

A Low Level Language with Precise Integer Types

Ian Shehadeh

Department of Math and Computer Science
St. Mary's College of Maryland
St. Mary's City, Maryland, USA

IRHSHEADEH@SMCM.EDU

Abstract

We present *Howlite* a language targeting RISC-V, with a similar level of abstraction to C. Howlite uses a single scalar type, *integer*, which allows users to specify exactly the set of values allowed. Collection types are checked with a simple, structural bi-directional type checker.

Keywords: programming language

1 Introduction

Memory safety in systems programming languages has garnered a lot of attention in the last several years. A compiler that enforces strict rules on object's lifetime and mutability is helpful in large projects, especially when security is a top concern. Checking these properties at compile time allows the compiler to omit parts of its runtime, like a garbage collector, while providing similar guarantees.

These innovations in language design fail to directly address a class of problems where direct memory manipulation is essential. These problems force the programmer to fully disable the compiler's checks, or encourage awkward solutions which trade clarity for small guarantees.

Howlite aims to address these problems. Howlite is not a language to write a web server, it is not for writing applications, it isn't even a language for writing programming languages. It is a language for writing a single module for a very specific data structure, wrapped in a python library. It is a language for writing a boot loaded, or the entrypoint to a kernel. The compiler does not impose strict requirements on how the programmer manages memory, or accesses data. Instead, the type systems gives a rich set of tools, allowing one to set their own constraints.

2 Syntax

```
func boundedAdd(a: u32, b: u32): u32 {  
  if U32_MAX - a > b {  
    U32_MAX  
  } else {  
    a + b  
  }  
}
```

Listing 1: Addition without Overflow

Howlite's syntax prioritizes familiarity, ease of parsing, and clarity. The syntax should be familiar, someone unfamiliar with the language should be able to immediately grasp the programmer's intent, even if they do not understand every line. In a similar vein, the programmer should be guided towards writing code that is easily legible by others. We approach this issue by providing language

constructs that clearly express intent. For example, flow control constructs, like if statements may have a value. This allows the programmer to clearly show a variable's value is

the result of some condition. In order to make tooling easier to write, we prioritize creating an unambiguous grammar, with no constructs that require unbounded look-ahead.

2.1 Familiarity

Howlite code should be recognizable to C programmers. For this reason, we use curly braces (“{” and “}”) to denote blocks of code. We use familiar imperative keywords: “if”, “else”, and “while”, and mathematical expressions follow typical infix notation. Howlite differs from C in that it requires a sigil character or keyword before beginning a new construct. Types do not lead in variable assignments or functions. Instead we use the “let” or “func” keywords, respectively. This simplifies parsing, since we know what type of statement or expression will follow, similarly, type ascriptions are always prefixed with `:`. These keywords and symbols were decided by surveying popular languages during design. For example, “let”, and `:` come from TypeScript, while “func” is a keyword in Go.

2.2 Clarity

TODO

3 Type Checking

Howlite’s implements a simple bi-directional type checker [Dunfield and Krishnaswami (2020)]. Every node in the AST is given a type. An AST node’s type is typically derived from it’s children’s types, through a process called *synthesis*, we call these types *synthesized types*. Many constructs in the language must be ascribed types by the programmer: variables declared with “let”, function parameters, and return values. Types which are declared explicitly are called *assumed types*.

```
let a: UInt32 = 1;
```

Listing 2: Simple Let statement

Here, `UInt32` is the assumed type of `x`. Where ever `x` is referenced, we can consider it of type `UInt32`. The literal `1` has no assumed type. Instead, we synthesize a type for `1` by following a set of rules. For literals, this rule is simple: *for a literal scalar N the synthesized type is $\{N\}$* . As expressions grow, synthesizing types becomes more complicated.

