasons of abruption are from the set $Abr = Break \mid Continue \mid Goto(Lab)$, to be refined for the extended models. We use an auxiliary dynamic function $values: Pos \rightarrow Result$ to store intermediate evaluation results from the set $Result = SimpleValue \cup Abr \cup \{Undef, Norm\}$. For the initial state we assume

- $mem(i) = Undef$ for every $i \in Adr$
- $pos = \text{root position of the abstract syntax tree}$
- $locals(x) \in Adr$ for every variable x ³

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 instruction⁴. 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

³This comes up to assume that the compiler chooses an address for each variable.

⁴The specializations can be regained instruction-wise by mere parameter expansion, a form of refinement that is easily proved to be correct.

has to be matched against the syntactically possible cases (in the textual order of the rule) to select the appropriate computation step. If the subtree at the position pos is already evaluated and pos carries its result in $values$, then $context(pos)$ is the parent node of pos together with its children where pos is marked with the cursor (\blacktriangleright); otherwise, $context(pos)$ consists of the node at pos together with its children. 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 $C\#_{\mathcal{I}}$ the memory value of a given type t at a given address adr is defined by $memValue(adr, t) = mem(adr)$). This is described by the following two macros:

```
YIELD( $val, p$ )  $\equiv$ 
   $values(p) := val$ 
   $pos := p$ 
```

```
YIELDINDIRECT( $adr, p$ )  $\equiv$ 
  if  $AddressPos(p)$  then YIELD( $adr, p$ ) else YIELD( $memValue(adr, type(p)), p$ )
```

We will use the macros in the two forms $YIELD(val) \equiv YIELD(val, pos)$ and $YIELDUP(val) \equiv YIELD(val, up(pos))$, similarly for $YIELDINDIRECT(adr)$ and $YIELDUPINDIRECT(adr)$.

We are now ready to define the machine $EXEC\#SHARPEXP_{\mathcal{I}}$ 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 macro $WRITEMEM(adr, t, val)$ denotes here $mem(adr) := val$; it will be refined in the model for $C\#_{\mathcal{O}}$.

```
EXEC\#SHARPEXP_{\mathcal{I}}  $\equiv$  match  $context(pos)$ 
   $lit \rightarrow$  YIELD( $ValueOfLiteral(lit)$ )
   $loc \rightarrow$  YIELDINDIRECT( $locals(loc)$ )
   $uop\ exp \rightarrow$   $pos := exp$ 
   $uop \blacktriangleright val \rightarrow$  if  $\neg UopException(uop, val)$  then YIELDUP( $Apply(uop, val)$ )
   $exp_1\ bop\ exp_2 \rightarrow$   $pos := exp_1$ 
   $\blacktriangleright val\ bop\ exp \rightarrow$   $pos := exp$ 
   $val_1\ bop \blacktriangleright val_2 \rightarrow$  if  $\neg BopException(bop, val_1, val_2)$  then YIELDUP( $Apply(bop, val_1, val_2)$ )
   $exp_0\ ?\ exp_1 : exp_2 \rightarrow$   $pos := exp_0$ 
   $\blacktriangleright val\ ?\ exp_1 : exp_2 \rightarrow$  if  $val$  then  $pos := exp_1$  else  $pos := exp_2$ 
   $True\ ?\ \blacktriangleright val : exp \rightarrow$  YIELDUP( $val$ )
   $False\ ?\ exp : \blacktriangleright val \rightarrow$  YIELDUP( $val$ )
   $loc = exp \rightarrow$   $pos := exp$ 
   $loc = \blacktriangleright val \rightarrow$  {WRITEMEM( $locals(loc), type(loc), val$ ), YIELDUP( $val$ )}
  ( $type$ )  $exp \rightarrow$   $pos := exp$ 
  ( $type$ )  $\blacktriangleright val \rightarrow$  if  $\neg UopException(type, val)$  then YIELDUP( $Convert(type, val)$ )
   $vexp\ op = exp \rightarrow$   $pos := vexp$ 
   $\blacktriangleright adr\ op = exp \rightarrow$   $pos := exp$ 
   $adr\ op = \blacktriangleright val \rightarrow$  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  $result = Convert(t, w)$  in
        {WRITEMEM( $adr, t, result$ ), YIELDUP( $result$ )}
   $vexp\ op \rightarrow$   $pos := vexp$  // where  $op \in \{++, --\}$ 
   $\blacktriangleright adr\ op \rightarrow$  let  $old = memValue(adr, type(pos))$  in
    if  $\neg UopException(op, old)$  then
      {WRITEMEM( $adr, type(up(pos)), Apply(op, old)$ ), YIELDUP( $old$ )}
  checked( $exp$ )  $\rightarrow$   $pos := exp$ 
  checked( $\blacktriangleright val$ )  $\rightarrow$  YIELDUP( $val$ )
  unchecked( $exp$ )  $\rightarrow$   $pos := exp$ 
  unchecked( $\blacktriangleright val$ )  $\rightarrow$  YIELDUP( $val$ )
```


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

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

Similarly for being in a checked context which is used to define whether operators throw an overflow exception (in which case a rule will be added in the model for $C\#\varepsilon$). 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.

$$\begin{aligned} \text{Checked}(\alpha) &\iff \text{label}(\alpha) = \text{Checked} \vee \\ &\quad (\text{label}(\alpha) \neq \text{Unchecked} \wedge \text{up}(\alpha) \neq \text{Undef} \wedge \text{Checked}(\text{up}(\alpha))) \\ \text{UopException}(\text{uop}, \text{val}) &\iff \text{Checked}(\text{pos}) \wedge \text{Overflow}(\text{uop}, \text{val}) \\ \text{BopException}(\text{bop}, \text{val}_1, \text{val}_2) &\iff \\ &\quad \text{DivisonByZero}(\text{bop}, \text{val}_2) \vee \text{DecimalOverflow}(\text{bop}, \text{val}_1, \text{val}_2) \vee \\ &\quad (\text{Checked}(\text{pos}) \wedge \text{Overflow}(\text{bop}, \text{val}_1, \text{val}_2)) \end{aligned}$$

