

# The L<sub>3</sub> project

Advanced Compiler Construction  
Michel Schinz – 2016-02-25

1

## Project overview

What you will get (as the semester progresses):

- parts of an L<sub>3</sub> compiler written in Scala, and
- parts of a virtual machine, written in C.

What you will have to do:

- one non-graded exercise to warm you up,
- complete the compiler,
- complete the virtual machine.

2

# The L<sub>3</sub> language

3

## The L<sub>3</sub> language

L<sub>3</sub> is a **Lisp-like** language. Its main characteristics are:

- it is “dynamically typed”,
- it is functional:
  - functions are first-class values, and can be nested,
  - there are few side-effects (exceptions: mutable blocks and I/O),
  - it automatically frees memory,
- it has six kinds of values: unit, booleans, characters, integers, blocks and functions,
- it is simple but quite powerful.

4

## A taste of L<sub>3</sub>

An L<sub>3</sub> function to compute  $x^y$  for  $x \in \mathbb{Z}, y \in \mathbb{N}$ :

```
(defrec pow
  (fun (x y)
    (cond ((= 0 y)
           1)
          ((even? y)
           (let ((t (pow x (/ y 2))))
             (* t t)))
          (#t
           (* x (pow x (- y 1)))))))
```

$$x^0 = 1$$

$$x^{2z} = (x^z)^2$$

$$x^{z+1} = x(x^z)$$

5

## Literal values

"c<sub>1</sub>...c<sub>n</sub>"

String literal (translated to a block expression, see later).

'c'

Character literal.

... -2 -1 0 1 2 3 ...

Integer literals (also available in base 16 with #x prefix, or in base 2 with #b prefix).

#t #f

Boolean literals (true and false, respectively).

#u

Unit literal.

6

## Top-level definitions

(def n e)

Top-level non-recursive definition. The expression e is evaluated and its value is bound to name n in the rest of the program. The name n is not visible in expression e.

(defrec n f)

Top-level recursive *function* definition. The function expression f is evaluated and its value is bound to name n in the rest of the program. The function can be recursive, i.e. the name n is visible in the function expression f.

7

## Local definitions

(let ((n<sub>1</sub> e<sub>1</sub>) ...) b<sub>1</sub> b<sub>2</sub> ...)

Parallel local value definition. The expressions e<sub>1</sub>, ... are evaluated in that order, and their values are then bound to names n<sub>1</sub>, ... in the body b<sub>1</sub>, b<sub>2</sub>, ... The value of the whole expression is the value of the last b<sub>i</sub>.

(let\* ((n<sub>1</sub> e<sub>1</sub>) ...) b<sub>1</sub> b<sub>2</sub> ...)

Sequential local value definition. Equivalent to a nested sequence of let: (let ((n<sub>1</sub> e<sub>1</sub>)) (let (...)) ...)

(letrec ((n<sub>1</sub> f<sub>1</sub>) ...) b<sub>1</sub> b<sub>2</sub> ...)

Recursive local function definition. The function expressions f<sub>1</sub>, ... are evaluated and bound to names n<sub>1</sub>, ... in the body b<sub>1</sub>, b<sub>2</sub> ... The functions can be mutually recursive.

8

## Functions

(**fun** (n<sub>1</sub> ...) b<sub>1</sub> b<sub>2</sub> ...)

Anonymous function with arguments n<sub>1</sub>, ... and body b<sub>1</sub>, b<sub>2</sub>, ... The return value of the function is the value of the last b<sub>i</sub>.

(e e<sub>1</sub> ...)

Function application. Expressions e, e<sub>1</sub>, ... are evaluated in that order, and then the value of e – which must be a function – is applied to the value of e<sub>1</sub>, ...

Note : if e is a simple identifier, a special form of name resolution, based on arity, is used – see later.

9

## Conditional expressions

(**if** e<sub>1</sub> e<sub>2</sub> e<sub>3</sub>)

Two-ways conditional. If e<sub>1</sub> evaluates to a true value (i.e. anything but **#f**), e<sub>2</sub> is evaluated, otherwise e<sub>3</sub> is evaluated.

The value of the whole expression is the value of the evaluated branch.

The else branch is optional and defaults to **#u** (unit).

(**cond** (c<sub>1</sub> b<sub>1,1</sub> b<sub>1,2</sub> ...) (c<sub>2</sub> b<sub>2,1</sub> b<sub>2,2</sub> ...) ...)

