

A MiniTriangle Grammars

This appendix contains the grammars that define the concrete and abstract syntax of the version of MiniTriangle used for this dissertation. The concrete syntax is divided into two parts: lexical syntax and context-free syntax. The grammars are derived from the book by Watt & Brown. However, the multithreading feature developed in this dissertation is not included in this appendix.

A.1 MiniTriangle Lexical Syntax

Non-terminals are typeset in italics, like *this*. Terminals are typeset in type-writer font, like **this**. Terminals whose spelling (the concrete character sequence) is different from what is shown in the grammar, such as names of special characters, are typeset in italics and underlined, like this. For simplicity, we resort to a slightly informal way of stating that the keywords are not valid identifiers.

<i>Program</i>	$\rightarrow (Token \mid Separator)^*$
<i>Token</i>	$\rightarrow Keyword \mid Identifier \mid IntegerLiteral \mid Operator$ $ (\mid) \mid [\mid] \mid \{ \mid \} \mid , \mid . \mid ; \mid : \mid := \mid = \mid \underline{eot}$
<i>Keyword</i>	$\rightarrow begin \mid const \mid do \mid else \mid end \mid fun \mid if \mid in$ $ let \mid out \mid proc \mid then \mid var \mid while$
<i>Identifier</i>	$\rightarrow Letter \mid Identifier\ Letter \mid Identifier\ Digit$ except <i>Keyword</i>
<i>IntegerLiteral</i>	$\rightarrow Digit \mid IntegerLiteral\ Digit$
<i>Operator</i>	$\rightarrow \sim \mid * \mid / \mid + \mid - \mid < \mid \leq \mid == \mid != \mid \geq \mid > \mid \&\& \mid \mid \mid !$
<i>Letter</i>	$\rightarrow A \mid B \mid \dots \mid Z \mid a \mid b \mid \dots \mid z$
<i>Digit</i>	$\rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$
<i>Separator</i>	$\rightarrow Comment \mid \underline{space} \mid \underline{eol}$
<i>Comment</i>	$\rightarrow // (\text{any character except } \underline{eol})^* \underline{eol}$

A.2 MiniTriangle Context-Free Syntax

Non-terminals are typeset in italics, like *this*. Terminals are typeset in type-writer font, like `this`. Terminals whose spelling (the concrete character sequence) is different from what is shown in the grammar are typeset in italics and underlined, like *Identifier* and *IntegerLiteral*. Their spelling is defined by the lexical grammar (where they are non-terminals!).

<i>Program</i>	\rightarrow	<i>Command</i>
<i>Commands</i>	\rightarrow	<i>Command</i>
		<i>Command</i> ; <i>Commands</i>
<i>Command</i>	\rightarrow	<i>VarExpression</i> := <i>Expression</i>
		<i>VarExpression</i> (<i>Expressions</i>)
		<i>if Expression then Command</i>
		<i>else Command</i>
		<i>while Expression do Command</i>
		<i>let Declarations in Command</i>
		<i>begin Commands end</i>
<i>Expressions</i>	\rightarrow	ϵ
		<i>Expressions</i> ₁
<i>Expressions</i> ₁	\rightarrow	<i>Expression</i>
		<i>Expression</i> , <i>Expressions</i> ₁
<i>Expression</i>	\rightarrow	<i>PrimaryExpression</i>
		<i>Expression BinaryOperator Expression</i>
<i>PrimaryExpression</i>	\rightarrow	<u><i>IntegerLiteral</i></u>
		<u><i>VarExpression</i></u>
		<i>UnaryOperator PrimaryExpression</i>
		<i>VarExpression</i> (<i>Expressions</i>)
		[<i>Expressions</i>]
		{ <i>FieldDefs</i> }
		(<i>Expression</i>)
<i>VarExpression</i>	\rightarrow	<u><i>Identifier</i></u>
		<i>VarExpression</i> [<i>Expression</i>]
		<i>VarExpression</i> . <u><i>Identifier</i></u>

<i>BinaryOperator</i>	\rightarrow	$\wedge \mid * \mid / \mid + \mid - \mid < \mid \leq \mid == \mid != \mid \geq \mid > \mid \&\& \mid \parallel$
<i>UnaryOperator</i>	\rightarrow	$- \mid !$
<i>FieldDefs</i>	\rightarrow	ϵ <i>FieldDefs</i> ₁
<i>FieldDefs</i> ₁	\rightarrow	<i>FieldDef</i> <i>FieldDef</i> , <i>FieldDefs</i> ₁
<i>FieldDef</i>	\rightarrow	<u>Identifier</u> = <i>Expression</i>
<i>Declarations</i>	\rightarrow	<i>Declaration</i> <i>Declaration</i> ; <i>Declarations</i>
<i>Declaration</i>	\rightarrow	const <u>Identifier</u> : <i>TypeDenoter</i> = <i>Expression</i> var <u>Identifier</u> : <i>TypeDenoter</i> var <u>Identifier</u> : <i>TypeDenoter</i> := <i>Expression</i> fun <u>Identifier</u> (<i>ArgDecls</i>) : <i>TypeDenoter</i> = <i>Expression</i> proc <u>Identifier</u> (<i>ArgDecls</i>) <i>Command</i>
<i>ArgDecls</i>	\rightarrow	ϵ <i>ArgDecls</i> ₁
<i>ArgDecls</i> ₁	\rightarrow	<i>ArgDecl</i> <i>ArgDecl</i> , <i>ArgDecls</i> ₁
<i>ArgDecl</i>	\rightarrow	<u>Identifier</u> : <i>TypeDenoter</i> in <u>Identifier</u> : <i>TypeDenoter</i> out <u>Identifier</u> : <i>TypeDenoter</i> var <u>Identifier</u> : <i>TypeDenoter</i>
<i>TypeDenoter</i>	\rightarrow	<u>Identifier</u> <i>TypeDenoter</i> [<u>IntegerLiteral</u>] { <i>FieldTypes</i> }

$$\begin{aligned}
FieldTypes &\rightarrow \epsilon \\
&| FieldTypes_1 \\
FieldTypes_1 &\rightarrow FieldType \\
&| FieldType , FieldTypes_1 \\
FieldTypes &\rightarrow \underline{Identifier} : TypeDenoter
\end{aligned}$$

Note that the productions for *Expression* makes the grammar as stated above ambiguous. Operator precedence and associativity for the *binary* operators as defined in the following table is used to disambiguate:

Operator	Precedence	Associativity
$\hat{}$	1	right
$*$ /	2	left
$+$ -	3	left
$<$ \leq $==$ $!=$ \geq $>$	4	non
$\&\&$	5	left
$\ $	6	left

A precedence level of 1 means the highest precedence, 2 means second highest, and so on.