Similarly, the machine $\text{EXEC}\text{CSHARP}\text{STM}_I$ 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 below in such a way that in the refined models it is easily extended to abruptions due to return from procedures or to exceptions. In case of a new execution of the body of a while statement, the previously computed intermediate results have to be cleared. For the sake of brevity we skip the analogous transition rules for statements `do`, `for`, `switch`, `goto case`, `goto default`.

$$\begin{aligned} \text{EXEC}\text{CSHARP}\text{STM}_I &\equiv \text{match } \text{context}(\text{pos}) \\ &\quad ; \rightarrow \text{YIELD}(\text{Norm}) \\ \text{exp}; &\rightarrow \text{pos} := \text{exp} \\ \blacktriangleright \text{val}; &\rightarrow \text{YIELDUP}(\text{Norm}) \\ \text{break}; &\rightarrow \text{YIELD}(\text{Break}) \\ \text{continue}; &\rightarrow \text{YIELD}(\text{Continue}) \\ \text{goto } \text{lab}; &\rightarrow \text{YIELD}(\text{Goto}(\text{lab})) \\ \text{if } (\text{exp}) \text{ } \text{stm}_1 \text{ else } \text{stm}_2 &\rightarrow \text{pos} := \text{exp} \\ \text{if } (\blacktriangleright \text{val}) \text{ } \text{stm}_1 \text{ else } \text{stm}_2 &\rightarrow \text{if } \text{val} \text{ then } \text{pos} := \text{stm}_1 \text{ else } \text{pos} := \text{stm}_2 \\ \text{if } (\text{True}) \blacktriangleright \text{Norm} \text{ else } \text{stm} &\rightarrow \text{YIELDUP}(\text{Norm}) \\ \text{if } (\text{False}) \text{ } \text{stm} \text{ else } \blacktriangleright \text{Norm} &\rightarrow \text{YIELDUP}(\text{Norm}) \\ \text{while } (\text{exp}) \text{ } \text{stm} &\rightarrow \text{pos} := \text{exp} \\ \text{while } (\blacktriangleright \text{val}) \text{ } \text{stm} &\rightarrow \text{if } \text{val} \text{ then } \text{pos} := \text{stm} \text{ else } \text{YIELDUP}(\text{Norm}) \\ \text{while } (\text{True}) \blacktriangleright \text{Norm} &\rightarrow \{\text{pos} := \text{up}(\text{pos}), \text{CLEARVALUES}(\text{up}(\text{pos}))\} \\ \text{while } (\text{True}) \blacktriangleright \text{Break} &\rightarrow \text{YIELDUP}(\text{Norm}) \\ \text{while } (\text{True}) \blacktriangleright \text{Continue} &\rightarrow \{\text{pos} := \text{up}(\text{pos}), \text{CLEARVALUES}(\text{up}(\text{pos}))\} \\ \text{while } (\text{True}) \blacktriangleright \text{abr} &\rightarrow \text{YIELDUP}(\text{abr}) \\ \text{type } \text{loc}; &\rightarrow \text{YIELD}(\text{Norm}) \\ \text{lab: } \text{stm} &\rightarrow \text{pos} := \text{stm} \\ \text{lab: } \blacktriangleright \text{Norm} &\rightarrow \text{YIELDUP}(\text{Norm}) \\ \text{checked } \text{block} &\rightarrow \text{pos} := \text{block} \\ \text{checked } \blacktriangleright \text{Norm} &\rightarrow \text{YIELDUP}(\text{Norm}) \\ \text{unchecked } \text{block} &\rightarrow \text{pos} := \text{block} \\ \text{unchecked } \blacktriangleright \text{Norm} &\rightarrow \text{YIELDUP}(\text{Norm}) \\ \dots \blacktriangleright \text{abr} \dots &\rightarrow \text{if } \text{up}(\text{pos}) \neq \text{Undef} \wedge \text{PropagatesAbr}(\text{up}(\text{pos})) \text{ then } \text{YIELDUP}(\text{abr}) \end{aligned}$$

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

```

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

```

The auxiliary macro CLEARVALUES(α) to clear all values in the subtree at position α can be defined by recursion as follows, proceeding from top to bottom and from left to right⁵:

```

CLEARVALUES( $\alpha$ )  $\equiv$ 
  values( $\alpha$ ) := 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 Refining $C\#\mathcal{I}$ by static class features

In this section we refine the imperative core $C\#\mathcal{I}$ to $C\#_c$ 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 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. $\text{PropagatesAbr}(\alpha)$, 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.

In $C\#_c$ 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 interfaces⁶) 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 [3]. 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 $C\#_c$. The new set *Arg* of arguments appearing here reflects that besides value parameters also **ref** and **out** parameters can be used.

