

# The $\mu$ Specification — r0.1

## Preface

$\mu$  (stylized  $\mu$ ) is an extremely simple programming language. This document contains a complete specification of the language followed by a series of appendices elucidating motivation, design decisions, implementation strategies, and so on. These appendices are not part of the formal language specification and in the event of a conflict between the specification and the content of an appendix the specification is authoritative.

This is revision r0.1 of the specification.

## 1 Source

### 1.1 Encoding

A  $\mu$  program (recommended extensions  $\mu$ ,  $\mu\mu$ ) consists of a series of bytes. The behaviors of bytes with the high bit set is implementation defined. As a consequence a conforming implementation MAY choose to interpret  $\mu$  code as text in ASCII, UTF-8, or another similar encoding. This document will adopt such an interpretation when displaying  $\mu$  source.

### 1.2 Syntax

#### 1.2.1 Grammar

```
<program> ::= <ws> <expr> <ws>

<expr> ::= <atom> | <list>

<atom> ::= <ordinary> { <ordinary> }

<list> ::= <lparen>
          { <ws> <expr> }
          [ <ws> <dot> <ws> <expr> ] <ws>
          <rparen>

<ws> ::= { <tab> | <lf> | <cr> | <sp> }

<lparen> ::= 0x28 ;; LEFT PARENTHESIS '('
<rparen> ::= 0x29 ;; RIGHT PARENTHESIS ')'
<dot> ::= 0x2e ;; FULL STOP '.'

<tab> ::= 0x09 ;; HORIZONTAL TAB '\t'
<lf> ::= 0x0a ;; LINE FEED '\n'
```

```
<cr> ::= 0x0d ;; CARRIAGE RETURN '\r'
<sp> ::= 0x20 ;; SPACE ' '
```

```
<ordinary> ::= <any> - <special>
```

```
<any> ::= 0x20-0x7e ;; ALL PRINTABLE ASCII
<special> ::= <lparen> | <rparen> | <dot>
              | <tab> | <lf> | <cr> | <sp>
```

#### 1.2.2 Interpretation

Each distinct  $\langle atom \rangle$  in the source is an *atom name*, a sequence of bytes, and is given a distinct positive *atom number*.<sup>[1][2]</sup> How different byte sequences are mapped to numbers is completely implementation defined with the only restrictions being that:

- Atom numbers are between 1 and  $2^{31} - 1$ .
- Two atoms with the same name within the same execution of a program get the same atom number.
- Two atoms with different names within the same execution of a program get different atom numbers.

For example, all instances of the atom with name `hello` (hex 68 65 6c 6c 6f) are to be assigned the same atom number (say, 3) which must be different from the atom number for `world` (hex 77 6f 72 6c 64) (say, 7) but what exactly the atom numbers are can vary between implementations, programs, or even different executions of the same program at the implementer's discretion.

An implementation MAY assume that there are no more than  $2^{11} - 1$  distinct atom names referred to in a given source text, that each atom name has a length of at most  $2^7 - 1$  bytes, and that the combined length of these names (excluding duplicates) is no more than  $2^{15} - 2^{11}$  bytes.<sup>[3]</sup>

<sup>[1]</sup>See 2.1 Atoms.

<sup>[2]</sup>See TODO B Implementation Strategies.

<sup>[3]</sup>These restrictions permit an implementation to allocate its atom names into a singular region of memory of size `i32::MAX` separated by NULLs, for instance.

