Values (or data) representation

Advanced Compiler Construction Michel Schinz – 2025-03-06



The problem

Values representation

- The values representation problem: how to represent the values of the source language in the target language? Trivial in C and similar languages that have:
 - no parametric polymorphism, and
 - types corresponding directly to those of the target language (e.g. int, long, double),

More difficult in languages that have either:

- parametric polymorphism, as exact types are not at compilation time, or
- dynamic types, for the same reason, or
- types not corresponding directly to those of the target.

Consider the following L₃ function: (**def** pair-make (**fun** (f s) (let ((p (@block-alloc #_pair 2))) (@block-set! p 0 f) (@block-set! p 1 s) p))) representation must be used.



The L₃ compiler knows nothing about the type of f and s, so some uniform



The same problem exists in Scala when using parametric polymorphism: def pairMake[T,U](f: T, s: U): Pair[T,U] = new Pair[T,U](f, s)

The solutions

block called a **box** and containing:

- the value,

- some information about its type. Pros and cons:

- simple,

- very costly for small values (e.g. integers).



Boxing: all values are represented uniformly by a pointer to a heap-allocated



Tagging: all values are represented uniformly by a pointer-sized word containing either:

- a pointer to a boxed value, as before, or - a small value (e.g. integer) with a tag identifying its type.
- Pros and cons:
 - simple,
 - less costly than boxing,
 - reduced range for some small values (e.g. integers).

Example: integer tagging

2n + 1.

- distinguishable from (aligned) pointer by LSB,
- slightly reduced range (1 bit less).

Integer tagging example: represent the source integer *n* as the target integer

IEEE 754 floating-point values (i.e. double) have special NaN values, returned on error, identified by top 12 bits:



12 bits (must be 1)

NaN tagging:

- represent doubles as themselves,
- use 52 lower bits of NaNs to store tagged values:
 - pointers,
 - integers,
 - etc.

Example: NaN tagging

52 arbitrary bits

On-demand boxing

- (Un)boxing can be done **on-demand** for statically-typed languages: - box when entering polymorphic context,
- unbox when returning to monomorphic context. Pros and cons:
 - no penalty for monomorphic code,
 - can be expensive at runtime.
- Also doable for dynamically-typed languages, but requires type inference.

Specialization

Specialization (or **monomorphization**): get back to simple case by

translating polymorphism away. integers is generated.

Pros and cons:

- avoids the cost of boxing and tagging,
- produces *lots* of code,
- can fail to terminate.

- For example, if List[Int] appears in a program, a class representing lists of

Partial specialization

Partial specialization:

- share specialized code as much as possible (e.g. specialize only once for all reference types), and/or
- allow the programmer to specify when to specialize, and box otherwise. Pros and cons:
 - can provide the performance of specialization for critical code without the cost.

Comparing solutions



(fully) specialized



must be adapted to the representation, e.g.:

- addition of boxed integers is done by:
 - 1. fetching the two integers from their box,
 - 2. adding them,
 - 3. allocating a new box, storing the result in it.
- addition of tagged integers is done by:
 - 1. untagging the two integers,
 - 2. adding them,
 - 3. tagging the result.

For tagging, one can do better though!

Translation of operations

- Independently of the chosen solution, operations acting on source values



Tagged integer arithmetic

- [n + m] = 2[([n] 1) / 2 + ([m] 1) / 2] + 1
 - $= (\llbracket n \rrbracket 1) + (\llbracket m \rrbracket 1) + 1$
 - = [n] + [m] 1
- [n m] = 2[([n] 1) / 2 ([m] 1) / 2] + 1
 - = ([n] 1) ([m] 1) + 1
 - = [n] [m] + 1
- $[n \times m] = 2[(([n] 1) / 2) \times (([m] 1) / 2)] + 1$
 - $= ([[n]] 1) \times (([[m]] 1) / 2) + 1$
 - $= ([n] 1) \times ([m] \gg 1) + 1$

L₃ values representation

Representation of L₃ values

L₃ has the following kinds of values:

- 1. functions,
- 2. tagged blocks,
- 3. integers,
- 4. characters,
- 5. booleans,
- 6. unit.