A.3 MiniTriangle Abstract Syntax

This is the MiniTriangle abstract syntax. It captures the tree structure of MiniTriangle programs as concisely as possible. For example, note that there is only one non-terminal for expressions as opposed to three in the grammar for the concrete syntax. Such “extra” non-terminals are helpful for specifying the exact details of the concrete syntax, and sometimes to avoid ambiguity. But once a program has been successfully parsed, its structure has been determined, and such extra detail no longer serve any purpose. Another difference is that concrete unary and concrete (infix) binary operator application are subsumed by function application, as such operators *are* functions of one and two arguments, respectively. As a consequence, a single “variable” terminal *Name* replaces *Identifier* and *Operator*; i.e., $\underline{Name} = \underline{Identifier} \cup \underline{Operator}$.

The rightmost column gives the node labels for drawing abstract syntax trees. They are also used as the names of the data constructors of the datatypes for representing MiniTriangle programs in the compiler. Note that some elements of concrete syntax, such as keywords, do occur in the productions. They are there to make the connection between the concrete and abstract syntax clear, and to provide an *alternative* textual representation for

the abstract syntax (e.g. for use in typing rules). However, these fragments of concrete syntax are *omitted* when drawing abstract syntax trees, as they are implied by the node labels and thus are superfluous. Also note that some of the productions make use of the EBNF $*$ -notation for sequences. When drawing an abstract syntax tree, that means that the corresponding nodes will have a varying number of children.

<i>Program</i>	\rightarrow	<i>Command</i>	Program
<i>Command</i>	\rightarrow	$Expression := Expression$ $Expression (Expression^*)$ $\text{begin } Command^* \text{ end}$ $\text{if } Expression \text{ then } Command$ $\text{else } Command$ $\text{while } Expression \text{ do } Command$ $\text{let } Declaration^* \text{ in } Command$	CmdAssign CmdCall CmdSeq CmdIf CmdWhile CmdLet
<i>Expression</i>	\rightarrow	<u><i>IntegerLiteral</i></u> <u><i>Name</i></u> $Expression (Expression^*)$ $[Expression^*]$ $Expression [Expression]$ $\{ (Name = Expression)^* \}$ $Expression . Name$	ExpLitInt ExpVar ExpApp ExpArY ExpIx ExpRcd ExpPrj
<i>Declaration</i>	\rightarrow	$\text{const } Name : TypeDenoter$ $= Expression$ $\text{var } Name : TypeDenoter$ $(:= Expression \mid \epsilon)$ $\text{fun } Name (ArgDecl^*)$ $: TypeDenoter = Expression$ $\text{proc } Name (ArgDecl^*) Command$	DeclConst DeclVar DeclFun DeclProc
<i>ArgDecl</i>	\rightarrow	<i>ArgMode</i> <u><i>Name</i></u> : <i>TypeDenoter</i>	ArgDecl
<i>ArgMode</i>	\rightarrow	ϵ <u>in</u> <u>out</u> <u>var</u>	ByVal ByRefIn ByRefOut ByRefVar

$$\begin{array}{lll}
 TypeDenoter & \rightarrow & \underline{Name} \qquad \qquad \qquad \text{TDBaseType} \\
 & \rightarrow & TypeDenoter [\, IntegerLiteral \,] \qquad \text{TDArray} \\
 & \rightarrow & \{ \, (\underline{Name} : TypeDenoter)^* \, \} \qquad \text{TDRecord}
 \end{array}$$

B MiniTriangle Standard Environment

The MiniTriangle standard environment provides the following types, constants, and procedures:

Name	Type	Description
<i>Types</i>		
Boolean	type	Boolean type; elements: <code>false</code> , <code>true</code>
Integer	type	Integer type; 32 bits
<i>Constants</i>		
<code>false</code>	Boolean	The truth value <code>false</code>
<code>true</code>	Boolean	The truth value <code>true</code>
<code>minint</code>	Integer	The smallest representable integer
<code>maxint</code>	Integer	The largest representable integer
<i>Procedures</i>		
<code>getInt</code>	<code>(Snk Integer) → Void</code>	Read integer from the terminal
<code>putInt</code>	<code>Integer → Void</code>	Write integer to the terminal
<code>skip</code>	<code>() → Void</code>	Do nothing

Additionally, the standard environment defines a number of functions and procedures used internally by the compiler, such as implementations of all operators and a procedure for reporting array indices out of bounds.

C MiniTriangle Type System

This appendix explains and specifies the type system for the HMTC version of MiniTriangle (henceforth just MiniTriangle). The notation essentially follows B.C. Pierce *Types and Programming Languages*, and the presentation is also inspired by that book. Naming conventions:

Symbol	Meaning/syntactic category
Γ	Type environment (or typing context)
S, T	<i>Type</i>
c	<i>Command</i>
e	<i>Expression</i>
d	<i>Declaration</i>
a	<i>ArgDecl</i>
x	Variable <i>Name</i>
p, f	Procedure/function <i>Name</i>
n	<i>IntegerLiteral</i>

In other words, we have $c \in \text{Command}$, $e \in \text{Expression}$, $x, p, f \in \text{Name}$, and so on. Subscripted and primed variants of these are also used with the same interpretation. For example, e_1 , e_x , e' all stand for expressions; i.e., $e_1 \in \text{Expression}$, $e_x \in \text{Expression}$, $e' \in \text{Expression}$.

The syntactic categories referred to above (*Command*, *Expression*, *Name*, etc.) are as specified by the MiniTriangle abstract syntax (appendix A.3), except for *Type* that is defined in the following. The syntactic category *TypeDenoter* in the abstract syntax corresponds to types that may occur as parts of declarations: at present the base types Boolean and Integer along with arrays of arbitrary types and records. However, to specify the MiniTriangle type system we need a refined notion of type that can express the exact type of *any* typed MiniTriangle entity. *Type* is defined such that $\text{TypeDenoter} \subseteq \text{Type}$. This simplifies the typing rules for declarations, but is not strictly speaking necessary.