A list in the source code is constructed of cons cells<sup>[4]</sup> according to the following algorithm:

```
def parse_list(tokens):
    eat a lparen from tokens
    list = parse_list_inner(tokens)
    eat a rparen from tokens
    return list

def parse_list_inner(tokens):
    if ( the next token in tokens is dot ):
        eat a dot from tokens
        return parse an expression from tokens
    elif ( the next token in tokens is rparen ):
        return the 0 atom
    else:
        head = parse an expression from tokens
        tail = parse_list_inner(tokens)
        return a new cons cell of head and tail
```

A few things to note about this algorithm:

- `()` becomes the 0 atom. For this reason the 0 atom is also called nil and is usually denoted `()`.
- A list like  $(i_1\ i_2\ i_3)$  denotes an ordinary linked list  $i_1 \rightarrow i_2 \rightarrow i_3 \rightarrow \mathbf{NIL}$ .
- A dotted pair like  $(l\ .\ r)$  denotes a simple cons cell  $(l, r)$ .
- An ordinary list like  $(i_1\ i_2\ i_3\ \dots\ i_n)$  is equivalent to the dotted  $(i_1\ .\ (i_2\ .\ (i_3\ .\ (\dots\ (i_n\ .\ ())))))$ .

## 2 Types

$\mu$  is dynamically typed with exactly two types: atoms and cons cells.

### 2.1 Atoms

An atom is fundamentally an *unsigned 31 bit integer*.<sup>[5]</sup>

The 0 atom represents the empty list in list-related contexts.

Some positive atoms correspond to particular byte sequences (atom names) in the source code.<sup>[6]</sup> Not all atom values have a corresponding name nor is there a standard mechanism for converting between atoms and names during the execution of a program. This correspondence is simply a mechanism to allow source code to provide human-readable names to arbitrary symbols rather than a property of the language's runtime.

<sup>[4]</sup>See 2.2 Cons Cells.

<sup>[5]</sup>See B Implementation Strategies for a rationale for this unusual choice.

<sup>[6]</sup>See 1.2.2 Interpretation.

### 2.2 Cons Cells

A cons cell is simply an ordered pair of items (each of which may be either an atom or another cons cell).

The first item in the pair is called the *head* and the second is called the *tail*.

A list refers to a particular structure formed of either the 0 atom or a cons cell whose tail is itself a list.<sup>[7]</sup>

## 3 Environments

At the core of  $\mu$ 's scoping rules is the environment, a mapping from atoms to arbitrary values. This mapping is defined as a list of pairs of positive atoms and values. The first element in each pair is a variable to bind and the second is an value to bind to that variable.

For example:

```
( (hello . hello)
  (someslist . (a b c))
  (a_variable . a_value) )
```

is an environment mapping the atom `hello` to itself, the atom `someslist` to the list `(a b c)`, and the atom `a_variable` to the atom `a_value`.

An environment may contain multiple bindings for the same atom in which case the earlier binding *shadows* the later one, effectively overriding it<sup>[8]</sup>.

`()` is a legal environment containing no mappings.

### 3.1 Lookup

Looking up the value corresponding to an atom in an environment is defined according to the following algorithm:

```
def lookup(sym, env):
    if ( env is empty ):
        return sym
    else:
        binding = the head of env
        if ( the head of binding is sym):
            return the tail of binding
        else:
            return lookup(sym, the tail of env)
```

Note that duplicate bindings are resolved in favor of the first, and missing bindings resolve to themselves.

<sup>[7]</sup>See 1.2.2 Interpretation for how these are denoted.

<sup>[8]</sup>The shadowed binding is still accessible via `~env`, though.

## 4 Pattern Matching

A pattern is a value which can be matched against another value (called the scrutinee) in some base environment to produce a new environment which contains everything in the base environment in addition to new bindings representing components of the value.

Every value which does not contain duplicate positive atoms is a valid pattern.

Matching a pattern against a value follows the following rules:

- `()` matches the value `()` and introduces no bindings. Attempting to match this pattern against any other object is undefined behavior.<sup>[9]</sup>
- A positive atom matches any value and introduces a binding from the pattern atom to the value.
- A cons cell matches a cons cell value, recursively pattern matching the value's head and tail against its head and tail. The order in which the bindings from the head and tail are included in the final environment is implementation defined and may even be non-deterministic. Attempting to match a cons cell pattern against any other object is undefined behavior.
- Introduced bindings shadow bindings in the base environment (i.e. are placed before) but do not remove them. This allows careful environment manipulation to recover the base environment which is necessary for writing fully hygienic macros.

