

Continuations are great for talking about semantics, but they are also useful more concretely, for programming and for compilation. We look at some of these uses here.

1 First-class continuations

Some languages expose continuations as first-class values. Examples of such languages include Scheme and SML/NJ. In the latter, there is a module that a continuation type α **cont**, representing a continuation expecting a value of type α . There are two functions for manipulating continuations:

callcc: $(\alpha \text{ cont} \rightarrow \alpha) \rightarrow \alpha$

(**callcc** f) passes the current continuation to the function f .

throw: $\alpha \text{ cont} \rightarrow \alpha \rightarrow \beta$

(**throw** k v) sends the value v to the continuation k .

Because **callcc** passes the current continuation, corresponding to the evaluation context of the **callcc** itself, invocation of the passed continuation makes the **callcc** expression itself seem to return again. It's up to the evaluation context of the **callcc** to decide whether it's seeing the original return from f or a later invocation of the passed continuation.

1.1 Semantics of first-class continuations

Using the translation approach we introduced earlier, we can easily describe these mechanisms. Suppose we represent a continuation value for the continuation k by tagging it with the integer 7. Then we can translate **callcc** and **throw** as follows:

$$\begin{aligned} \llbracket \text{callcc } e \rrbracket \rho k &= \llbracket e \rrbracket \rho (\text{CHECK-FUN}(\lambda f. f(7, k) k)) \\ \llbracket \text{throw } e_1 e_2 \rrbracket \rho k &= \llbracket e_1 \rrbracket \rho (\text{CHECK-CONT}(\lambda k'. \llbracket e_2 \rrbracket \rho k')) \end{aligned}$$

The key to the added power is the non-linear use of k in the **callcc** rule. This allows k to be reused any number of times.

1.2 Implementing threads with continuations

Once we have first-class continuations, we can use them to implement all the different control structures we might want. We can even use them to implement (non-preemptive) threads, as in the following SML/NJ-like code that explains how Concurrent ML (CML) is implemented:

```
type thread = unit cont
ready: thread queue = new_queue (* a mutable FIFO queue *)
enqueue(t) = insert ready t
dispatch() = throw (dequeue ready) ()
spawn(f: unit→unit): unit =
  callcc( fn(k) => (enqueue k; f(); dispatch()))
yield(): unit = callcc (fn(k) => enqueue k; dispatch())
```

The interface to threads is the functions **spawn** and **yield**. The **spawn** function expects a function **f** containing the work to be done in the newly spawned thread. The **yield** function causes the current thread to relinquish control to the next thread on the ready queue. Control also transfers to a new thread when one thread finishes evaluating. To complete the implementation of this thread package, we just need a queue implementation. CML has preemptive threads, in which threads implicitly yield automatically after a certain amount of time; this requires just a little help from the operating system.

2 Compiling with continuations

Because continuations expose control explicitly, they make a good intermediate language for compilation, because control is exposed explicitly in machine language as well. We can show this by writing a translation from a stripped-down version of uML to a language similar to assembly.

The result of doing such a translation is that we will have a fairly complete recipe for compiling any of the languages we have talked about into the class down to the hardware.

3 Source language

Our source language looks like the lambda calculus with tuples and numbers, with the obvious (call-by-value) semantics:

$$e ::= n \mid x \mid \lambda x.e \mid e_0 e_1 \mid (e_0, e_1, \dots, e_n) \mid (\#n e) \mid e_0 + e_1$$

The target language looks more like assembly language:

$$\begin{aligned} p &::= bb_1; bb_2; \dots; bb_n \\ bb &::= lb : c_1; c_2; \dots; c_n; \mathbf{jump} x \\ c &::= \mathbf{mov} x_1, x_2 \\ &\quad \mid \mathbf{mov} x, n \\ &\quad \mid \mathbf{mov} x, lb \\ &\quad \mid \mathbf{add} x_1, x_2, x_3 \\ &\quad \mid \mathbf{load} x_1, x_2[n] \\ &\quad \mid \mathbf{cons} x_0, x_1, \dots, x_n \end{aligned}$$

A program p consists of a series of *basic blocks* bb , each with a distinct label lb . Each basic block contains a sequence of commands c , and ends with a jump instruction. Commands correspond to assembly language instructions and are largely self-evident; the only one that is high-level is the **cons** instruction, which allocates n words of space, places the address of the space into x_0 , and puts the contents of variables (registers) $x_1 \dots x_n$ into those words.