For now, we assume (incorrectly!) that functions are simple code pointers. Tagged blocks are represented as pointers to themselves. Integers, characters, booleans and the unit value are tagged.

L₃ tagging scheme

In L_3 , we require the two LSBs of pointers to be 0, in order to use the tagging scheme below:

Kind of v

Integer

Block (point

Character

Boolean

Unit

value	LSBs
	1 ₂
ter)	002
	110 ₂
	10102
	00102

Values representation phase

The values representation phase of the L₃ compiler:

- takes a "high-level" CPS program:
 - values: all L₃ values,
 - primitives: all L₃ primitives,
- produces an equivalent "low-level" CPS program:
 - values: bit vectors and pointers (both 32 bits),
 - primitives: instructions of the VM (similar to typical processor).
- Specified as usual as a transformation function called [.], mapping high-level CPS terms to their low-level equivalent.



```
[[n]] where n is a name =
   n
[[i]] where i is an integer literal =
   2i+1
[[c]] where c is a character literal =
   (code-point(c) \ll 3) | 110_2
[[#t]] =
   11010<sub>2</sub>
[#f] =
   01010<sub>2</sub>
[#u] =
   0010<sub>2</sub>
```

Continuations & functions

 $[(let_c ((c_1 (cnt (n_{1,1}...) e_1))...) e)] =$ $(let_{c} ((c_{1} (cnt (n_{1,1}...) [e_{1}]))...) [e])$ $[(app_c n a_1...)] =$

(app_c n [[a₁]]...)

assume the following *incorrect* translation: $[(let_f ((f_1 (f_1 (f_1 (r_1 (r_1 (r_1)) e_1)) ...) e)]] =$ $(let_f ((f_1 (f_1 (f_1 (r_1 n_{1,1} ...) [e_1])) ...) [e]))$ $[(app_f a n_c a_1 ...)] =$ $(app_f [a] n_c [a_1]...)$

Continuations are restricted enough that they don't need to be translated:

- Functions must be translated, but we ignore it for now (see next lecture) and

 $[(if (int? a) c_t c_e)] =$ (let_p ((<u>t1</u> (& [a] 1))) $(if (= t1 1) c_t c_e))$ $[(let_p ((n (+ a_1 a_2))) e)] =$ $(let* ((t1 (+ [a_1] [a_2]))))$ (n (- t1 1))) [[e]]) ... other arithmetic primitives are similar. $[(if (< a_1 a_2) c_t c_e)] =$ (if (< [[a₁]] [[a₂]]) C_t C_e) ... other integer comparison primitives are similar.



& is bit-wise and

Integers (2)

 $[(let_p ((n (block-alloc a_1 a_2))) e)] =$ (let* ((<u>t1</u> (shift-right [[a₁]] 1)) (<u>t2</u> (shift-right [[a₂]] 1)) (n (block-alloc t1 t2))) [e]) $[(let_p ((n (block-tag a))) e)] =$ (let* ((<u>t1</u> (block-tag [[a]])) $(\underline{t2} (shift-left t1 1))$ (n (+ t2 1))) [[e]]) ... other block primitives are similar.

Integers (3)

[(letp ((n (byte-read))) e)] = (let* ((t1 (byte-read))) $(\underline{t2} (shift-left t1 1))$ (n (+ t2 1))) [[e]]) [(letp ((n (byte-write a))) e)] = left as an exercise

 $[(let_p ((n (char -> int a))) e)] =$ (letp ((n (shift-right [a] 2))) [[e]]) $[(let_p ((n (int->char a))) e)] =$ (let* ((<u>t1</u> (shift-left [[a]] 2)) (n (+ t1 2))) [[e]]) $[(if (char? v) c_t c_e)] =$ left as an exercise

Characters

Booleans, unit, etc.

[(if (bool? a) ct ce)] =
 (letp ((r (& [a] 11112)))
 (if (= r 10102) ct ce))
 [(if (unit? a) ct ce)] =
 left as an exercise
 [(halt a)] =
 left as an exercise



How does the values representation phase translate the following CPS/L₃ version of the successor function? (let_f ((succ (fun (c x) (**let**_p ((t1 (+ x 1))) (**app**_c c t1))))) succ)