

# Data representation

Advanced Compiler Construction  
Michel Schinz - 2013-03-07

# The problem

# Data representation

A compiler must map the datatypes of the source language to the datatypes provided by the target language. This mapping must be as efficient as possible, in terms of both memory consumption and execution speed.

For simple languages like C, the mapping is trivial: most source datatypes (e.g. `int`, `float`, `double`, etc.) are directly mapped to the corresponding datatypes of the target machine.

For more complex languages, things are not so easy...

# The problem

Most high-level languages have the ability to manipulate values whose exact type is unknown at compilation time. This is trivially true of all “dynamically-typed” languages, but also of statically-typed languages that offer parametric polymorphism - e.g. Scala, Java 5 and later, Haskell, etc. Generating code to manipulate such values is problematic: how should values whose size is unknown be stored in variables, passed as arguments, etc. ?

# Example

To illustrate the problem of data representation, consider the following  $L_3$  function:

```
(def pair-make  
  (fun (f s) (let ((p (@block-alloc-0 2)))  
                (@block-set! p 0 f)  
                (@block-set! p 1 s) p)))
```

Obviously, nothing is known about the type of  $f$  and  $s$  at compilation time. How can the compiler generate code to pass  $f$  and  $s$  around, copy them in the allocated block, etc. given that their size is unknown?

The same problem appears in statically-typed languages offering parametric polymorphism, e.g. Scala:

```
def pairMake[T,U](f: T, s: U): Pair[T,U] =  
  new Pair[T,U](f, s)
```

# The solutions

# Uniform data representation

The simplest solution to the data representation problem is to use a **uniform data representation**. That is, all source datatypes are represented by the same target datatype. Typically, the target datatype in question is a pointer to some form of tagged block, allocated in the heap. The block contains the representation of the source data, and the tag identifies the type of that data. The data is then said to be **boxed**.

While simple, this data representation is very costly, since even small values like integers or booleans are stored in the heap. A simple operation like addition requires fetching the integers to add from their box, adding them and storing the result in a newly-allocated box...

# Tagging

To avoid paying the hefty price of uniform data representation, **tagging** can be used to represent small values like integers, booleans, characters, etc.

Tagging takes advantage of the fact that, on most target architectures, heap blocks are aligned on multiple of 2, 4 or more addressable units. Therefore, some of their least significant bits (LSBs) are always zero.

As a consequence, it is possible to use values with non-zero LSBs to represent small values.



# Integer tagging example

A popular technique to represent a source integer  $n$  as a tagged target value is to encode it as  $2n + 1$ . This ensures that the LSB of the encoded integer is always 1, which makes it distinguishable from pointers at run time.

The drawback of this encoding is that the source integers have one bit less than the integers of the target architecture. While rarely a problem, some applications become very cumbersome to write in such conditions - e.g. the computation of a 32 bits checksum.

For that reason, some languages offer several kinds of integer-like types: fast, tagged integers with reduced range, and slower, boxed integers with the full range of the target architecture.

# NaN tagging

With the current move towards 64 bits architectures, it can make sense to use 64 bits words to represent values.

Although standard tagging techniques can be used in that case, it is also possible to take advantage of a characteristic of 64 bits IEEE 754 floating-point values (i.e. **double** in Java): so-called *not-a-number* values (NaN), which are used to indicate errors, are identified only by the 12 most-significant bits, which must be 1. The 52 least-significant bits can be arbitrary.

Therefore, floating-point values can be represented naturally, and other values (pointers, integers, booleans, etc.) as specific NaN values.

# (Full) specialization

For statically-typed, polymorphic languages, **specialization** (or **monomorphization**) is another way to represent data. Specialization consists in translating polymorphism away by producing several specialized versions of polymorphic code.

For example, when the type `List[Int]` appears in a program, the compiler produces a special class that represents lists of integers - and of nothing else.

# (Full) specialization issues

Full specialization removes the need for a uniform representation of data. However, this is achieved at a considerable cost in terms of code size.

Moreover, the specialization process can loop for ever in pathological cases like:

```
class C[T]  
class D[T] extends C[D[D[T]]]
```

# Partial specialization

To reap some of the benefits of specialization without paying its full cost, it is possible to perform **partial specialization**.