The typing rules make use of vector notation for conciseness. For example, \overline{T} is a sequence of zero or more types; i.e. $\overline{T} \in \text{Type}^*$. However, we will also use vector notation as a shorthand to avoid distracting repetition when the meaning is clear. For example, we take $\overline{e} : \overline{T}$ to be a shorthand for $\overline{e} = e_1, e_2, \dots, e_n$; $\overline{T} = T_1, T_2, \dots, T_n$; and $e_1 : T_1, e_2 : T_2, \dots, e_n : T_n$ for some $n \in \mathbb{N}$.

C.1 MiniTriangle Types

The syntax of types is defined by the following context-free grammar:

<i>Type</i>	\rightarrow	<i>Types:</i>
	Void	<i>The empty type (procedure return type)</i>
	Boolean	<i>The Boolean type</i>
	Integer	<i>The Integer type</i>
	Src Type	<i>Read-only variable reference (source)</i>
	Snk Type	<i>Write-only variable reference (sink)</i>
	Ref Type	<i>Variable reference</i>
	Type [IntegerLiteral]	<i>Array type</i>
	{(Name : Type)*}	<i>Record type</i>
	Type* → Type	<i>Type of procedures and functions (arrow)</i>

The intended meaning of these types should be clear, except perhaps for reference types that will be explained later.

The set of types specified above is more general than what is needed for the present version of MiniTriangle. For example, the type syntax does not rule out long chains of reference types (such as Ref (Ref (Ref Integer))) or higher-order procedures and functions, neither of which is possible since the MiniTriangle grammar simply does not provide any way to express programs making use of such types (and nor are they necessarily supported by the later stages of the compiler or abstract machine). However, the above grammar is much simpler than a more precise account of the types that actually can occur, and it facilitates future generalisations.

C.2 Imperative Variables and Dereferencing

One feature common to most imperative languages is that dereferencing is implicit when variables are read. For example, consider the following C-like declarations:

```
int x;
int y;
```

and the code fragment:

```
y = x + 1;
```

The variables `x` and `y` are each really a *reference* to a memory location where an integer can be stored, and the code above actually says:

Fetch the integer stored at the address `x` refers to. Add one to this integer. Store the result at the address `y` refers to.

However, note that whereas the addition (`+ 1`) and storing the result (`=`) are both operations that are explicitly mentioned in the code fragment, fetching is not: it is tacitly assumed that `x` stands for the *value* stored in the memory location associated with `x`, *not* the address of this location. In contrast, `y`, on the left-hand side of the assignment, *does* stand for the address of the memory location associated with `y`.

Thus we see that it is the usage context of a variable occurrence that determines if this occurrence is to be understood as the address of the variable or the value stored there. Not needing to explicitly indicate when a value has to be fetched from memory is what is meant by implicit dereferencing.

C.3 MiniTriangle Reference Types

The fact that variables are references to memory locations is made explicit in the MiniTriangle type system. For example, consider a declaration of a variable `x` of type `Integer`:

```
let
    var x : Integer
in
    ...
```

The type attributed to `x` in the body of the `let`-construct becomes `Ref Integer`; that is, a reference to a memory location that can hold a value of type `Integer`, or “reference of type `Integer`” for short. References are always typed to make the type of the referenced value clear.

There are actually three kinds of references in MiniTriangle:

- `Ref T`: read/write reference; i.e., values of type `T` can be written to and read from the referenced memory location.
- `Src T`: *source*; i.e., read-only reference.
- `Snk T`: *sink*; i.e., write-only reference.

These three types are collectively referred to as “reference types” or just “references”. We formalise this through the following predicate:

$$\begin{aligned} \text{reftype}(\text{Src } T) &\quad (1) \\ \text{reftype}(\text{Snk } T) &\quad (2) \\ \text{reftype}(\text{Ref } T) &\quad (3) \end{aligned}$$

Perhaps somewhat confusingly, we will also use “reference of type T ” in the narrower sense of `Ref T`. Usually the context will make it clear what is meant; otherwise we will write “read/write reference” to distinguish from “source” and “sink”.

Read/write references are used for variables, as illustrated by the example above, while sources are used for constants. For example, a definition:

```
const c : Integer = ...
```

results in the type `Src Integer` being attributed to `c`.

The three reference types are also used for passing arguments to procedures and functions by reference, allowing the specification of input/output, input, and output arguments. For example, the type of the procedure `getInt` in the MiniTriangle standard library, which is used to input an integer from the terminal, is `Snk Integer → Void`, meaning that it needs to be passed a reference to a memory location to which the integer read from the terminal can be written.

However, as in most imperative languages, dereferencing is still implicit when values of reference type are used. The way this works is that the MiniTriangle system is set up so that an entity of reference type *also* is considered to have the type of the referenced entity. To illustrate, when the integer variable `x` from the example above is used in an expression like `x + 1`, `x` has *both* type `Ref Integer` and `Integer`. As the type of the operator `+` is `(Integer, Integer) → Integer`, we can see that the first argument has to have type `Integer`. But as that is *one* of the possible types for `x`, the expression `x + 1` is well-typed in this case, with the overall type of the expression being `Integer`.

Behind the scenes, the type checker, in addition to checking that a program is well-formed in the sense defined by the MiniTriangle type system, inserts dereferencing operations to ensure that an expression of reference type actually gets the type it needs to have to fit with the usage context. To continue the example, the type checker will transform `x + 1` into something like `deref(x) + 1` in this case¹.

Should there be multiple levels of referencing, the type checker would insert two or more dereferencing operations as needed. For example, assume the type of `x` had been `Src (Ref Integer)` instead. Then `x + 1` would be transformed to `deref(deref(x)) + 1`. Ultimately, each `deref`-operation is translated into an instruction that reads the memory contents at the referenced location (which, in the case of nested referencing, will be another reference that may be further dereferenced in turn).

¹Note that `deref` is a “hidden” language construct, only used internally by the compiler at present, not a function that can be used by the programmer.

C.4 Subtyping

One type S is said to be a subtype of another type T , written $S <: T$, if a value of type S can be used wherever a value of type T is expected. Object-oriented languages, like Java and C#, are important examples of languages with type systems based on subtyping. Recall that an object that is an instance of a class C can be used wherever instances of any of C 's superclasses are expected. I.e. C is a subtype of C 's superclasses.