### 4.1 Examples

- Matching the expression `()` against the pattern `()` in the base environment `((a . b))` produces the environment `((a . b))`.
- Matching the expression `(a b c d)` against the pattern `(x y . z)` in the base environment `()` produces some permutation of the environment `((x . a) (y . b) (z . (c d)))`
- Matching the expression `(a b c)` against the pattern `(x y)` is undefined behavior.
- Matching the expression `(a b)` against the pattern `(x y)` in the base environment `((x . xx))` may produce either the environment `((x . a) (y . b) (x . xx))` or `((y . b) (x . a) (x . xx))` but MUST NOT produce the environment `((x . a) (y . b))`.

<sup>[9]</sup>See 7 Undefined Behavior.

## 5 Evaluation

The heart of the language is evaluation. In fact, executing a program simply consists of parsing it then evaluating it in an empty environment.

An *expression* is any value to be evaluated.

Evaluation takes place in an environment  $e$ .<sup>[10]</sup>

Evaluating an expression behaves differently depending on whether the expression is the 0 atom, a positive atom, or a non-empty list.

### 5.1 Zero Atom

The 0 atom `()` evaluates to itself.

### 5.2 Positive Atoms

A positive atom  $a$  evaluates to `lookup( $a$ ,  $e$ )`.<sup>[11]</sup>

### 5.3 Lists

If a cons cell is evaluated it should be a valid list. An implementation MAY reject attempting to evaluate a non-list cons cell or MAY try to interpret such an expression cohesively.

A list is evaluated as a call by first evaluating its head and then, based on the result (referred to as the *receiver*), performing an operation on the remaining elements.

We will use the meta syntax of angle brackets  $\langle expr \rangle_{env}$  to denote an expression which evaluates to  $expr$  in the environment  $env$  and the syntax of an arrow with environment above ( $\xrightarrow{env}$ ) to denote the evaluation of an expression to a value in an environment.

#### 5.3.1 Builtins

If the receiver is an atom it should belong to the following list of atoms with builtin behavior. Notice that with the exception of `()` which has no corresponding name these builtin receivers all have names starting `~` (hex 7e 7e).

`()` Called "quote", but represented by the 0 atom for technical reasons.<sup>[12]</sup> Evaluates to its first argument *unevaluated*.

$$(\langle () \rangle_e i) \xrightarrow{e} i$$

<sup>[10]</sup>See 3 Environments.

<sup>[11]</sup>See 3.1 Lookup for the definition of `lookup()`.

<sup>[12]</sup>Using the 0 atom as the quote builtin macro ensures that the quote macro is always accessible as the 0 atom is the only atom which cannot be rebound. This, in turn, ensures that any value is always accessible by invoking quote. The ability to access arbitrary content reliably regardless of environment is critical for hygienic macros to be possible.