```

Vexp ::= ... | Field | Class ‘.’ Field
Sexp ::= ... | Meth ( [Args] ) | Class ‘.’ Meth ( [Args] )
Arg  ::= Exp | ‘ref’ Vexp | ‘out’ Vexp
Args ::= Arg { ‘,’ Arg }
Stm  ::= ... | ‘return’ Exp ‘;’ | ‘return’ ‘;’

```

⁵Intuitively 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 [2] for a precise definition.

⁶Note that struct and enum types and delegates are introduced by further refinement steps below.

The type constraints for the new expressions and the return statement are spelt out in [3]. 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.

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):

$$\begin{aligned}
Meth &= Type \times Msig \\
TypeState &= Linked \mid InProgress \mid Initialized \\
Frame &= Meth \times Pos \times Loc \times Values, \quad \text{where } Values = (Pos \rightarrow Result)
\end{aligned}$$

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:

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

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:

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

We extend the stipulations for the initial state as follows:

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

The submachine $\text{EXEC}\text{C}\text{SHARP}_C$ extends the machine $\text{EXEC}\text{C}\text{SHARP}_I$ for $\text{C}\#\mathcal{I}$ by additional rules for the evaluation of the new expressions and for the execution of return statements. In the same way the further refinements of $\text{EXEC}\text{C}\text{SHARP}$ in the sections below consist only in the parallel addition of submachines $\text{EXEC}\text{C}\text{SHARP}_O$, $\text{EXEC}\text{C}\text{SHARP}_E$, $\text{EXEC}\text{C}\text{SHARP}_D$, $\text{EXEC}\text{C}\text{SHARP}_T$, and $\text{EXEC}\text{C}\text{SHARP}_U$.

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

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

$$\begin{aligned} \text{EXEC}\text{C}\text{SHARP}\text{EXP}_C &\equiv \mathbf{match} \text{ context}(pos) \\ c.f \rightarrow &\mathbf{if} \text{ Initialized}(c) \mathbf{then} \text{ YIELD}\text{INDIRECT}(\text{globals}(c::f)) \mathbf{else} \text{ INITIALIZE}(c) \\ c.f = \text{exp} &\rightarrow pos := \text{exp} \\ c.f = \blacktriangleright \text{val} &\rightarrow \mathbf{if} \text{ Initialized}(c) \mathbf{then} \\ &\quad \text{WRITEMEM}(\text{globals}(c::f), \text{type}(c::f), \text{val}) \\ &\quad \text{YIELDUP}(\text{val}) \\ &\mathbf{else} \text{ INITIALIZE}(c) \\ c.m(\text{args}) &\rightarrow pos := (\text{args}) \\ c.m \blacktriangleright (\text{vals}) &\rightarrow \text{INVOKE}\text{STATIC}(c::m, \text{vals}) \\ \mathbf{ref} \text{ vexp} &\rightarrow pos := \text{vexp} \\ \mathbf{ref} \blacktriangleright \text{adr} &\rightarrow \text{YIELDUP}(\text{adr}) \\ \mathbf{out} \text{ vexp} &\rightarrow pos := \text{vexp} \\ \mathbf{out} \blacktriangleright \text{adr} &\rightarrow \text{YIELDUP}(\text{adr}) \\ () &\rightarrow \text{YIELD}([\] \\ (\text{arg}, \dots) &\rightarrow pos := \text{arg} \\ (\text{val}_1, \dots, \blacktriangleright \text{val}_n) &\rightarrow \text{YIELDUP}([\text{val}_1, \dots, \text{val}_n]) \\ (\dots \blacktriangleright \text{val}, \text{arg} \dots) &\rightarrow pos := \text{arg} \end{aligned}$$

The macro $\text{INVOKE}\text{STATIC}$ invokes the method – if the class is initialized, otherwise it initializes the class. For methods which are not declared external, $\text{INVOKE}\text{METHOD}$ 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 $\text{C}\#$ – we provide here for their invocation a submachine $\text{INVOKE}\text{EXTERN}$, to be defined separately depending on the class of external (e.g. library) methods⁷.

$$\begin{aligned} \text{INVOKE}\text{STATIC}(c::m, \text{vals}) &\equiv \\ &\mathbf{if} \text{ Initialized}(c) \mathbf{then} \text{ INVOKE}\text{METHOD}(c::m, \text{vals}) \mathbf{else} \text{ INITIALIZE}(c) \\ \text{INVOKE}\text{METHOD}(c::m, \text{vals}) &\equiv \\ &\mathbf{if} \text{ extern} \in \text{modifiers}(c::m) \mathbf{then} \text{ INVOKE}\text{EXTERN}(c::m, \text{vals}) \\ &\mathbf{else} \mathbf{let} p = \mathbf{if} \text{ StaticCtor}(c::m) \mathbf{then} pos \mathbf{else} up(pos) \mathbf{in} \\ &\quad \text{frames} := \text{push}(\text{frames}, (\text{meth}, p, \text{locals}, \text{values})) \\ &\quad \text{meth} := c::m \\ &\quad pos := \text{body}(c::m) \\ &\quad \text{values} := \emptyset \\ &\quad \text{INIT}\text{LOCALS}(c::m, \text{vals}) \end{aligned}$$

⁷For an illustration of this use of external methods see below the model for delegates.

We remind the reader that 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 [9] and [5, 2.4.4] for a justification of this assumption. See also the model for $C\#_C$ where we provide a complete 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 and out arguments
    locals( $x$ ) := vals( $paramIndex(c::m, x)$ )

```

The rules for method return in $EXEC\text{SHARPSTM}_C$ trigger an abruptio upon returning from a method, resulting in the execution of EXITMETHOD .

```

EXEC\text{SHARPSTM}_C  $\equiv$  match context( $pos$ )
  return  $exp$ ;  $\rightarrow pos := exp$ 
  return  $\blacktriangleright val$ ;  $\rightarrow \text{YIELDUP}(Return(val))$ 
  return;  $\rightarrow \text{YIELD}(Return)$ 
  Return  $\rightarrow$  if  $pos = body(meth) \wedge \neg Empty(frames)$  then  $\text{EXITMETHOD}(Norm)$ 
  Return( $val$ )  $\rightarrow$  if  $pos = body(meth) \wedge \neg Empty(frames)$  then  $\text{EXITMETHOD}(val)$ 
   $\blacktriangleright Norm$ ;  $\rightarrow \text{YIELDUP}(Norm)$ 

```

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.⁸

```

EXITMETHOD( $result$ )  $\equiv$ 
  let ( $oldMeth, oldPos, oldLocals, oldValues$ ) = top( $frames$ ) in
    meth :=  $oldMeth$ 
    pos :=  $oldPos$ 
    locals :=  $oldLocals$ 
    frames := pop( $frames$ )
    if  $StaticCtor(meth) \wedge result = Norm$  then
      typeState( $type(meth)$ ) :=  $Initialized$ 
      values :=  $oldValues$ 
    else
      values :=  $oldValues \oplus \{oldPos \mapsto result\}$ 
  FREELOCALS

```

```

FREELOCALS  $\equiv$ 
  forall  $x \in LocalVars(meth) \cup ValueParams(meth)$  do
    FREEMEMORY( $locals(x), type(x)$ )

```

Following [7, §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:

$$Initialized(c) \iff typeState(c) = Initialized \vee typeState(c) = InProgress$$

⁸Under 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 .

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

Note that the initialization of a class does not trigger the initialization of its direct base class (as it is the case for Java).

With respect to the execution of initializers of static class fields the Ecma standard [7, §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. We do not model the last behavior, since Microsoft’s $C\#$ compiler currently creates a static constructor in this case.

4 Refinement $C\#\mathcal{O}$ of $C\#c$ by object related features

In this section we refine the static class features of $C\#c$ by adding objects (for class instances, comprising arrays and structs) together with *instance* fields, methods and constructors⁹ as well as inheritance (including overriding and overloading of methods). Accordingly we extend the ASM universes and functions of $C\#c$ 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 $C\#\mathcal{O}$. 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 [3].