3.1.1 Typechecking an AST

To better illustrate this process, we’ll walk through synthesizing a tree.

```
func average(x : 0..10, y : 0..10, z : 0..10) : 0..10 {
  (x + y + z) / 3
}
```

The function parameters: `x`, `y`, and `z` have each been given the assumed types `UInt32`. An assumed type is analogous to the the statement “no matter the value of `x`, we can always assume it is a `UInt32`”. The function’s assumed return type is `UInt32`. This allows any caller to treat the expression `average(a, b, c)` as a `UInt32`, even if the operations performed by the function are unknown. An assumed type is a promise; it allows the references to entity to *assume* the type of that entity, without knowing anything else about it.

To illustrate how these assumed types interact with synthesized types, we’ll manually type check the function.

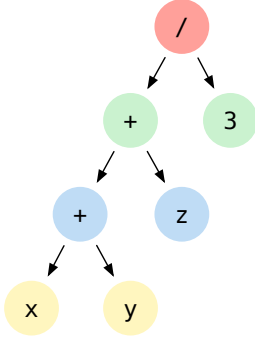
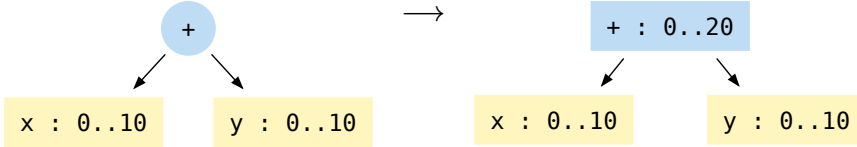


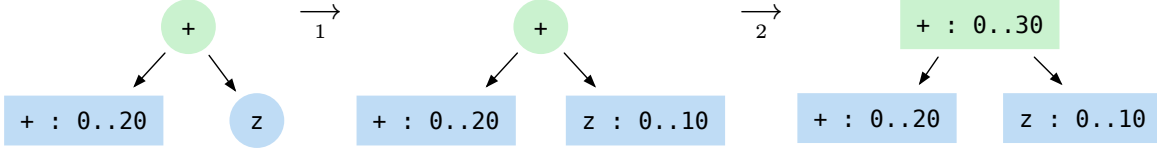
Figure 1: AST

The function body, $(x + y + z) / 3$, has the syntax tree seen in Figure 1. The type checker works bottom-up, left-to-right. So, we begin with the leaves of the tree: x , and y . Identifier AST node's synthesized type is the assumed type of the symbol they include. So x is synthesized to type $0..10$ (the assumed type of x), and y is synthesized to type $0..10$ (the assumed type of y).

This information is added to the tree, and we reference it synthesize $+$. An operator node's synthesized type is constructed by applying the given operation to the synthesized types of each operand. Types may be constructed using arithmetic operations, this process will be defined more formally in Section 3.2. For now, take for granted that $0..10 + 0..10 : 0..20$.

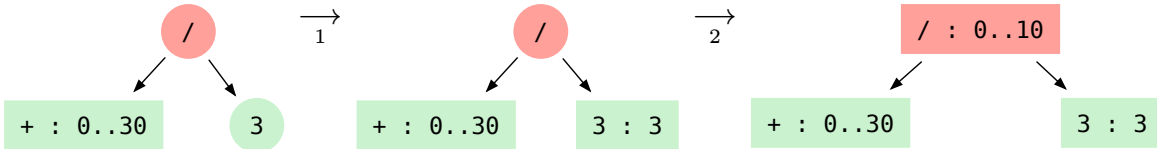


Now, we move up the tree, to synthesize the right hand side of $+$, then finally $+$ itself.



In (1) we synthesize the node's type from the assumed type of z . In (2) we used this information, and the type of $+$ to synthesize a type for $+$.

Finally, we again move up the tree, now to $/$.



Due to the the functions return value, the assumed type of the body is $0..10$. Function body's type is synthesized based on the possible return values. So the synthesized type of this function's body is the the type of $/$.

Type checking is the process of comparing assumed and synthesized types. If a synthesized is not a subset of the assumed type, then a type error is attached to that node.