N-ways conditional. If c<sub>1</sub> evaluates to a true value, evaluate b<sub>1,1</sub>, b<sub>1,2</sub> ...; else, if c<sub>2</sub> evaluates to a true value, evaluate b<sub>2,1</sub>, b<sub>2,2</sub> ...; etc. The value of the whole expression is the value of the evaluated branch or **#u** if none of the conditions are true.

10

## Logical expressions

(**and** e<sub>1</sub> e<sub>2</sub> e<sub>3</sub> ...)

Short-cutting conjunction. If e<sub>1</sub> evaluates to a true value, proceed with the evaluation of e<sub>2</sub>, and so on. The value of the whole expression is that of the last evaluated e<sub>i</sub>.

(**or** e<sub>1</sub> e<sub>2</sub> e<sub>3</sub> ...)

Short-cutting disjunction. If e<sub>1</sub> evaluates to a true value, produce that value. Otherwise, proceed with the evaluation of e<sub>2</sub>, and so on.

(**not** e)

Negation. If e evaluates to a true value, produce the value **#f**. Otherwise, produce the value **#t**.

11

## Loops and blocks

(**rec** n ((n<sub>1</sub> e<sub>1</sub>) ...) b<sub>1</sub> b<sub>2</sub> ...)

General loop. Equivalent to:

```
(letrec ((n (fun (n1 ...) b1 b2 ...)))  
  (n e1 ...))
```

(**begin** b<sub>1</sub> b<sub>2</sub> ...)

Sequential evaluation. First evaluate expression b<sub>1</sub>, discarding its value, then b<sub>2</sub>, etc. The value of the whole expression is the value of the last b<sub>i</sub>.

12

## Arity-based name lookup

A special name lookup rule is used when analyzing a function application in which the function is a simple name:

`(n e1 e2 ... ek)`

In such a case, the name `n@k` (i.e. the name itself, followed by `@`, followed by the arity in base 10) is first looked up, and used instead of `n` instead if it exists. Otherwise, name analysis proceeds as usual.

This allows a kind of overloading based on arity (although it is *not* overloading per se).

13

## Arity-based name lookup

Arity-based name lookup can for example be used to define several functions to create lists of different lengths:

```
(def list-make@1 (fun (e1) ...))  
(def list-make@2 (fun (e1 e2) ...))  
and so on for list-make@3, list-make@4, etc.
```

With these definitions, the following two function applications are both valid:

1. `(list-make 1)` (invokes `list-make@1`),
2. `(list-make 1 (+ 2 3))` (invokes `list-make@2`).

However, the following one is *not* valid, unless a definition for the bare name `list-make` also appears in scope:

```
(map list-make l)
```

14

## Values

L<sub>3</sub> integers are represented using 31 (!) bits, in two's complement, and therefore range from  $-2^{31}$  to  $2^{31} - 1$ .

L<sub>3</sub> characters represent Unicode code points, i.e. 21 bits positive integers in one of the following ranges:

- from `000016` to `D7FF16`, or
- from `E00016` to `10FFFF16`.

(Notice that characters and integers are different types, and conversion between the two must be done using the primitives `int->char` and `char->int`.)

15

## Primitives

`(@ p e1 e2 ...)`

Primitive application. First evaluate expressions `e1`, `e2`, ... in that order, and then apply primitive `p` to the value of these expressions.

L<sub>3</sub> offers the following primitives:

- integer: `<` `<=` `>` `>=` `+` `-` `*` `/` `%`
- bit vectors (integers): `<<` `>>` `&` `|` `^`
- polymorphic: `=` `!=` `id`
- type tests: `block?` `int?` `char?` `bool?` `unit?`
- character: `char->int` `int->char`
- I/O: `byte-read` `byte-write`
- tagged blocks: `block-alloc-n`
- `block-tag` `block-length` `block-get` `block-set!`

Floored integer division, e.g.  
`(/ -5 2) ⇒ -3`

identity

$0 \leq n \leq 255$

16

## Tagged blocks

L<sub>3</sub> offers a single structured datatype: tagged blocks. They are manipulated with the following primitives:

(@ `block-alloc-n s`)

Allocates an uninitialized block with tag `n` and length `s`.

(@ `block-tag b`)

Returns the tag of block `b` (as an integer).

(@ `block-length b`)

Returns the length of block `b`.

(@ `block-get b n`)

Returns the  $n^{\text{th}}$  element (0-based) of block `b`.