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 $C\#\mathcal{D}$ to comprise also delegates. Value types represent a feature that distinguishes $C\#$ from Java. A *RefExp* is an expression of a reference type and an *ArrayExp* is an expression of an array type.

```
Exp ::= ... | ‘null’ | ‘this’ | ‘typeof’ (‘ RetType ‘) | Exp ‘is’ Type | Exp ‘as’ RefType
      | (‘ Type ‘) Exp | ‘new’ NonArrayType [‘ Exps ‘] {Rank} [ArrayInitializer]
Vexp ::= ... | Vexp ‘.’ Field | RefExp ‘.’ Field | ‘base’ ‘.’ Field | ArrayExp [‘ Exps ‘]
Sexp ::= ... | ‘new’ Type ( [Args] ) | Exp ‘.’ Meth ( [Args] ) | ‘base’ ‘.’ Meth ( [Args] )
Exps ::= Exp {‘,’ Exp}
Rank ::= [‘ { ‘,’ } ‘]’
```

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 reassumed by Fig. 4. We refer the reader to [3] 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 follows:

- T any type $\implies T \preceq \text{object}$ and $T \preceq T$
- class S derived from $T \implies S \preceq T$

⁹Destructors or finalizers which relate to garbage collection are not modeled here.

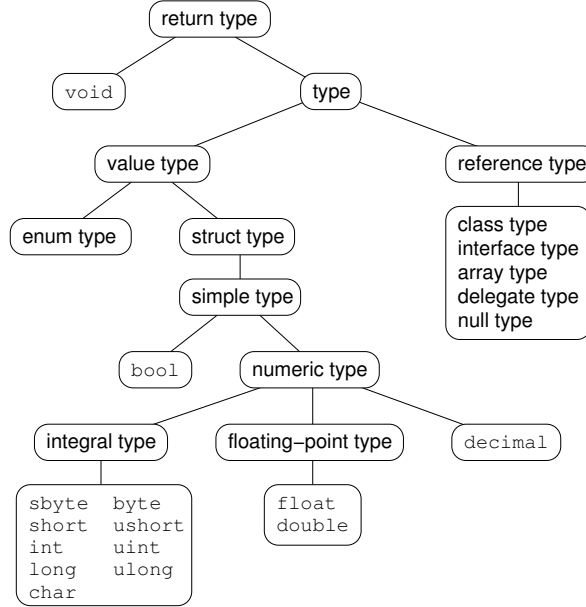


Figure 4: The classification of types of C#.

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

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 abstract syntax tree is annotated with 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, certain expressions are reduced to basic expressions at compile-time.

For the base access of fields and methods we have:

- $\text{base}.f$ in class C is replaced by $\text{this}.B::f$, where B is the first base class of C where a field f is declared.
- $\text{base}.m(args)$ in class C is replaced by $\text{this}.B::M(args)$, 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 [3], formalizing the conditions stated in [7, §14.4.2/3].

For instance field access and class instance creation we have:

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

Hence we split an instance creation expression into a creation part and an invocation of an instance constructor. We assume that class instance constructors return the value of **this**.

Instance constructors of structs need an address for **this**.

- Let S be a struct type. Then $vexp = \text{new } S::M(args)$ is replaced by $vexp.S::M(args)$.

- Otherwise, $\mathbf{new} S::M(args)$ is replaced by $x.S::M(args)$, where x is a new temporary local variable of type S . We assume that constructors of structs return the value of \mathbf{this} .

For automatic boxing we have:

- $vexp = exp$ 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 $type(exp) \preceq T$.
- 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 $type(arg) \preceq T$.

We are now ready to describe the extension of the dynamic state for the model of $C\#_O$. The domain of values is extended to contain also references (assuming $Ref \cap Adr = \emptyset$) 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. $structField(val, T, f) = val(f)$.

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. 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 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 $runTimeType(ref) = c$ 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. \mathbf{this} is treated as first parameter and is passed by value. Therefore $paramIndex(c::m, \mathbf{this}) = 0$ and \mathbf{this} is element of both $LocalVars(c::m)$ and $ValueParams(c::m)$.

4.1 Transition rules for $C\#_O$

For the refinement of the EXECCSHARP transition rules it suffices to add the EXECCSHARP_O rules, defined by two submachines to evaluate the new expressions and to execute the new statements respectively:

$$\begin{aligned} EXECCSHARP_O &\equiv \\ &EXECCSHARPEXP_O \\ &EXECCSHARPSTM_O \end{aligned}$$

EXECCSHARPEXP_O contains rules for each of the numerous forms of new expressions. For better readability we organize them 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). The rules below for calls of instance methods contain no class initialization test, since an instance method which is not an instance constructor can be called only when the corresponding class is already initialized.

$$\begin{aligned} EXECCSHARPEXP_O &\equiv \mathbf{match} \ context(pos) \\ &\mathbf{null} \rightarrow \mathbf{null} \\ &\mathbf{this} \rightarrow \mathbf{YIELDINDIRECT}(locals(\mathbf{this})) \\ &\mathbf{TESTCASTEXP}_O \\ &\mathbf{FIELDEXP}_O \end{aligned}$$


```

new c → let ref = new(Ref, c) in
  runTimeType(ref) := c
  forall f ∈ instanceFields(c) do
    let adr = fieldAdr(ref, f) and t = type(f) in
      WRITEMEM(adr, t, defaultValue(t))
  YIELD(ref)

exp . T::M(args) → pos := exp
► val . T::M(args) → if StructValueInvocation(up(pos)) then
  let adr = new(Adr, type(pos)) in // create home for struct value
  WRITEMEM(adr, type(pos), val)
  values(pos) := adr
  pos := (args)
val . T::M►(vals) → if InstanceCtor(M) ∧ ¬Initialized(T) then INITIALIZE(T)
  elseif val ≠ null then INVOKEINSTANCE(T::M, val, vals)

ARRAYEXPO

```

A struct value invocation is a method invocation on a struct value.

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

The rules for casting in TESTCASTEXP_O use the new macro YIELDUPBOX defined below. Note that in expressions ‘ exp 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 $(t)exp$ 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
    if t = type(pos) then YIELDUP(val) // compile-time identity
    if t ∈ RefType then YIELDUPBOX(type(pos), val) // box value
  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
      YIELDUP(memValue(valueAdr(val), t)) // un-box a copy of the value

```

The rules for instance field access and assignment in FIELDEXP_O 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. We use $\text{type}(exp . t::f) = \text{type}(t::f)$.