There are other possibilities as well. For example, in a language that has a type `Nat` for natural numbers, `Nat` might be considered to be a subtype of an integer type `Int`, as any natural number is also an integer. In this case, it is clear that there is a close connection between subsets and subtypes. That is, at least at a conceptual level: internally, it is not necessarily the case that elements of a subtype (in this case `Nat`) have the same representation as the corresponding elements of the supertype (in this case the elements of `Int` that are natural numbers).

However, if it is possible to make representations coincide in an implementation, this is advantageous as it makes the implementation of subtyping both easier and more efficient as there is no need to convert between representations at runtime. Object-oriented languages are thus typically designed to make this possible. For example, in Java, there is no runtime overhead associated with viewing an instance of some class C as an instance of a superclass of C ².

C.5 Subtyping in MiniTriangle

In MiniTriangle, the reference types naturally induces a subtyping relation as a read/write reference can be used in place of either a source (read-only) or a sink (write-only). The following inference rules define the MiniTriangle subtyping relation:

²*Downcasting* is a different matter. As there is no static guarantee that an entity that statically is known to have some type T at runtime actually has some more refined type S as dictated by the cast, with $S <: T$, a runtime check that ensures that the entity actually has the more refined type needs to be inserted to guarantee language safety.

$$\begin{array}{ll}
T <: T & (1) \\
\frac{S <: T}{\mathbf{Src} \ S <: \mathbf{Src} \ T} & (2) \\
\frac{T <: S}{\mathbf{Snk} \ S <: \mathbf{Snk} \ T} & (3) \\
\frac{S <: T \quad T <: S}{\mathbf{Ref} \ S <: \mathbf{Ref} \ T} & (4) \\
\frac{S <: T}{\mathbf{Ref} \ S <: \mathbf{Src} \ T} & (5) \\
\frac{T <: S}{\mathbf{Ref} \ S <: \mathbf{Snk} \ T} & (6) \\
\frac{\overline{T} <: \overline{S} \quad S' <: T'}{\overline{S} \rightarrow S' <: \overline{T} \rightarrow T'} & (7)
\end{array}$$

Rule (1) says that any type is a subtype of itself; i.e., subtyping is *reflexive*. In fact, subtyping is also *transitive*, which is essential, making it a *preorder*. However, transitivity is not manifest from these rules but has to be proved (by induction).

If we have a value of type $\mathbf{Src} \ S$, then it can be dereferenced to obtain a value of type S . Clearly, if $S <: T$, then this value of type S can be used wherever a value of type T is expected. But then it follows that we can use a value of type $\mathbf{Src} \ S$ wherever a value of type $\mathbf{Src} \ T$ is expected; i.e., it should be the case that $\mathbf{Src} \ S <: \mathbf{Src} \ T$. Rule (2) formalises this. Note that \mathbf{Src} is *covariant*: it preserves the subtyping ordering.

By a similar line of reasoning, we expect a subtyping relationship to hold between $\mathbf{Snk} \ S$ and $\mathbf{Snk} \ T$ if S and T are related by subtyping. Rule (3) captures this. Note that this time, the requirement is that $T <: S$; i.e., \mathbf{Snk} is a *contravariant* type constructor (reverses the subtyping ordering). Why? Well, if we are using a sink of type S at type $\mathbf{Snk} \ T$, then any value of type T written into this sink could potentially be used at type S through some different route. But that will only work if we insist that T is a subtype of S .

Rule (4) is essentially a combination of rule (2) and (3) as a reference simultaneously is a source and a sink.

Rules (5) and (6) formalises that a reference is both a source and sink and thus can be used in place of either. Note that rule (5) is covariant while rule (6) is contravariant for the same reasons as above.

Rule (7), finally, formalises when a procedure or function type is a subtype of another. Note the covariant subtyping relationship between the result types and the contravariant one between the arguments. Intuitively, this is because returning a value is akin to a source, while passing arguments is akin to writing them to a sink.

All reference types in MiniTriangle ($\mathbf{Src} \ T$, $\mathbf{Snk} \ T$, $\mathbf{Ref} \ T$) share a common representation: they are just pointers. The differentiation between sources, sinks, and read/write references only serve to keep track of read and write

“permissions”. No change of representation is needed to use a value of some type T at one of T ’s supertypes, and thus no runtime overhead is incurred.

C.6 Implicit Dereferencing in MiniTriangle

As discussed in Appendix C.3, dereferencing is implicit in MiniTriangle. At the type level, this is manifested through an implicit coercion from a reference type to the type of the referenced entity, meaning that the type system ascribes more than one type to such entities. For example, an expression of type **Ref Integer** is implicitly coerced to **Integer** when necessary, meaning it can be typed at both types.

However, unlike subtyping, this is not just a matter of viewing an entity as having a more refined type, but these coercions actually involve a representational change: from a reference (pointer) to an entity to the referenced entity itself by following the reference. The MiniTriangle type checker therefore has to insert *explicit* dereferencing operations wherever the type system makes use of an implicit coercion from a reference type to the referenced type.

The predicates sources and sinks defined below account for this. If it is the case that $\text{sources}(S, T)$ holds, this means that a value of type S can “source” a value of type T through zero or more dereferencing operations. For example, from a value of type **Ref (Src (Snk Boolean))** we can obtain a value of type **Snk Boolean** by dereferencing twice. If $\text{sinks}(S, T)$ holds, it means that it is possible to obtain a sink to which a value of type T can be written from a value of type S through zero or more dereferencing operations. For example, from a value of type **Ref (Snk Integer)** it is possible to obtain an integer sink by dereferencing once.

$$\begin{array}{c}
 \frac{S <: T}{\text{sources}(S, T)} \quad (1) \qquad \frac{T <: S}{\text{sinks}(\text{Snk } S, T)} \quad (1) \\
 \frac{\text{sources}(S, T)}{\text{sources}(\text{Src } S, T)} \quad (2) \qquad \frac{T <: S}{\text{sinks}(\text{Ref } S, T)} \quad (2) \\
 \frac{\text{sources}(S, T)}{\text{sources}(\text{Ref } S, T)} \quad (3) \qquad \frac{\text{sinks}(S, T)}{\text{sinks}(\text{Src } S, T)} \quad (3) \\
 \frac{\text{sinks}(S, T)}{\text{sinks}(\text{Ref } S, T)} \quad (4)
 \end{array}$$

Rules (2) and (3) of the definition of the predicate sources and rules (3) and (4) of the definition of the predicate sinks are justified by dereferencing. Thus, when the type checker carries out a typing derivation, it will have to

insert a dereferencing operation exactly once each time one of these rules are used³. By contrast, rule (1) of the definition of the predicate sources and rules (1) and (2) of the definition of the predicate sinks are justified by subtyping alone and no dereferencing operations are thus inserted in those cases. For example, if $S <: T$, then it is clearly possible to obtain a value of type T from a value of type S as the latter also has type T by virtue of subtyping.