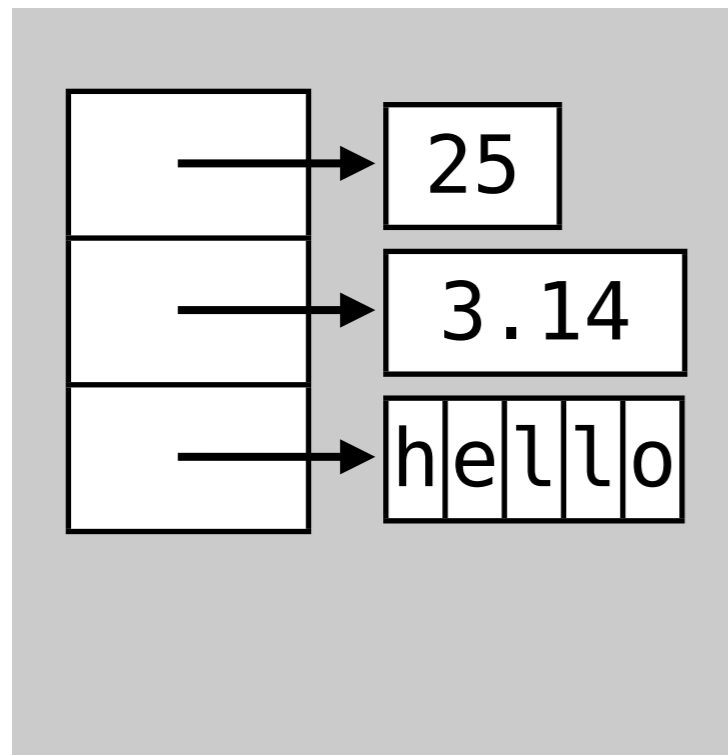
For example, since in Java all objects are manipulated by reference, there is no point in specializing for more than one kind of reference (e.g. `String` and `Object`).

It is also possible to generate specialized code only for types that are deemed important – e.g. `double` in a numerical application – and fall back to non-specialized code (with a uniform representation) for the other types. See for example Scala's `@specialized` annotation.

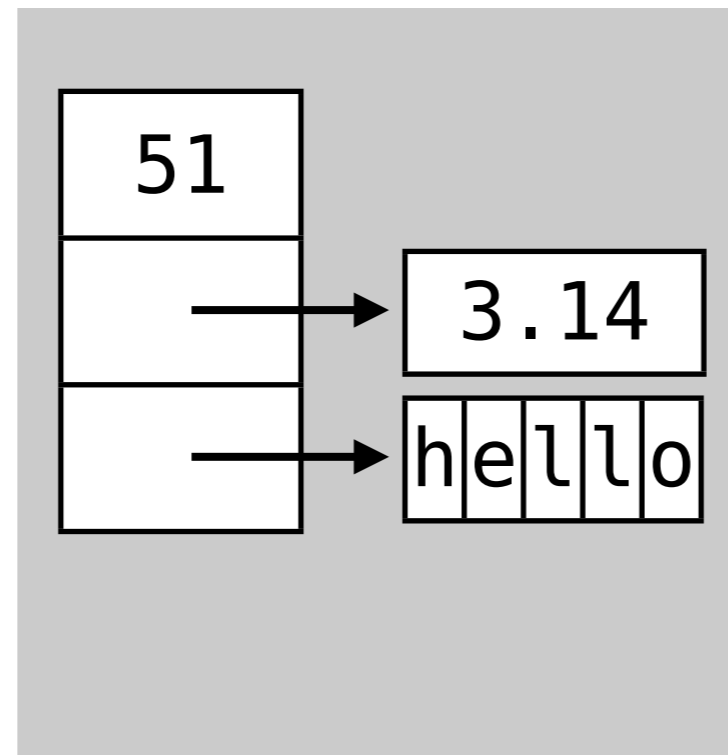
# Comparison of solutions

The following drawings show how an object containing the integer 25, the real 3.14 and the string *hello* could be represented using the three techniques previously described.