```

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 $dimLength(ref, i)$.

```

ARRAYEXPO ≡
new T[exp1, ..., expn][R1] ... [Rk] → pos := exp1
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, default Value(S))
            YIELDUP(ref)
exp0[exp1, ..., expn] → pos := exp0
► ref[exp1, ..., expn] → pos := exp1
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] = expn+1 → pos := exp0
► ref[exp1, ..., expn] = exp → pos := exp1
ref[i1, ..., ►in] = exp → pos := exp
ref[i1, ..., in] = ►val →
let T = elementTyperunTimeType(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)

```

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 run-time 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#_O the notion of reading from the memory is refined by extending the simple equation $memValue(adr, t) = mem(adr)$ of C#_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))

```

Also writing to memory is refined from $\text{WRITEMEM}(adr, t, val) \equiv \text{mem}(adr) := val$ in $\text{C}\#_{\mathcal{I}}$, 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
    for all 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 $\text{C}\#_{\mathcal{I}}$ is refined to include also lvalue nodes of *StructType*.

```

AddressPos(α) ⇔ FirstChild(α) ∧
  label(up(α)) ∈ {ref, out, ++, --} ∨ label(up(α)) ∈ Aop ∨
  (label(up(α)) = '.' ∧ α ∈ Vexp ∧ type(α) ∈ StructType)

```

Address positions are: $\text{ref } \square$, $\text{out } \square$, $\square++$, $\square--$, $\square \text{op} = \text{exp}$, $\square.f$, $\square.m(\text{args})$.

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 defined to be t . The struct is copied in both cases, when it is boxed and when it is un-boxed.

```

YIELDUPBOX(t, val) ≡ let ref = new(Ref) and adr = new(Adr, t) in
  runTimeType(ref) := t
  valueAdr(ref) := adr
  WRITEMEM(adr, t, val)
  YIELDUP(ref)

```

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

```

let adr = new(Adr, T) in P
let ref = new(Ref, T) in P
let ref = new(Ref, [l1, ..., ln], T) in P

```

In the context of the machine EXECCSHARP this comes up to specify an abstract memory management. In fact $\text{let } adr = \text{new}(Adr, T) \text{ in } P$ stands for the sequential execution of a new address allocation followed by P :

```

let adr = new(Adr, T) in P ≡ (import adr do ALLOCADR(adr, T)) seq P

```

where the operator seq for sequential execution of two ASMs M, N is to be understood as defined for turbo ASMs in [4] (alternatively see [5, 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 $\text{let } ref = \text{new}(Ref, T) \text{ in } P$ stands for the sequential execution of address allocation for all instance fields of a given type followed by P :

```

let ref = new(Ref, T) in P ≡
  import ref do
    Ref(ref) := True
    ALLOCFIELDS(ref, instanceFields(T))
  seq P

```

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

```

let  $ref = new(Ref, [l_1, \dots, l_n], T)$  in  $P \equiv$ 
  import  $ref$  do
     $Ref(ref) := True$ 
    forall  $\alpha \in [0..l_1 - 1] \times \dots \times [0..l_n - 1]$  do
      import  $adr$  do
         $elemAdr(ref, \alpha) := adr$ 
         $ALLOCADR(adr, T)$ 
    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 [2] (or see alternatively [5, Ch.4.1.2]):

```

 $ALLOCADR(adr, T) \equiv$ 
   $Adr(adr) := True$ 
  if  $T \in StructType$  then  $ALLOCFIELDS(adr, instanceFields(T))$ 

```

```

 $ALLOCFIELDS(x, fs) \equiv$ 
  forall  $f \in fs$  import  $adr$  do
     $fieldAdr(x, f) := adr$ 
     $ALLOCADR(adr, type(f))$ 

```

5 Refinement $C\#\varepsilon$ of $C\#\emptyset$ by exception handling

In this section we extend $C\#\emptyset$ 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 abrupted.

For the refinement of EXECCSHARP by exceptions, as in the previous section it suffices to add the rules for EXECCSHARP_E and to extend the static semantics. The set of statements is extended by throw and try-catch-finally statements satisfying the following constraints:

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

- 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 $i < j \implies E_j \not\preceq E_i$ (and obviously $E_i \preceq \text{System.Exception}$) holds for every try-catch statement `try block catch (E1 x1) block1 ... catch (En xn) blockn`

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.

$$Abr = \dots \mid Exc(Ref), \quad TypeState = \dots \mid Exc(Ref), \quad excStack: List(Ref)$$

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:

<pre> try TryBlock catch (E₁ x₁) CatchBlock₁ ⋮ catch (E_n x_n) CatchBlock_n finally FinallyBlock </pre>	\Longrightarrow	<pre> try { try TryBlock catch (E₁ x₁) CatchBlock₁ ⋮ catch (E_n x_n) CatchBlock_n } finally FinallyBlock </pre>
--	-------------------	--

Unhandled exceptions in a static constructor are wrapped into a `TypeInitializationException` by translating `static T() { BlockStatements }` into

```

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

```

For $stm \equiv \text{try } tryBlock \text{ catch } (\dots) catchBlock_1 \dots \text{catch } (\dots) catchBlock_n$ the reachability rules and the definite assignment constraints 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 $normal(stm)$.
- $before(tryBlock) = before(stm)$
- $before(catchBlock_i) = before(stm)$ 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 $normal(stm)$.
- $before(tryBlock) = before(stm)$
- $before(finallyBlock) = before(stm)$
- $after(stm) = after(tryBlock) \cup after(finallyBlock)$

5.1 Transition rules for $C\#\varepsilon$

The transition rules for $EXEC\text{SHARP}_E$ are defined by two submachines. 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.

```

EXEC\text{SHARP}_E \equiv
  EXEC\text{SHARP}\text{EXP}_E
  EXEC\text{SHARP}\text{STM}_E