C.7 MiniTriangle Typing Relations

We are now in a position to specify the MiniTriangle type system as such. This is done through the following typing relations:

$\Gamma \vdash c$	Command c is well-formed in type environment Γ
$\Gamma \vdash e : T$	Expression e has type T in type environment Γ
$\Gamma ; \Gamma_B \vdash \bar{d} \Gamma'$	Declarations \bar{d} are well-formed in type environments Γ and Γ_B , extending the environment Γ to Γ'
$\Gamma \vdash \bar{a} \Gamma'$	Argument declarations \bar{a} are well-formed in type environment Γ , extending the environment to Γ'

A type environment associates names with types, allowing the type of a named entity to be found (if it is in scope). However, it is also necessary to keep track of the current scope level and the scope level at which a named entity has been declared. An environment is therefore taken to be a pair of scope level and a list of pairs of the form

$$x_{(n)} : T$$

This can be read “ x at scope level n has type T ”. Additionally, for an environment to be well-formed, the scope levels of the names must not exceed the scope level of the environment. Thus, an environment Γ has the form:

$$\Gamma = (n; \langle x_{1(n_1)} : T_1, x_{2(n_2)} : T_2, \dots x_{k(n_k)} : T_k \rangle)$$

where $n_i \leq n$ for $1 \leq i \leq k$.

To keep the typing rules clear and concise, we adopt some notation, conventions, and abbreviations in relation to environments.

- $\Gamma_{(n)}$: A subscript within brackets is used to refer to the scope level of an environment. If two environments with the same name but different scope levels occur in a rule, e.g. $\Gamma_{(n)}$ and $\Gamma_{(n+1)}$, it is understood that the environments are the same, except for the scope level.

³This correspondence is particularly clear in the definition of the function **sources** in the type checker. The definition is not much more than a transliteration of the rules defining the predicate sources, with the addition of code for inserting the dereferencing operations.

- $\Gamma(x)$: The type of x in Γ ; if there are more than one entity named x in scope, then the type of one with the highest scope level.
- $x \in \Gamma$: True if $x_{(n)} : T$ for some n and T is in Γ .
- $x_{(n)} \in \Gamma$: True if $x_{(n)} : T$ for some T is in Γ .
- $\Gamma, x : T$: Given

$$\Gamma = (n; \langle x_{1(n_1)} : T_1, x_{2(n_2)} : T_2, \dots x_{k(n_k)} : T_k \rangle)$$

$\Gamma, x : T$ is shorthand for the extended environment

$$(n; \langle x_{1(n_1)} : T_1, x_{2(n_2)} : T_2, \dots x_{k(n_k)} : T_k, x_{(n)} : T \rangle)$$

and an additional side condition:

$$x_{(n)} \notin \Gamma$$

In other words, the environment Γ is extended with the information that an entity named x at the present scope level n has type T , but only if there is not already a declaration for x at this scope level.

C.7.1 Commands

The following inference rules define the typing relation specifying well-formed commands:

$$\frac{\Gamma \vdash e_x : S \quad \Gamma \vdash e : T \quad \neg \text{reftype}(T) \quad \text{sinks}(S, T)}{\Gamma \vdash e_x := e} \quad (\text{T-ASSIGN})$$

$$\frac{\Gamma \vdash e_p : \overline{T} \rightarrow \text{Void} \quad \Gamma \vdash \overline{e} : \overline{T}}{\Gamma \vdash e_p(\overline{e})} \quad (\text{T-CALL})$$

$$\frac{\Gamma \vdash \overline{c}}{\Gamma \vdash \text{begin } \overline{c} \text{ end}} \quad (\text{T-SEQ})$$

$$\frac{\Gamma \vdash e : \text{Boolean} \quad \Gamma \vdash c_1 \quad \Gamma \vdash c_2}{\Gamma \vdash \text{if } e \text{ then } c_1 \text{ else } c_2} \quad (\text{T-IF})$$

$$\frac{\Gamma \vdash e : \text{Boolean} \quad \Gamma \vdash c}{\Gamma \vdash \text{while } e \text{ do } c} \quad (\text{T-WHILE})$$

$$\frac{\Gamma_{(n+1)} ; \Gamma' \vdash \overline{d} \mid \Gamma' \quad \Gamma' \vdash c}{\Gamma_{(n)} \vdash \text{let } \overline{d} \text{ in } c} \quad (\text{T-LET})$$

Most of the rules above are straightforward and fairly standard. However, T-ASSIGN and T-LET deserve some comments. The rule T-ASSIGN says that the expression e_x denoting the variable that will be written to must have a type S such that S can sink (possibly after implicit dereferencing) a value of type T , where T is the type of the expression yielding the value to be written. This should be fairly intuitive. The reason that the rule further insists that the type T must not be a reference type is to avoid ambiguity. Because an expression can have more than one type due to implicit dereferencing (and subtyping; see Appendix C.6), it could be that it would not be clear just how many dereferencing steps should be carried out otherwise.

For example, suppose $e_x : \text{Ref } (\text{Snk Integer})$ and $e : \text{Ref Integer}$. Note that $\text{sinks}(\text{Ref } (\text{Snk Integer}), \text{Ref Integer})$ holds because $\text{Ref } (\text{Integer}) <: \text{Snk } (\text{Integer})$ (rule (2) in the definition of sinks). This typing corresponds to storing a reference (pointer) to an integer into the location referred to by the value of e_x .