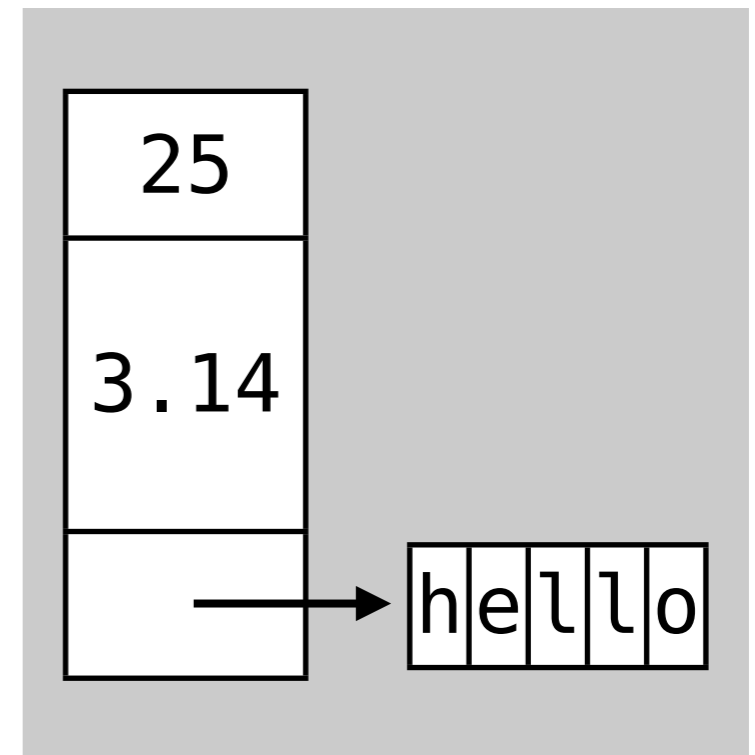
uniform



uniform with  
integer tagging



(fully) specialized



# Translation of operations

When a source datatype is encoded to a target datatype, the operations on the source datatype have to be compiled according to the encoding.

For example, if integers are boxed, addition of two source integers has to be compiled as the following sequence: the two boxed integers are unboxed, the sum of these unboxed values is computed, and finally a new box is allocated and filled with that result. This new box is the result of the addition.

Similarly, if integers are tagged, the tags must be removed before the addition is performed, and added back afterwards - at least in principle. In the case of tagging, it is however possible to do better for several operations...

# Tagged integer arithmetic

The table below illustrates how the encoded version of the four basic arithmetic primitives can be derived for tagged integers. Source terms are in red, target terms in blue.

$$\begin{aligned} \llbracket n + m \rrbracket &= 2[(\llbracket n \rrbracket - 1) / 2 + (\llbracket m \rrbracket - 1) / 2] + 1 \\ &= (\llbracket n \rrbracket - 1) + (\llbracket m \rrbracket - 1) + 1 \\ &= \llbracket n \rrbracket + \llbracket m \rrbracket - 1 \end{aligned}$$

$$\begin{aligned} \llbracket n - m \rrbracket &= 2[(\llbracket n \rrbracket - 1) / 2 - (\llbracket m \rrbracket - 1) / 2] + 1 \\ &= (\llbracket n \rrbracket - 1) - (\llbracket m \rrbracket - 1) + 1 \\ &= \llbracket n \rrbracket - \llbracket m \rrbracket + 1 \end{aligned}$$

$$\begin{aligned} \llbracket n \cdot m \rrbracket &= 2[(\llbracket n \rrbracket - 1) / 2 \cdot ((\llbracket m \rrbracket - 1) / 2)] + 1 \\ &= (\llbracket n \rrbracket - 1) \cdot ((\llbracket m \rrbracket - 1) / 2) + 1 \\ &= (\llbracket n \rrbracket - 1) \cdot (\llbracket m \rrbracket \ggg 1) + 1 \end{aligned}$$

$$\begin{aligned} \llbracket n / m \rrbracket &= 2[(\llbracket n \rrbracket - 1) / 2 / ((\llbracket m \rrbracket - 1) / 2)] + 1 \\ &= 2[(\llbracket n \rrbracket - 1) / (\llbracket m \rrbracket - 1)] + 1 \end{aligned}$$



# **L<sub>3</sub> data representation**

# L<sub>3</sub> datatypes

L<sub>3</sub> has the following datatypes: tagged blocks, functions, integers, characters, booleans, unit.

Tagged blocks are represented as pointers to themselves.

Functions are currently represented as code pointers, although we will see later that this is incorrect!

Integers, characters, booleans and the unit value are represented as tagged values.

# L<sub>3</sub> tagging scheme

For L<sub>3</sub>, we require that the two least-significant bits (LSBs) of pointers are always 0. This enables us to represent integers, booleans and unit as tagged values.

The tagging scheme for L<sub>3</sub> is given by the table below.

Datatype	LSBs
Integer	...1 <sub>2</sub>
Block (pointer)	...00 <sub>2</sub>
Character	...110 <sub>2</sub>
Boolean	...1010 <sub>2</sub>
Unit	...0010 <sub>2</sub>

# Data representation phase

The data representation phase of the  $L_3$  compiler takes as input a CPS program where all values and primitives are “high-level”, in that they have the semantics of the  $L_3$  language gives them. It produces a “low-level” version of that program as output, in which all values are either integers or pointers and primitives correspond directly to instructions of a typical microprocessor.

As usual, we will specify this phase as a transformation function called  $\llbracket \cdot \rrbracket$ , which maps high-level CPS terms to their low-level equivalent.

# Representing L<sub>3</sub> functions

L<sub>3</sub> functions also need to be represented specially so that the operations permitted on them - e.g. `function?` - can be implemented in the target language.

The phase that takes care of representing functions is commonly known as **closure conversion**. While it logically belongs to the data conversion phase, it will be presented separately in the next lecture.