3.2 Scalars

There is a single scalar type in Howlite, this simplifies the type checking by condensing many cases into a single, generic case. There are no distinct enumerable types, true boolean types, or even a unit type in the language. Instead of distinct types, we have the scalar type “Integer” (floating point number are out of scope). A scalar may be any set of Integers.

3.2.1 Synthesis of Scalars

As seen above, a scalar may be synthensized from a single value, for example the type of -5 is $\{-5\}$. We can also construct new scalars using arithmetic operations:

Given a scalar type $T = \{t_1, t_2, t_3 \dots t_n\}$, where $\forall i : t_i \in \mathbb{Z}$, and a scalar type $U = \{u_1, u_2, u_3 \dots u_n\}$ where $\forall j : u_j \in \mathbb{Z}$. (i.e T, U are subsets of the integers). We can construct the following types:

- $T \times U = \{tu : \forall t \in T, \forall u \in U\}$
- $T + U = \{t + u : \forall t \in T, \forall u \in U\}$
- $T - U = \{t - u : \forall t \in T, \forall u \in U\}$
- $T \div U = \{t \div u : \forall t \in T, \forall u \in U\}$

For example, given $T = \{1, 2, 3\}$ and $U = \{-5, -7\}$, we’d compute the following:

- $T \times U = \{1(-5), 2(-5), 3(-5), 1(-7), 2(-7), 3(-7))\} = \{-5, -10, -15, -7, -14, -21\}$
- $T + U = \{1 + -5, 2 + -5, 3 + -5, 1 + -7, 2 + -7, 3 + -7\} = \{-4, -3, -2, -6, -5, -4\}$
- $T - U = \{1 - (-5), 2 - (-5), 3 - (-5), 1 - (-7), 2 - (-7), 3 - (-7)) = \{6, 7, 8, 9, 10\}$
- $T \div U = \{1 \div (-5), 2 \div (-5), 3 \div (-5), 1 \div (-7), 2 \div (-7), 3 \div (-7)) = \{0\}$

3.2.2 Storage Classes

Scalar types belong to a *storage class* that identifies how they are encoded in memory. Storage classes are organized by size, whether or not they include a sign bit. The signed storage classes are **s8**, **s16**, **s32**, **s64**, and the unsigned are **u8**, **u16**, **u32**, **u64**. Going forward, we will identify the storage class of a scalar T using the notation **u32[T]**.

The storage class of a number influences how arithmetic and bitwise operations behave on the inner type.

3.2.2.1 Unsigned Storage Classes

given a storage class **uN**, where N is the width in bits, and variables **a** : **uN[T]**, and **b** : **uN[T]**

- $a + b = (a + b) \bmod 2^N$
- $a - b = 2^N - |a - b| \bmod 2^N$
- $a * b = (a * b) \bmod 2^N$
- $\frac{a}{b} = \frac{a - (a \bmod b)}{b}$ (i.e. division is always rounded down)
- $\sim a = (2^N - 1) - a$
- TODO other bitwise ops defined in terms of the above operations
- TODO except xor, maybe?

3.2.2.2 Signed Storage Classes

given a storage class \mathbf{uN} , where N is the width in bits, and variables $\mathbf{a} : \mathbf{sN[T]}$, and $\mathbf{b} : \mathbf{sN[T]}$

- $a + b = (a + b) \bmod 2^N$
- $a - b = 2^N - |a - b| \bmod 2^N$
- $a * b = (a * b) \bmod 2^N$
- $\frac{a}{b} = \frac{a - (a \bmod b)}{b}$ (i.e. division is always rounded down)
- $\sim a = (2^N - 1) - a$
- TODO other bitwise ops defined in terms of the above operations
- TODO except xor, maybe?

References

Dunfield, J. and Krishnaswami, N. (2020) *Bidirectional Typing*, doi: 10.48550/arXiv.1908.05839