However, there is another possibility. Note that e_x can source a integer sink (one dereferencing operation). Thus we also have $e_x : \text{Snk Integer}$. Similarly, e can source an integer (one dereferencing operation). Thus we also have $e : \text{Integer}$. And clearly, $\text{sinks}(\text{Snk Integer}, \text{Integer})$ holds. This typing corresponds to fetching the integer referred to be e and storing in the location *indirectly* referred to be e_x . This is a very different semantics from above. By insisting that the assigned values are as “dereferenced as possible”, we avoid this ambiguity (at the expense of losing flexibility not needed in the present version of MiniTriangle anyway), opting for the latter interpretation.

Regarding T-LET, note how the scope level is increased by one before extending the environment according to the declarations. Also note that Γ' , the extended environment, is not only used for checking that the body of the `let-command` is well-formed, but also when checking that the declarations themselves are well-formed. The latter allows for recursive procedures and functions. See the rules for declarations below.

C.7.2 Expressions

The typing rules for expressions are as follows:

$$\frac{\Gamma \vdash e : S \quad \text{sources}(S, T)}{\Gamma \vdash e : T} \quad (\text{T-SOURCES})$$

$$\Gamma \vdash n : \text{Integer} \quad (\text{T-LITINT})$$

$$\frac{\Gamma(x) = T}{\Gamma \vdash x : T} \quad (\text{T-VAR})$$

$$\frac{\Gamma \vdash e_f : \overline{T} \rightarrow T' \quad \Gamma \vdash \overline{e} : \overline{T}}{\Gamma \vdash e_f(\overline{e}) : T'} \quad (\text{T-APP})$$

$$\frac{\Gamma \vdash \overline{e} : T \quad n = \text{length}(\overline{e})}{\Gamma \vdash [\overline{e}] : T[n]} \quad (\text{T-ARY})$$

$$\frac{\Gamma \vdash e_a : R(T[n]) \quad \Gamma \vdash e_i : \text{Integer} \quad R \in \{\text{Src}, \text{Snk}, \text{Ref}\}}{\Gamma \vdash e_a[e_i] : R T} \quad (\text{T-IX})$$

$$\frac{\Gamma \vdash \overline{e} : \overline{T} \quad \text{alldistinct}(\overline{x})}{\Gamma \vdash \{\overline{x} = e\} : \{\overline{x} : \overline{T}\}} \quad (\text{T-RCD})$$

$$\frac{\Gamma \vdash e : R(\{\overline{x : T}\}) \quad x_n : T_n \in \overline{x : T} \quad R \in \{\text{Src}, \text{Snk}, \text{Ref}\}}{\Gamma \vdash e.x_n : R T_n} \quad (\text{T-PRJ})$$

The rules T-LITINT, T-VAR, T-APP are standard. For T-ARY, we take the common type T of the array elements to be arbitrary if $n = 0$. T-IX specifies that the array expression must denote a *reference* to an array, and that the result of indexing this array is a reference (of the same kind) to an individual element of that array. (This ensures arrays are accessed in place, rather than needlessly being copied to the stack first.) The constraint $\text{alldistinct}(\overline{x})$ in the typing rule T-RCD for record expressions means that each name in the vector \overline{x} of field names must be distinct from the others. The typing rule T-PRJ for record projection is similar to the rule for array indexing (T-IX) in that the expression must denote a reference to a record and that the result is a reference (of the same kind) to the specified field of the record. (The rationale for the design is the same as that for arrays.) The notation $x_n : T_n \in \{\overline{x : T}\}$ means that the field name and type pair $x_n : T_n$ must be one of field name and type pairs of the type of the record. Finally, T-SOURCES is the rule that allows implicit dereferencing. It says that any expression that has type S in the environment Γ also has type T in the same environment if S can source T ; i.e., if a value of type S can be dereferenced zero or more times

to obtain a value of type T . See the definition of sources in Section C.6, and note how the sources predicate also takes care of subtyping thanks to rule (1) of its definition.

C.7.3 Declarations

Finally we turn to the typing rules for declarations. The relation to determine if a list of declarations is well-typed is defined inductively by simply recursing down this list, extending the environment in which the declarations are checked along the way, thus ensuring that each declared entity is brought into scope in the following declarations. Once the end of the list is reached, the final environment is returned as the one in which to check the body of a `let`-command; see rule T-LET above. However, to allow for (mutually) recursive functions and procedures, the bodies of these are checked in a separate “body” environment, Γ_B . The rule T-LET is set up so that Γ_B is the same as the environment in which the body of the `let`-command is checked, meaning that *all* declared entities are in scope in the bodies of all functions and procedures.

$$\begin{array}{c}
 \frac{\Gamma \vdash e : T \quad \Gamma, x : \text{Src } T ; \Gamma_B \vdash \bar{d} | \Gamma'}{\Gamma \vdash \text{wellinit}(n, e)} \\
 \hline
 \Gamma_{(n)} ; \Gamma_B \vdash \text{const } x : T = e ; \bar{d} | \Gamma' \quad (\text{T-DECLCONST})
 \end{array}$$

$$\frac{\Gamma, x : \text{Ref } T ; \Gamma_B \vdash \bar{d} | \Gamma'}{\Gamma ; \Gamma_B \vdash \text{var } x : T ; \bar{d} | \Gamma'} \quad (\text{T-DECLVAR})$$

$$\begin{array}{c}
 \frac{\Gamma \vdash e : T \quad \Gamma, x : \text{Ref } T ; \Gamma_B \vdash \bar{d} | \Gamma'}{\Gamma \vdash \text{wellinit}(n, e)} \\
 \hline
 \Gamma_{(n)} ; \Gamma_B \vdash \text{var } x : T := e ; \bar{d} | \Gamma' \quad (\text{T-DECLINITVAR})
 \end{array}$$

$$\frac{\Gamma_{B(n+1)} \vdash \bar{a} | \Gamma'_B \quad \Gamma'_B \vdash c}{\Gamma ; \Gamma_{B(n)} \vdash \text{proc } p(\bar{a}) c ; \bar{d} | \Gamma'} \quad (\text{T-DECLPROC})$$

$$\frac{\Gamma_{B(n+1)} \vdash \bar{a} | \Gamma'_B \quad \Gamma'_B \vdash e : T \quad \Gamma, f : \text{funtype}(\bar{a}, T) ; \Gamma_B \vdash \bar{d} | \Gamma'}{\Gamma ; \Gamma_{B(n)} \vdash \text{fun } f(\bar{a}) : T = e ; \bar{d} | \Gamma'} \quad (\text{T-DECLFUN})$$