<b>~true</b>	Takes two arguments and tail evaluates <sup>[13]</sup> to its first argument. DOES NOT evaluate its second argument.	<b>~sub</b>	Takes two arguments and evaluates them. They should both be atoms. Returns the difference of the two atoms modulo $2^{31}$ .
	$(\langle \sim\text{true} \rangle_e \langle t \rangle_e f_e) \xrightarrow{e} t$		$(\langle \sim\text{sub} \rangle_e \langle a \rangle_e \langle b \rangle_e) \xrightarrow{e} (a - b) \bmod 2^{31}$
<b>~false</b>	Takes two arguments and tail evaluates to its second argument. DOES NOT evaluate its first argument.	<b>~and</b>	Takes two arguments and evaluates them. They should both be atoms. Returns the bitwise conjunction (bitwise and) of the two atoms
	$(\langle \sim\text{false} \rangle_e t_e \langle f \rangle_e) \xrightarrow{e} f$		$(\langle \sim\text{and} \rangle_e \langle a \rangle_e \langle b \rangle_e) \xrightarrow{e} (a \wedge b)$
<b>~head</b>	Takes an argument and evaluates it. The result should be a cons cell. Returns the head of that cons cell.	<b>~or</b>	Takes two arguments and evaluates them. They should both be atoms. Returns the bitwise disjunction (bitwise or) of the two atoms
	$(\langle \sim\text{head} \rangle_e \langle (h . t) \rangle_e) \xrightarrow{e} h$		$(\langle \sim\text{or} \rangle_e \langle a \rangle_e \langle b \rangle_e) \xrightarrow{e} (a \vee b)$
<b>~tail</b>	Takes an argument and evaluates it. The result should be a cons cell. Returns the tail of that cons cell.	<b>~not</b>	Takes one argument and evaluates it. It should be an atom. Returns the bitwise negation (bitwise not) of the atom
	$(\langle \sim\text{tail} \rangle_e \langle (h . t) \rangle_e) \xrightarrow{e} t$		$(\langle \sim\text{not} \rangle_e \langle a \rangle_e) \xrightarrow{e} (\neg a)$
<b>~cons</b>	Takes two arguments, evaluates them, and returns a new cons cell constructed from the first and second arguments.	<b>~sl</b>	Takes two arguments and evaluates them. They should both be atoms. The second argument should be less than 31. Returns the first shifted left by the second number of bits.
	$(\langle \sim\text{cons} \rangle_e \langle h \rangle_e \langle t \rangle_e) \xrightarrow{e} (h . t)$		$(\langle \sim\text{sl} \rangle_e \langle a \rangle_e \langle b \rangle_e) \xrightarrow{e} (a << b)$
<b>~lte</b>	Takes two arguments and evaluates them. Returns ~true if the first is less than or equal to the second and ~false otherwise. Atoms are compared according to their numbers. Cons cells are always less than atoms. Two distinguishable cons cells do not compare equal. Two indistinguishable cons cells may or may not compare equal.	<b>~sr</b>	Takes two arguments and evaluates them. They should both be atoms. The second argument should be less than 31. Returns the first shifted right by the second number of bits.
	$(\langle \sim\text{lte} \rangle_e \langle a \rangle_e \langle b \rangle_e) \xrightarrow{e} \begin{cases} \sim\text{true} & \text{if } a \leq b \\ \sim\text{false} & \text{otherwise} \end{cases}$		$(\langle \sim\text{sr} \rangle_e \langle a \rangle_e \langle b \rangle_e) \xrightarrow{e} (a >> b)$
<b>~eq</b>	Takes two arguments, evaluates them, and returns ~true if the first is equal to the second and ~false otherwise. Atoms are compared according to their numbers. Cons cells are always less than atoms. Two distinguishable cons cells do not compare equal. Two indistinguishable cons cells may or may not compare equal.	<b>~env</b>	Takes no arguments. Returns the environment.
	$(\langle \sim\text{eq} \rangle_e \langle a \rangle_e \langle b \rangle_e) \xrightarrow{e} \begin{cases} \sim\text{true} & \text{if } a = b \\ \sim\text{false} & \text{otherwise} \end{cases}$		$(\langle \sim\text{env} \rangle_e) \xrightarrow{e} e$
<b>~add</b>	Takes two arguments and evaluates them. They should both be atoms. Returns the sum of the two atoms modulo $2^{31}$ .	<b>~sys</b>	Takes two arguments. DOES NOT evaluate them. Based on the value of the first argument perform a different implementation defined operation on the second argument. See 5.3.3 The ~sys Builtin for a complete definition.
	$(\langle \sim\text{add} \rangle_e \langle a \rangle_e \langle b \rangle_e) \xrightarrow{e} (a + b) \bmod 2^{31}$		$(\langle \sim\text{sys} \rangle_e \text{sym } arg) \xrightarrow{e} \text{See 5.3.3}$

<sup>[13]</sup>See 6 Required Optimizations.

<sup>[14]</sup>See 7 Undefined Behavior.

Each builtin takes a fixed number of arguments. Providing the wrong number of arguments to a builtin is undefined behavior.<sup>[14]</sup>

Attempting to use any atom that is not on this list of receivers is undefined behavior as is passing an argument of (or that evaluates to, where relevant) a type other than that expected.