(@ `block-set! b n v`)

Sets the  $n^{\text{th}}$  element (0-based) of block `b` to `v`.

17

## Using tagged blocks

Tagged blocks are a low-level data structure. They are not meant to be used directly in programs, but rather as a means to implement more sophisticated data structures like strings, arrays, lists, etc.

The valid tags range from 0 to 255, inclusive. Tags  $\geq 200$  are reserved by the compiler, while the others are available for general use. (For example, our L<sub>3</sub> library uses a few tags to represent arrays, lists, etc.)

18

## Valid primitive arguments

Primitives only work correctly when applied to certain types of arguments, otherwise their behavior is undefined.

`+` `-` `*` `<<` `>>` `&` `|` `^` :  $int \times int \Rightarrow int$

`/` `%` :  $int \times \text{non-zero } int \Rightarrow int$

`<` `<=` `>` `>=` :  $int \times int \Rightarrow bool$

`=` `!=` :  $\forall \alpha, \beta. \alpha \times \beta \Rightarrow bool$

`id` :  $\forall \alpha. \alpha \Rightarrow \alpha$

`int->char` :  $int \text{ (valid Unicode code-point)} \Rightarrow char$

`char->int` :  $char \Rightarrow int$

`block?` `int?` `char?` `bool?` `unit?` :  $\forall \alpha. \alpha \Rightarrow bool$

19

## Valid primitive arguments

`byte-read` :  $\Rightarrow int \text{ between } 0 \text{ and } 255, \text{ or } -1$

`byte-write` :  $\exists \alpha. int \text{ between } 0 \text{ and } 255 \Rightarrow \alpha$

`block-alloc-n` :  $int \Rightarrow block$

`block-tag` `block-length` :  $block \Rightarrow int$

`block-get` :  $\forall \alpha. block \times int \Rightarrow \alpha$

`block-set!` :  $\forall \alpha \exists \beta. block \times int \times \alpha \Rightarrow \beta$

the return value  
is arbitrary

20

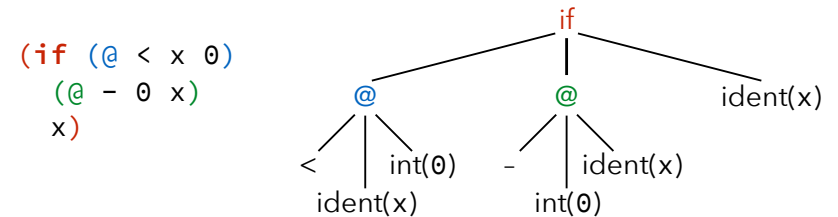
# Undefined behavior

The fact that primitives have undefined behavior when applied to invalid arguments means that they can do *anything* in such a case.  
For example, division by zero can produce an error, crash the program, or produce an arbitrary value.

21

# Grasping the syntax

Like all Lisp-like languages, L<sub>3</sub> “has no syntax”, in that its concrete syntax is very close to its abstract syntax.  
For example, the L<sub>3</sub> expression on the left is almost a direct transcription of a pre-order traversal of its AST on the right, in which nodes are parenthesized and tagged, while leaves are unadorned.



22

# L<sub>3</sub> EBNF grammar (1)

```
program ::= { def | defrec | expr } expr
def ::= (def ident expr)
defrec ::= (defrec ident fun)
expr ::= fun | let | let* | letrec | rec | begin | if | cond | and | or
        | not | app | prim | ident | num | str | chr | bool | unit
exprs ::= expr { expr }
fun ::= (fun ({ ident }) exprs)
let ::= (let ({ (ident expr) }) exprs)
let* ::= (let* ({ (ident expr) }) exprs)
letrec ::= (letrec ({ (ident fun) }) exprs)
rec ::= (rec ident ({ (ident expr) }) exprs)
begin ::= (begin exprs)
```

23

# L<sub>3</sub> EBNF grammar (2)

```
if ::= (if expr expr [ expr ])
cond ::= (cond (expr exprs) {(expr exprs)})
and ::= (and expr expr { expr })
or ::= (or expr expr { expr })
not ::= (not expr)
app ::= (expr { expr })
prim ::= (@ prim-name { expr })
```