$$\Gamma ; \Gamma_B \vdash \epsilon | \Gamma \quad (\text{T-DECLEMPTY})$$

There are five main cases: constant definition, uninitialised variable declaration, initialised variable declaration, procedure declaration, and function

declaration. Note how the declared type of the constant or variable is checked against the defining or initialisation expression in case there is one. Further, note how the type of the declared entity is added to the environment as $\text{Src } T$ in the case of a constant and $\text{Ref } T$ in the case of a variable. Thus in all cases, the type of the declared entity is a reference type, making it clear that a dereferencing operation is needed to get the actual value. This is because both constants and variables are stored on the stack, with only their addresses (offsets with respect to the stack base or local base), not their values, known at compile time. The type of a constant is $\text{Src } T$ as constants can only be read (after their initial definition), while the type of variables is $\text{Ref } T$ as they can be both read and written. Finally, note how Γ_B , rather than Γ , as discussed above, is extended by the declarations of formal arguments to Γ'_B , and then used for checking the well-formedness of the bodies of declared procedures and functions.

The types of declared procedures and functions are computed by the following auxiliary functions:

$$\begin{aligned}
 \text{proctype}(\bar{a}) &= \text{argtypes}(\bar{a}) \rightarrow \text{Void} \\
 \text{funtype}(\bar{a}, T) &= \text{argtypes}(\bar{a}) \rightarrow T \\
 \\
 \text{argtypes}(x : T, \bar{a}) &= T, \text{argtypes}(\bar{a}) \\
 \text{argtypes}(\text{in } x : T, \bar{a}) &= \text{Src } T, \text{argtypes}(\bar{a}) \\
 \text{argtypes}(\text{out } x : T, \bar{a}) &= \text{Snk } T, \text{argtypes}(\bar{a}) \\
 \text{argtypes}(\text{var } x : T, \bar{a}) &= \text{Ref } T, \text{argtypes}(\bar{a}) \\
 \text{argtypes}(\epsilon) &= \epsilon
 \end{aligned}$$

The defining expressions for constants and the initialisation expressions for variables must be “well-initialised”, meaning that they must not use functions defined in the same `let`-block. This make a straightforward implementation of constant and variable allocation and initialisation possible. (The type system as such would work fine without it.) The relation $\Gamma \vdash \text{wellinit}(n, e)$ formalises this requirement:

$$\begin{array}{ll}
\Gamma \vdash \text{wellinit}(_, n) & (\text{WI-LITINT}) \\
\\
\Gamma \vdash \text{wellinit}(_, x) & (\text{WI-VAR}) \\
\\
\frac{f(n) \notin \Gamma \quad \Gamma \vdash \text{wellinit}(n, \bar{e})}{\Gamma \vdash \text{wellinit}(n, f(\bar{e}))} & (\text{WI-APP}) \\
\\
\frac{\Gamma \vdash \text{wellinit}(n, \bar{e})}{\Gamma \vdash \text{wellinit}(n, [\bar{e}])} & (\text{WI-ARY}) \\
\\
\frac{\Gamma \vdash \text{wellinit}(n, e_a) \quad \Gamma \vdash \text{wellinit}(n, e_i)}{\Gamma \vdash \text{wellinit}(n, e_a[e_i])} & (\text{WI-IX}) \\
\\
\frac{\Gamma \vdash \text{wellinit}(n, \bar{e})}{\Gamma \vdash \text{wellinit}(n, \{\bar{x}=\bar{e}\})} & (\text{WI-RCD}) \\
\\
\frac{\Gamma \vdash \text{wellinit}(n, e)}{\Gamma \vdash \text{wellinit}(n, e.x_n)} & (\text{WI-PRJ})
\end{array}$$

Note: if an unknown function is being called (although not currently possible), then that would conservatively not be considered well-initialised.

Finally, the relation for checking argument declarations and extending the given environment has a similar structure to that for checking declarations. Note how all formal arguments are treated as constants (sources).

$$\begin{array}{ll}
\frac{\Gamma, x : \text{Src } T \vdash \bar{d} \mid \Gamma'}{\Gamma \vdash x : T ; \bar{d} \mid \Gamma'} & (\text{T-DECLARG}) \\
\\
\frac{\Gamma, x : \text{Src } (\text{Src } T) \vdash \bar{d} \mid \Gamma'}{\Gamma \vdash \text{in } x : T ; \bar{d} \mid \Gamma'} & (\text{T-DECLINARG}) \\
\\
\frac{\Gamma, x : \text{Src } (\text{Snk } T) \vdash \bar{d} \mid \Gamma'}{\Gamma \vdash \text{out } x : T ; \bar{d} \mid \Gamma'} & (\text{T-DECLOUTARG}) \\
\\
\frac{\Gamma, x : \text{Src } (\text{Ref } T) \vdash \bar{d} \mid \Gamma'}{\Gamma \vdash \text{var } x : T ; \bar{d} \mid \Gamma'} & (\text{T-DECLVARARG}) \\
\\
\Gamma \vdash \epsilon \mid \Gamma & (\text{T-DECLARGEMPTY})
\end{array}$$

C.8 Notes on MiniTriangle Type Checker Implementation

As far as languages and type systems go, MiniTriangle is a small language and its type system is fairly simple. Yet, turning the specification of its type system into a type checker is not entirely straightforward. The main issue is that the typing rules are not algorithmic; i.e., they can't directly be read as specifying a function, where the inputs uniquely determine the output, but just a relation that can be one to many. As a case in point, consider the rule T-SOURCES that allows a single expression to have many possible types. This was deliberate: the typing rules have been written for clarity and ease of understanding, keeping details about exactly *how* to go about checking types from obscuring the specification of the type system.

There are a number of approaches for handling this. A principled one would be to allow for many possible types, which then gradually would be narrowed down to, hopefully, only a single type once all constraints have been taken into account.

The MiniTriangle type system has, however, been designed to permit a simpler approach. By carefully considering what is known and unknown in each rule, each rule can be turned into one, or, in case what is known and not differ in different contexts, more functions, with the knowns as inputs and the unknowns as outputs, uniquely determined by the inputs. As a case in point, the relation $\Gamma \vdash e : T$ is implemented as a number of functions, including `chkTpExp` to be used to check that an expression e has a known type T in a known environment Γ (environments are always known), and `infTpExp` that infers the “least dereferenced” type for e , if well-typed, in the known environment Γ .