```

$EXEC\text{SHARP}\text{EXP}_E$ 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 paramter of FAILUP.

```

EXEC\text{SHARP}\text{EXP}_E \equiv \text{match } context(pos)
  uop \blacktriangleright val \rightarrow \text{if } Checked(pos) \wedge Overflow(uop, val) \text{ then FAILUP(OverflowException)
  val_1 bop \blacktriangleright val_2 \rightarrow
    \text{if } DivisionByZero(bop, val_2) \text{ then FAILUP(DivideByZeroException)
    \text{elseif } DecimalOverflow(bop, val_1, val_2) \vee (Checked(pos) \wedge Overflow(bop, val_1, val_2))
    \text{then FAILUP(OverflowException)
CAST\text{EXCEPTIONS}
NULL\text{REF}\text{EXCEPTIONS}
ARRAY\text{EXCEPTIONS}

```

FAILUP(E) \equiv INVOKEMETHOD(*ExcSupport::ThrowE*, [])

CASTEXCEPTIONS \equiv **match** *context(pos)*

(t) \blacktriangleright *val* \rightarrow

if *type(pos)* \in *RefType* **then**

if $t \in \text{RefType} \wedge \text{val} \neq \text{Null} \wedge \text{runTimeType}(\text{val}) \not\leq t$ **then**

FAILUP(*InvalidCastException*)

if $t \in \text{ValueType}$ **then**

if $\text{val} = \text{Null}$ **then** FAILUP(*NullReferenceException*)

elseif $t \neq \text{runTimeType}(\text{val})$ **then** FAILUP(*InvalidCastException*)

if *type(pos)* \in *SimpleType* $\wedge t \in \text{SimpleType} \wedge \text{Checked}(\text{pos}) \wedge \text{Overflow}(t, \text{val})$

then FAILUP(*OverflowException*)

// attempt to unbox

NULLREFEXCEPTIONS \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 $T[l_1, \dots, \blacktriangleright l_n][R_1] \dots [R_k] \rightarrow$

if $\exists i \in [1..n] (l_i < 0)$ **then** FAILUP(*OverflowException*)

ref [$i_1, \dots, \blacktriangleright i_n$] \rightarrow

if *ref* = *Null* **then** FAILUP(*NullReferenceException*)

elseif $\exists k \in [1..n] (i_k < 0 \vee \text{dimLength}(\text{ref}, k-1) \leq i_k)$ **then**

FAILUP(*IndexOutOfRangeException*)

elseif $\text{RefOrOutArg}(\text{up}(\text{pos})) \wedge \text{type}(\text{up}(\text{pos})) \in \text{RefType} \wedge$

$\text{elementType}(\text{runTimeType}(\text{ref})) \neq \text{type}(\text{up}(\text{pos}))$

then FAILUP(*ArrayTypeMismatchException*)

ref [i_1, \dots, i_n] = \blacktriangleright *val* \rightarrow

if *ref* = *Null* **then** FAILUP(*NullReferenceException*)

elseif $\exists k \in [1..n] (i_k < 0 \vee \text{dimLength}(\text{ref}, k-1) \leq i_k)$ **then**

FAILUP(*IndexOutOfRangeException*)

elseif $\text{type}(\text{pos}) \in \text{RefType} \wedge \text{runTimeType}(\text{val}) \not\leq \text{elementType}(\text{runTimeType}(\text{ref}))$

then FAILUP(*ArrayTypeMismatchException*)

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. The semantics of **throw**; is explained in terms of the exception stack *excStack*. When an exception is caught, it is pushed on top of the exception stack. 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.

EXECSSHARPSTM_E \equiv **match** *context(pos)*

throw *exp*; \rightarrow *pos* := *exp*

throw \blacktriangleright *ref*; \rightarrow **if** *ref* = *Null* **then** FAILUP(*NullReferenceException*)

else {INITSTACKTRACE(*ref, meth*), YIELDUP(*Exc(ref)*)}

throw; \rightarrow YIELD(*Exc(top(excStack))*)

try *block* **catch** (*E x*) *stm* ... \rightarrow *pos* := *block*

try \blacktriangleright *Norm* **catch** (*E x*) *stm* ... \rightarrow YIELDUP(*Norm*)

try \blacktriangleright *Exc(ref)* **catch** (*E*₁ *x*₁) *stm*₁ ... **catch** (*E*_{*n*} *x*_{*n*}) *stm*_{*n*} \rightarrow

```

if  $\exists i \in [1..n]$   $runTimeType(ref) \preceq E_i$  then
  let  $j = \min\{i \in [1..n] \mid runTimeType(ref) \preceq E_i\}$  in
     $pos := stm_j$ 
     $excStack := push(ref, excStack)$ 
     $WRITEMEM(locals(x_j), object, ref)$ 
  else  $YIELDUP(Exc(ref))$ 
try  $\blacktriangleright$   $abr$  catch  $(E_1 x_1) stm_1 \dots \text{catch}(E_n x_n) stm_n \rightarrow YIELDUP(abr)$ 
try  $Exc(ref) \dots \text{catch}(\dots) \blacktriangleright res \dots \rightarrow \{excStack := pop(excStack), YIELDUP(res)\}$ 
try  $tryBlock$  finally  $finallyBlock \rightarrow pos := tryBlock$ 
try  $\blacktriangleright res$  finally  $finallyBlock \rightarrow pos := finallyBlock$ 
try  $res$  finally  $\blacktriangleright Norm \rightarrow YIELDUP(res)$ 
try  $res$  finally  $\blacktriangleright Exc(ref) \rightarrow YIELDUP(Exc(ref))$ 
 $Exc(ref) \rightarrow$  if  $pos = body(meth) \wedge \neg Empty(frames)$  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#c` re-throws the exception object of a type which had an exception in the static constructor, thus preventing its re-initialization.