For now, we assume - incorrectly - that functions are not translated specially:

$$\llbracket (\text{let}_f ((f_1 (\text{fun } (k_1 n_{1,1} \dots) e_1)) \dots) e) \rrbracket = \\ (\text{let}_f ((f_1 (\text{fun } (k_1 n_{1,1} \dots) \llbracket e_1 \rrbracket)) \dots) \llbracket e \rrbracket)$$
$$\llbracket (\text{app}_f n n_1 \dots) \rrbracket = \\ (\text{app}_f n n_1 \dots)$$

# Representing $L_3$ continuations

Continuations, unlike functions, are limited enough that they do not need to be transformed!

Their body must still be transformed recursively:

$$\llbracket (\text{let}_k ((n (\text{cont } (n_1 \dots) b))) e) \rrbracket = \\ (\text{let}_k ((n (\text{cont } (n_1 \dots) \llbracket b \rrbracket))) \llbracket e \rrbracket)$$
$$\llbracket (\text{app}_k n n_1 \dots) \rrbracket = \\ (\text{app}_k n n_1 \dots)$$

# Representing $L_3$ integers (1)

$\llbracket (\text{let}_i ((n\ i))\ e) \rrbracket$  where  $i$  is an integer literal =

$(\text{let}_i ((n\ 2^{i+1}))\ \llbracket e \rrbracket)$

$\llbracket (\text{if}\ (\text{int?}\ n)\ kt\ kf) \rrbracket =$

$(\text{let}^* ((\underline{c1}\ 1)$   
           $(\underline{t1}\ (\&\ n\ c1)))$   
           $(\text{if}\ (= t1\ c1)\ kt\ kf))$

& is bit-wise and

# Representing L<sub>3</sub> integers (2)

```
[[letp ((n (+ n1 n2))) e]] =  
  (let* ((c1 1)  
         (t1 (+ n1 n2))  
         (n (- t1 c1)))  
    [[e]])
```

... other arithmetic primitives are similar (see slide about tagged integer arithmetic).

```
[[if (< n1 n2) kt kf]] =  
  (if (< n1 n2) kt kf)
```

... other integer comparison primitives are similar.



# Representing L<sub>3</sub> integers (3)

```
[[letp ((n (block-alloc-k n1))) e]] =  
  (let* ((c1 1)  
         (t1 (>> n1 c1))  
         (n (block-alloc-k t1)))  
    [[e]])
```

>> is right shift

```
[[letp ((n (block-tag n1))) e]] =  
  (let* ((c1 1)  
         (t1 (block-tag n1))  
         (t2 (<< t1 c1))  
         (n (+ t2 c1)))  
    [[e]])
```

<< is left shift

... other block primitives are similar.

# Representing L<sub>3</sub> characters

$\llbracket (\text{let}_i ((n\ c))\ e) \rrbracket$  where  $i$  is an character literal =

$(\text{let}_i ((n\ [(char \rightarrow int(c) \ll 3) \mid 110_2]))\ \llbracket e \rrbracket)$

$\llbracket (\text{let}_p ((n\ (char-read)))\ e) \rrbracket =$

$(\text{let}^* ((\underline{c3}\ 3)$   
           $(\underline{ct}\ 110_2)$   
           $(\underline{t1}\ (char-read))$   
           $(\underline{t2}\ (<<\ t1\ c3))$   
           $(n\ (\mid\ t2\ ct)))$

$\mid$  is bit-wise or

$\llbracket e \rrbracket$ )

$\llbracket (\text{let}_p ((m\ (char-print\ n)))\ e) \rrbracket =$

*left as an exercise*

# Representing $L_3$ booleans

```
[[ (let1 ((n #t)) e ) ] ] =  
  (let1 ((n 110102)) [[e]])
```

```
[[ (let1 ((n #f)) e ) ] ] =  
  (let1 ((n 010102)) [[e]])
```

```
[[ (if (bool? n) kt kf) ] ] =  
  (let* ((m 11112)  
         (t 10102)  
         (r (& n m)))  
    (if (= r t) kt kf))
```

# Representing L<sub>3</sub> unit, halt

$\llbracket (\text{let}_1 ((n \#u)) e) \rrbracket =$   
 $(\text{let}_1 ((n \text{0010}_2)) \llbracket e \rrbracket)$

$\llbracket (\text{if (unit? n) kt kf}) \rrbracket =$

*left as an exercise*

The halt statement is left untouched by the data representation transformation.

$\llbracket (\text{halt}) \rrbracket =$   
 $(\text{halt})$

# Exercise

How does the data-representation phase translate the following CPS version of the successor function?

```
(letf ((succ (fun (k x)
              (let* ((c1 1)
                    (t1 (+ x c1))
                    (appk k t1))))))
  succ)
```