The notion of scope level is refined into major and minor scope levels in the implementation. The major scope level is incremented for bodies of procedures and functions. This is because a new activation record is created when a function or procedure is called. The difference in major scope level between the context in which a variable reference occurs and the scope level of the referenced variable determines how many static links that have to be followed to get to the correct activation record (see the lecture notes). In contrast, variables declared by a `let`-command gets allocated in the current activation record. Therefore only the minor scope levels needs incrementing.

D Triangle Abstract Machine Instructions

Meta variable	Meaning
a	Address: one of the forms specified by table below when part of an instruction, specific stack address when on the stack
b	Boolean value (false = 0 or true = 1)
ca	Code address; address to routine in the code segment
d	Displacement; i.e., offset w.r.t. address in register or on the stack
l	Label name
m, n	Integer
x, y	Any kind of stack data
x^n	Vector of n items, in this case any kind

Address form	Description
$[SB + d]$	Address given by contents of register SB (Stack Base) $+/-$ displacement d
$[LB + d]$	Address given by contents of register LB (Local Base) $+/-$ displacement d
$[ST + d]$	Address given by contents of register ST (Stack Top) $+/-$ displacement d
$[LB - d]$	
$[ST - d]$	

Instruction	Stack effect	Description
<i>Label</i>		
Label l	—	Pseudo instruction: symbolic location
<i>Load and store</i>		
LOADL n	$\dots \Rightarrow n, \dots$	Push literal integer n onto stack
LOADCA l	$\dots \Rightarrow \text{addr}(l), \dots$	Push address of label l (code segment) onto stack
LOAD a	$\dots \Rightarrow [a], \dots$	Push contents at address a onto stack
LOADA a	$\dots \Rightarrow a, \dots$	Push address a onto stack
LOADI d	$a, \dots \Rightarrow [a + d], \dots$	Load indirectly; push contents at address $a + d$ onto stack
STORE a	$n, \dots \Rightarrow \dots$	Pop value n from stack and store at address a
STOREI d	$a, n, \dots \Rightarrow \dots$	Store indirectly; store n at address $a + d$

Instruction	Stack effect	Description
<i>Block operations</i>		
LOADLB $m\ n$	$\dots \Rightarrow m^n, \dots$	Push block of n literal integers m onto stack
LOADIB n	$a, \dots \Rightarrow [a + (n - 1)], \dots, [a + 0], \dots$	Load block of size n indirectly
STOREIB n	$a, x^n, \dots \Rightarrow \dots$	Store block of size n indirectly
POP $m\ n$	$x^m, y^n, \dots \Rightarrow x^m, \dots$	Pop n values below top m values
<i>Arithmetic operations</i>		
ADD	$n_2, n_1, \dots \Rightarrow n_1 + n_2, \dots$	Add n_1 and n_2 , replacing n_1 and n_2 with the sum
SUB	$n_2, n_1, \dots \Rightarrow n_1 - n_2, \dots$	Subtract n_2 from n_1 , replacing n_1 and n_2 with the difference
MUL	$n_2, n_1, \dots \Rightarrow n_1 \cdot n_2, \dots$	Multiply n_1 by n_2 , replacing n_1 and n_2 with the product
DIV	$n_2, n_1, \dots \Rightarrow n_1/n_2, \dots$	Divide n_1 by n_2 , replacing n_1 and n_2 with the (integer) quotient
NEG	$n, \dots \Rightarrow -n, \dots$	Negate n , replacing n with the result
<i>Comparison & logical operations</i> (false = 0, true = 1)		
LSS	$n_2, n_1, \dots \Rightarrow n_1 < n_2, \dots$	Check if n_1 is smaller than n_2 , replacing n_1 and n_2 with the Boolean result
EQL	$n_2, n_1, \dots \Rightarrow n_1 = n_2, \dots$	Check if n_1 is equal to n_2 , replacing n_1 and n_2 with the Boolean result
GTR	$n_2, n_1, \dots \Rightarrow n_1 > n_2, \dots$	Check if n_1 is greater than n_2 , replacing n_1 and n_2 with the Boolean result
AND	$b_2, b_1, \dots \Rightarrow b_1 \wedge b_2, \dots$	Logical conjunction of b_1 and b_2 , replacing b_1 and b_2 with the Boolean result
OR	$b_2, b_1, \dots \Rightarrow b_1 \vee b_2, \dots$	Logical disjunction of b_1 and b_2 , replacing b_1 and b_2 with the Boolean result
NOT	$b, \dots \Rightarrow \neg b, \dots$	Logical negation of b , replacing b with the result

Instruction	Stack effect	Description
<i>Control transfer</i>		
JUMP l	—	Jump unconditionally to location identified by label l
JUMPIFZ l	$n, \dots \Rightarrow \dots$	Jump to location identified by label l if $n = 0$ (i.e., n is false)
JUMPIFNZ l	$n, \dots \Rightarrow \dots$	Jump to location identified by label l if $n \neq 0$ (i.e., n is true)
CALL l	$\dots \Rightarrow \text{PC} + 1, \text{LB}, 0, \dots$	Call global subroutine at location l : Activation record set up by pushing static link (0 for global level), dynamic link (value of LB), and return address ($\text{PC}+1$, address of instruction after the call instruction) onto the stack; $\text{PC} = l$; LB = start of activation record (address of static link)
CALLI	$ca, sl, \dots \Rightarrow \text{PC} + 1, \text{LB}, sl, \dots$	Call subroutine indirectly: address of routine (ca) and static link to use (sl) on top of the stack; activation record and new PC and LB as for CALL
RETURN $m n$	$x^m, ra, olb, sl, y^n \dots \Rightarrow x^m, \dots$	Return from subroutine, replacing activation record by result, jumping to return address ($\text{PC} = ra$), and restoring the old local base ($\text{LB} = olb$)
<i>Input/Output</i>		
PUTINT	$n, \dots \Rightarrow \dots$	Print n to the terminal as a decimal integer
PUTCHR	$n, \dots \Rightarrow \dots$	Print the character with character code n to the terminal
GETINT	$\dots \Rightarrow n, \dots$	Read decimal integer n from the terminal and push onto the stack
GETCHR	$\dots \Rightarrow n, \dots$	Read character from the terminal and push its character code n onto the stack
<i>TAM Control</i>		
HALT	—	Stop execution and halt the machine