24

## L<sub>3</sub> EBNF grammar (3)

```
str ::= "{any character except newline}"
chr ::= 'any character'
bool ::= #t | #f
unit ::= #u
ident ::= identstart { identstart | digit }[@ digit { digit }]
identstart ::= a | ... | z | A | ... | Z | | ! | % | & | * | + | -
| . | / | : | < | = | > | ? | ^ | _ | ~
prim-name ::= block-tag | block-alloc-n | etc.
```

$0 \leq n < 200$

25

## L<sub>3</sub> EBNF grammar (4)

```
num ::= num2 | num10 | num16
num2 ::= #b [-] digit2 { digit2 }
num10 ::= [-] digit10 { digit10 }
num16 ::= #x [-] digit16 { digit16 }
digit2 ::= 0 | 1
digit10 ::= digit2 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9
digit16 ::= digit10 | A | B | C | D | E | F | a | b | c | d | e | f
```

26

## Exercise

Write the L<sub>3</sub> version of the factorial function, defined as:

fact(0) = 1

fact(n) = n · fact(n - 1) [if n > 0]

What does the following (valid) L<sub>3</sub> program compute?

```
((fun (f x) (f x))
 (fun (x) (@+ x 1))
 20)
```

27

## L<sub>3</sub> syntactic sugar

28

## L<sub>3</sub> syntactic sugar

L<sub>3</sub> has a substantial amount of **syntactic sugar**: constructs that can be syntactically translated to other existing constructs. Syntactic sugar does not offer additional expressive power to the programmer, but some syntactical convenience.

For example, L<sub>3</sub> allows `if` expressions without an `else` branch, which is implicitly taken to be the unit value `#u`:  
`(if e1 e2) ⇔ (if e1 e2 #u)`

29

## Desugaring

Syntactic sugar is typically removed very early in the compilation process – e.g. during parsing – to simplify the language that the compiler has to handle.

This process is known as **desugaring**.

Desugaring can be specified as a function denoted by `[[·]]` taking an L<sub>3</sub> term and producing a desugared CL<sub>3</sub> term (CL<sub>3</sub> is *Core L<sub>3</sub>*, the desugared version of L<sub>3</sub>). To clarify the presentation, L<sub>3</sub> terms appear in orange, CL<sub>3</sub> terms in green, and meta-terms in black.

30

## L<sub>3</sub> desugaring (1)

To simplify the specification of desugaring for whole programs, we assume that all top-level expressions are wrapped sequentially in a single `(program ...)` expression.

```
[[ (program (def n e) s1 s2 ...) ]] =  
  (let ((n [[e]])) [[ (program s1 s2 ...) ]])  
[[ (program (defrec n e) s1 s2 ...) ]] =  
  (letrec ((n [[e]])) [[ (program s1 s2 ...) ]])  
[[ (program e s1 s2 ...) ]] =  
  [[ (begin e (program s1 s2 ...)) ]]  
[[ (program e) ]] =  
  [[e]]
```

31

## L<sub>3</sub> desugaring (2)

Desugaring sometimes requires the creation of **fresh names**, i.e. names that do not appear anywhere else in the program. Their binding occurrence is underlined in the rules, as illustrated by the one below.

```
[[ (begin b1 b2 b3 ...) ]] =  
  (let ((t [[b1]])) [[ (begin b2 b3 ...) ]])  
[[ (begin b) ]] =  
  [[b]]
```

32



## L<sub>3</sub> desugaring (3)

```
[[let ((n1 e1) ...) b1 b2 ...]] =  
  (let ((n1 [[e1]] ...) [(begin b1 b2 ...)]))  
[[let* ((n1 e1) (n2 e2) ...) b1 b2 ...]] =  
  [[let ((n1 e1)) (let* ((n2 e2) ...) b1 b2 ...)]]  
[[let* () b1 b2 ...]] =  
  [(begin b1 b2 ...)]  
[[letrec ((f1 (fun (n1,1 ...) b1,1 b1,2 ...) ...) ...) b1 b2 ...]] =  
  (letrec ((f1 (fun (n1,1 ...) [(begin b1,1 b1,2 ...)])  
            ...))  
    [(begin b1 b2 ...)]])
```

33

## L<sub>3</sub> desugaring (4)

```
[[fun (n1 ...) b1 b2 ...]] =  
  (letrec ((f (fun (n1 ...) [(begin b1 b2 ...)])))  
    f)  
[[rec n ((n1 e1) ...) b1 b2 ...]] =  
  (letrec ((n (fun (n1 ...) [(begin b1 b2 ...)]))  
          (n [[e1]] ...))  
    ...)  
[[e e1 ...]] =  
  [[e]] [[e1]] ...  
[[@ p e1 ...]] =  
  (@ p [[e1]] ...)
```