### 5.3.2 User Defined Receivers

If the receiver is a list then it represents a user-defined macro or function and should have exactly 2 or 3 elements, respectively.

The first of these elements is a pattern<sup>[15]</sup> and the second is the body. The third element has two purposes:

1. Its presence indicates that this receiver is a function rather than a macro.
2. It defines an environment to use when expanding the function body. Combined with the `~env` builtin, this allows lexical scoping of functions, scoped evaluation as a derived object, and more.<sup>[16]</sup>

To evaluate a call involving a user-defined receiver preform the following steps:

1. If the receiver is a function, evaluate all of its arguments, otherwise, leave them unevaluated.
2. Match the list of arguments against the included pattern to produce a new environment.<sup>[17]</sup>
  - If the receiver is a function, use the included environment as the base when constructing the new environment.
  - Otherwise, use the calling environment as the base when constructing the new environment.
3. Tail evaluate<sup>[18]</sup> the included body in the constructed environment.

Attempting to use any cons-cell which does not form a valid list of either exactly two or three elements where the first element is a valid pattern and the third (if present) is a valid environment as a receiver is undefined behavior.<sup>[19]</sup>

### 5.3.3 The `~sys` Builtin

`~sys` is the main point of extensibility built into `μ`.

Calling the `~sys` with two 0 atoms, e.g. `((() ~sys) () ())`, should return a list of pairs mapping *system operation names* to *system operation codes*.<sup>[20]</sup> This mapping should remain constant throughout a program's lifetime.

The heads of these pairs, the *system operation names*, are atoms which correspond to human readable names of a sort a programmer may write in their source code. The tails,

meanwhile, the *system operation codes*, should also be atoms but which need not be human readable names.

In order to call a system operation, the system operation code is used as the first argument to `~sys` and the desired argument is used as the second argument. The code, argument, and calling environment are then handled in some implementation-specific way. Neither the code nor argument is evaluated prior to calling, but the system operation has all the information it needs to evaluate the argument if it wishes to.

Using both a readable name and a code and delegating the mapping to the user's program this way allows for an extremely simple implementation of `~sys` using a jump table or similar structure while still allowing readable names to be used to refer to these operations.

## 6 Required Optimizations

Several points in 5 Evaluation specify to tail evaluate some expression. In such cases implementations MUST perform proper tail-call optimization. This means that evaluating an expression which causes a tail evaluation of another expression which itself causes the tail evaluation of another expression and so on should only take up  $O(1)$  space for storage of the current evaluation state (e.g. call stack). This optimization is necessary to ensure programmers can write loops without fear of running out of space.

Implementations must not needlessly duplicate cons cells. Acquiring a cons cell from a binding multiple times should not require additional space for each instance.

Implementations must have a garbage collector. Concretely, repeatedly creating and then no longer using cons cells should not cause memory expenditure to grow without bound.

E Tests offers a list of test cases which, in addition to testing correctness of an implementation, test that these characteristics are met appropriately.

## 7 Undefined Behavior

Several points throughout this specification refer to certain occurrences as undefined behavior. Implementations are free to adopt arbitrary behavior in these instances including but not limited to:

- Extending the specification behavior
- Producing some kind of error
- Producing a nonsense result
- Halting and catching fire

<sup>[15]</sup>See 4 Pattern Matching.

<sup>[16]</sup>See TODO.

<sup>[17]</sup>See 4 Pattern Matching for the definition of this operation.

<sup>[18]</sup>See 6 Required Optimizations.

<sup>[19]</sup>See 7 Undefined Behavior.

<sup>[20]</sup>For a system which does not provide implementation specific functionality through this mechanism just return `()`.

## Appendix A

# Background

## **Appendix B**

# **Implementation Strategies**

## Appendix C

# Derivation of Higher-Level Functionality



## Appendix D

# Reference Implementation (WebAssembly)

## Appendix E

# Tests