```

INITIALIZE( $c$ )  $\equiv$ 
  ...
  if  $typeState(c) = Exc(ref)$  then  $YIELD(Exc(ref))$ 

```

6 Refinement `C#D` of `C# ε` by delegates

In this section we extend `C# ε` 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.

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 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 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 [3].

We use the model EXECCSHARPSTM_I, which includes a description of the `for` statement of C#_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 .¹⁰ It then suffices to refine the ASM rule INVOKEEXTERN defined in EXECCSHARPEXP_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 analogously of `_length`.

$$\begin{aligned}
\text{INVOKEEXTERN}(T::M, vals) & \equiv \\
& \text{if } T \in \text{DelegateType} \text{ then} \\
& \quad \text{if } name(M) = _length \text{ then DELEGATELENGTH}(vals(0)) \\
& \quad \text{if } name(M) = _invoke \text{ then INVOKEDELEGATE}(vals(0), vals(1), drop(vals, 2)) \\
\text{DELEGATELENGTH}(ref) & \equiv \\
& \text{YIELDUP}(length(invocationList(ref))) \\
\text{INVOKEDELEGATE}(ref, i, vals) & \equiv \\
& \text{match } invocationList(ref)(i) \\
& \quad T::M \quad \rightarrow \text{INVOKESTATIC}(T::M, vals) \\
& \quad (target, T::M) \rightarrow \text{INVOKEINSTANCE}(T::M, target, vals)
\end{aligned}$$

Since there are no new statements appearing in C#_D, the addition of EXECCSHARP_D consists in the following ASM EXECCSHARPEXP_D, which defines the meaning of delegate instance creation.

$$\begin{aligned}
\text{EXECCSHARPEXP}_D & \equiv \text{match } context(pos) \\
& \text{new } D(T::M) \rightarrow \\
& \quad \text{let } d = \text{new}(Ref, D) \text{ in} \\
& \quad \quad runTimeType(d) := D \\
& \quad \quad invocationList(d) := [T::M] \\
& \quad \quad \text{YIELD}(d) \\
& \text{new } D(exp.T::M) \rightarrow pos := exp \\
& \text{new } D(\blacktriangleright ref.T::M) \rightarrow \\
& \quad \text{if } ref = \text{Null} \text{ then FAILUP}(\text{NullReferenceException}) \\
& \quad \text{else let } d = \text{new}(Ref, D) \text{ in} \\
& \quad \quad runTimeType(d) := D \\
& \quad \quad invocationList(d) := [(ref, T::M)]
\end{aligned}$$

¹⁰In [7, §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.


```

    YIELDUP(d)
new D(exp) → pos := exp
new D(► ref) →
  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(S1, ..., Sn))]
    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, and they can be *combined* (concatenated with ‘+’) and elements 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 *suffix*(*u*, *v*) the part of *v* after the last occurrence of *u* in *v*.

```

DELEGATEEQUAL(r1, r2) ≡
  if r1 = Null ∨ r2 = Null then YIELDUP(r1 = r2)
  else let l1 = invocationList(r1) and l2 = invocationList(r2) in
    YIELDUP(length(l1) = length(l2) ∧ ∀ i < length(l1) (l1(i) = l2(i)))

```

6.2 Properties, events and further features in C#_D

In this section we add further language features of C# whose semantics can be easily defined in terms of the model developed so far, essentially by simple syntactical reductions.

Properties. Collections of a read and/or a write method for attributes of a class or struct are called *properties* in C# 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 identifier¹¹, 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 **get**_ P () ; and **void** **set**_ P (T **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:

$$\begin{array}{l}
 T \text{ } P \{ \\
 \quad \text{get} \{ \text{getAccessor} \} \\
 \quad \text{set} \{ \text{setAccessor} \} \\
 \} \\
 \\
 \text{exp} . P \implies \text{exp} . \text{get} _ P () \quad \text{exp}_1 . P = \text{exp}_2 ; \implies \text{exp}_1 . \text{set} _ P (\text{exp}_2) ;
 \end{array}
 \implies
 \begin{array}{l}
 T \text{ get} _ P () \{ \\
 \quad \text{getAccessor} \\
 \} \\
 \text{void set} _ P (T \text{ value}) \{ \\
 \quad \text{setAccessor} \\
 \}
 \end{array}$$

This translation comprises also expressions of the form $\text{exp}_1 . P \text{ op} = \text{exp}_2$, since they can be assumed to be compiled to $\langle x = \text{exp}_1, y = x . \text{get} _ P () \text{ op} \text{exp}_2, x . \text{set} _ P (y), y \rangle$ 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:

$$\begin{array}{l}
 \langle \text{exp}, \dots \rangle \rightarrow \text{pos} := \text{exp} \\
 \langle \text{val}_1, \dots, \blacktriangleright \text{val}_n \rangle \rightarrow \text{YIELDUP}(\text{val}_n) \\
 \langle \dots \blacktriangleright \text{val}, \text{exp} \dots \rangle \rightarrow \text{pos} := \text{exp}
 \end{array}$$

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 properties is reduced to the semantics of arrays and methods via the following compile-time translation (and corresponding operator expression translation as explained for properties):

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

Events. Events can be declared in C# like fields, in the form ‘**event**’ *DelegateType Identifier* ‘;’ (we omit the modifiers), or like properties, in the form

‘**event**’ *DelegateType Identifier* ‘{’ ‘**add**’ *Block* ‘**remove**’ *Block* ‘}’.

¹¹Without knowing whether it is accessed directly or whether an accessor method is being called.

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#* 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_Args 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_Args 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_Args e) { if (X != null) X(this, e); }`.

It suffices to assign a meaning to `void add_X(D value)` and `void remove_X(D value)`, which are reserved signatures 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#_T* which is explained in [6]¹².

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

Further constructs. For similar syntactical reductions to those given above, which 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 see [3].

7 Refinement *C#_u* by pointers in unsafe code

In this section we add the features *C#* 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. This 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. As an alternative to pinning, data of unmanaged type can also be ‘`stackalloc`’ed, instead of using the heap.

Java has no such unsafe extension. The refinement consists mainly in a definition of the *memory* function in terms of byte sequences, using an encoding of simple types and a corresponding refinement of the function *structField*.

¹²If one prefers not to depend on the thread model *C#_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#_T* for the multiple thread execution model.