34

## L<sub>3</sub> desugaring (5)

```
[[if e e1]] =  
  [[if e e1 #u]]  
[[if e e1 e2]] =  
  (if [[e]] [[e1]] [[e2]])  
[[cond (e1 b1,1 b1,2 ...) (e2 b2,1 b2,2 ...) ...]] =  
  [[if e1 (begin b1,1 b1,2) (cond (e2 b2,1 b2,2) ...)]]  
[[cond ()]] =  
  #u
```

35

## L<sub>3</sub> desugaring (6)

```
[[and e1 e2 e3 ...]] =  
  [[if e1 (and e2 e3 ...) #f]]  
[[and e]] =  
  [[e]]  
[[or e1 e2 e3 ...]] =  
  [[let ((v e1)) (if v v (or e2 e3 ...))]]  
[[or e]] =  
  [[e]]  
[[not e]] =  
  [[if e #f #t]]
```

36

## L<sub>3</sub> desugaring (7)

L<sub>3</sub> does not have a string type. It offers string literals, though, which are desugared to blocks of characters.

```
["c1...cn"] =  
  (let ((s (@block-alloc-200 n)))  
    (@block-set! s 0 'c1)  
    ...  
    s)
```

[[l]] = if l is a (non-string) literal

|

[[n]] = if n is a name

n

the (reserved)  
tag 200 is used for  
strings

37

## L<sub>3</sub> desugaring example

```
[(program (@byte-write (if #t 79 75))  
           (@byte-write (if #f 79 75)))]  
= [(begin (@byte-write (if #t 79 75))  
          (program  
            (@byte-write (if #f 79 75')))))]  
= (let ((t [(@byte-write (if #t 79 75'))]))  
     [(begin  
        (program  
          (@byte-write (if #f 79 75'))))])])  
= (let ((t (@byte-write (if #t 79 75))))  
     (@byte-write (if #f 79 75)))
```

38

## Exercise

Desugar the following L<sub>3</sub> expression :

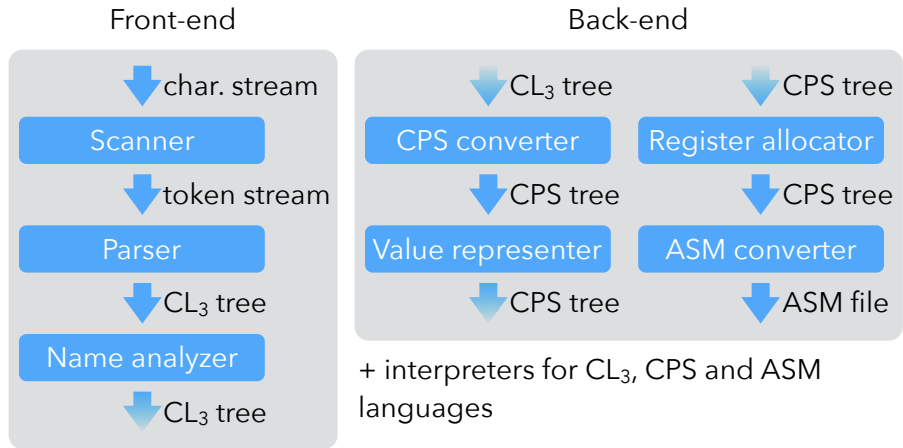
```
(rec loop ((i 1))  
  (int-print i)  
  (if (< i 9)  
      (loop (+ i 1))))
```

39

## The L<sub>3</sub> compiler

40

# L<sub>3</sub> compiler architecture



Note: CL<sub>3</sub>, CPS and ASM each designate a *family* of very similar languages, with minor differences between them.

# Intermediate languages

The L<sub>3</sub> compiler manipulates a total of four languages:

1. L<sub>3</sub> is the source language that is parsed, but never exists as a tree – it is desugared to CL<sub>3</sub> immediately,
2. CL<sub>3</sub> – a.k.a. CoreL<sub>3</sub> – is the desugared version of L<sub>3</sub>,
3. CPS is the main intermediate language, on which optimizations are performed,
4. ASM is the assembly language of the target (virtual) machine.

The compiler contains interpreters for the last three languages, which is useful to check that a program behaves in the same way as it undergoes transformation. These interpreters also serve as semantics for their language.