7.1 Signature refinement for C# \mathcal{U}

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

$$\textit{Type} ::= \textit{ValueType} \mid \textit{RefType} \mid \textit{PointerType} \quad \textit{PointerType} ::= \textit{UnmanagedType} \textit{'*'} \mid \textit{'void'} \textit{'*'}$$

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 `'stackalloc'` 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.

$$\begin{aligned} \textit{Exp} & ::= \dots \mid \textit{'\&'} \textit{Vexp} \mid \textit{Exp} \textit{'->'} \textit{Meth} ([\textit{Args}]) \mid \textit{Exp} \textit{'->'} \textit{Field} \\ & \quad \mid \textit{'sizeof'} \textit{'('} \textit{UnmanagedType} \textit{'\>'} \textit{'\>'} \\ \textit{Vexp} & ::= \dots \mid \textit{'*' } \textit{Exp} \\ \textit{Stm} & ::= \dots \mid \textit{'unsafe'} \textit{Block} \mid \textit{'fixed'} \textit{'('} \textit{PointerType} \textit{Loc} = \textit{Exp} \textit{'\>'} \textit{'\>'} \textit{Stm} \\ \textit{Bstm} & ::= \dots \mid \textit{PointerType} \textit{Loc} \textit{'='} \textit{'stackalloc'} \textit{UnmanagedType} \textit{'['} \textit{Exp} \textit{'\>'} \textit{'\>'} \textit{'\>'} \textit{'\>'} \end{aligned}$$

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):

- $++p, --p, p++, p--, p + i, i + p, p - i, p - q, p == q, p != q, p < q, p <= q, p > q, p >= q$
- $(T^*)i, (T^*)p, (\textit{int})p, (\textit{uint})p, (\textit{long})p, (\textit{ulong})p$

On the types of the new expressions the following constraints are imposed.

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

We assume `e -> m` to be translated to `(*e).m` and `e[i]` to `*e + i`.

For statements the following type constraints are assumed:

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

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. Hence, the garbage collection can determine the pinned objects by looking at the values of pinned local variables on the frames stack.

The principal refinement in the ASM extension EXECCSHARP \mathcal{U} for C# \mathcal{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 ($\textit{Adr} = \mathbb{N}$):

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

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 $sizeOf(T)$, satisfy for values val the equations

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

For every pointer type T^* holds $sizeOf(T^*) = sizeOf(void^*)$.

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

$$fieldOffset(T, f) + sizeOf(type(f)) \leq sizeOf(T)$$

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

To determine the layout of arrays with unmanaged element type we stipulate the following refinement of the function $elemAdr$ 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 $runTimeType(ref) = T[dim(n)]$, where T is an unmanaged type and $l_i = dimLength(ref, i - 1)$ for $i \in [1..n]$, we assume the following:

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

7.2 Transition rule refinement for unsafe code

Besides the rules below which define the semantics of the new expressions and statements we have to refine the notions of reading from and writing to memory for values of unmanaged type.

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

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

$$\begin{aligned} WRITEMEM(adr, t, val) \equiv & \\ & \text{if } t \in RefType \text{ then } mem(adr) := val \\ & \text{elseif } t \in UnmanagedType \text{ then} \\ & \quad \text{forall } i \in [0..sizeOf(t) - 1] \text{ do } mem(adr + i) := val(i) \\ & \text{elseif } t \in StructType \text{ then} \\ & \quad \text{forall } f \in instanceFields(t) \text{ do } SETFIELD(adr, f, val(f)) \end{aligned}$$

$$SETFIELD(adr, t::f, val) \equiv WRITEMEM(fieldAdr(adr, t::f), type(t::f), val)$$

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:

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

In the rules for EXEC $C_{SHARPEXP}_U$ 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.

$$\begin{aligned} EXEC_{SHARPEXP}_U \equiv & \text{match } context(pos) \\ & sizeof(T) \rightarrow YIELD(sizeOf(T)) \\ & \&exp \rightarrow pos := exp \\ & \&\blacktriangleright adr \rightarrow YIELDUP(adr) \end{aligned}$$

```

*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
fixed (char* p = exp) stm	{ char* p; p = Cstring(exp); stm }
fixed (T* p = exp) stm	{ T* p; p = &exp[0]; stm }
fixed (T* p = &vexp) stm	{ 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 [7].

References

- [1] T. Archer and A. Whitechapel. *Inside C#*. Microsoft Press, 2002.
- [2] E. Börger and T. Bolognesi. Remarks on turbo ASMs for computing functional equations and recursion schemes. In E. Börger, A. Gargantini, and E. Riccobene, editors, *Abstract State Machines 2003 – Advances in Theory and Applications*, volume 2589 of *Lecture Notes in Computer Science*, pages 218–228. Springer-Verlag, 2003.
- [3] E. Börger, G. Fruja, V. Gervasi, and R. Stärk. A complete formal definition of the semantics of C#. Technical report, In preparation, 2003.
- [4] E. Börger and J. Schmid. Composition and submachine concepts for sequential ASMs. In P. Clote and H. Schwichtenberg, editors, *Computer Science Logic (Proceedings of CSL 2000)*, volume 1862 of *Lecture Notes in Computer Science*, pages 41–60. Springer-Verlag, 2000.
- [5] E. Börger and R. F. Stärk. *Abstract State Machines. A Method for High-Level System Design and Analysis*. Springer, 2003.
- [6] E. Börger and R. F. Stärk. An ASM model for C# threads. Technical report, Computer Science Department, ETH Zürich, 2003.

- [7] C# Language Specification. Standard ECMA-334, 2001.
<http://www.ecma-international.org/publications/standards/ECMA-334.HTM>.
- [8] Foundations of Software Engineering Group, Microsoft Research. AsmL.
Web pages at <http://research.microsoft.com/foundations/AsmL/>, 2001.
- [9] N. G. Fruja and R. F. Stärk. The hidden computation steps of turbo Abstract State Machines. In E. Börger, A. Gargantini, and E. Riccobene, editors, *Abstract State Machines 2003 – Advances in Theory and Applications*, volume 2589 of *Lecture Notes in Computer Science*, pages 244–262. Springer-Verlag, 2003.
- [10] A. Hejlsberg, S. Wiltamuth, and P. Golde. *C# Language Specification*. Addison-Wesley, 2003.
- [11] J. Prosise. *Programming Microsoft .NET*. Microsoft Press, 2002.
- [12] J. Richter. *Applied Microsoft .NET Framework Programming*. Microsoft Press, 2002.
- [13] J. Schmid. Executing ASM specifications with AsmGofer.
Web pages at <http://www.tydo.de/AsmGofer>.
- [14] J. Schmid. *Refinement and Implementation Techniques for Abstract State Machines*. PhD thesis, University of Ulm, Germany, 2002.
- [15] R. F. Stärk, J. Schmid, and E. Börger. *Java and the Java Virtual Machine—Definition, Verification, Validation*. Springer-Verlag, 2001.