We don't need a more general **store** instruction because the language is functional. And we aren't worrying about memory management. The **jump** instruction is an indirect jump. It makes the program counter take the value of the argument register.

4 Intermediate language 1

The first intermediate language is in continuation-passing style:

$$\begin{aligned} v &::= n \mid x \mid \lambda kx.c \mid \mathbf{halt} \\ &\quad \mid \overline{\lambda x}.c \\ e &::= v \mid v_0 + v_1 \mid (v_1, v_2, \dots, v_n) \mid (\#n v) \\ c &::= \mathbf{let} x = e \mathbf{in} c \\ &\quad \mid v_0 v_1 v_2 \\ &\quad \mid v_0 v_1 \end{aligned}$$

Some things to note about the intermediate language:

- Lambda abstractions corresponding to continuations are marked with a overline. These are considered *administrative* lambdas that we will eliminate at compile time, either by reducing them or by converting them to real lambdas.

- There are no subexpressions in the language (e does not occur in its own definition).
- Commands c look a lot like basic blocks:

```

let  $x_1 = e_1$  in
  let  $x_2 = e_2$  in
    ...
    let  $x_n = e_n$  in
       $v_0 v_1 v_2$ 

```

- Lambdas are not closed and can occur inside other lambdas.

The contract of the translation is that $\llbracket e \rrbracket k$ will evaluate e and pass its result to the continuation k . To translate an entire program, we use $k = \mathbf{halt}$, where \mathbf{halt} is the continuation to send the result of the entire program to. Here is the translation from the source to the first intermediate language:

$$\begin{aligned}
\llbracket x \rrbracket k &= k x \\
\llbracket n \rrbracket k &= k n \\
\llbracket \lambda x. e \rrbracket k &= k (\lambda x k'. (\llbracket e \rrbracket k')) \\
\llbracket e_0 e_1 \rrbracket k &= \llbracket e_0 \rrbracket (\bar{\lambda} f. \llbracket e_1 \rrbracket (\bar{\lambda} v. (f v k))) \\
\llbracket (e_1, e_2, \dots, e_n) \rrbracket k &= \llbracket e_1 \rrbracket (\bar{\lambda} x_1. \dots \llbracket e_n \rrbracket (\bar{\lambda} x_n. \mathbf{let} t = (x_1, x_2, \dots, x_n) \mathbf{in} (k t))) \\
\llbracket \#n e \rrbracket k &= \llbracket e \rrbracket (\bar{\lambda} t. \mathbf{let} y = \#n t \mathbf{in} (k y)) \\
\llbracket (e_1 + e_2) \rrbracket k &= \llbracket e_1 \rrbracket (\bar{\lambda} x_1. \llbracket e_2 \rrbracket (\bar{\lambda} x_2. \mathbf{let} z = x_1 + x_2 \mathbf{in} (k z)))
\end{aligned}$$

Let's see an example. We translate the expression $\llbracket (\lambda a. (\#1 a)) (3, 4) \rrbracket k$, using $k = \mathbf{halt}$.

$$\begin{aligned}
&\llbracket (\lambda a. (\#1 a)) (3, 4) \rrbracket k \\
&= \llbracket \lambda a. (\#1 a) \rrbracket (\bar{\lambda} f. \llbracket (3, 4) \rrbracket (\bar{\lambda} v. (f v k))) \\
&= (\bar{\lambda} f. \llbracket (3, 4) \rrbracket (\bar{\lambda} v. (f v k))) (\lambda a k'. \llbracket \#1 a \rrbracket k') \\
&= (\bar{\lambda} f. \llbracket 3 \rrbracket (\bar{\lambda} x_1. \llbracket 4 \rrbracket (\bar{\lambda} x_2. \mathbf{let} b = (x_1, x_2) \mathbf{in} (\bar{\lambda} v. (f v k) b)))) (\lambda a k'. \llbracket \#1 a \rrbracket k') \\
&= (\bar{\lambda} f. (\bar{\lambda} x_1. (\bar{\lambda} x_2. \mathbf{let} b = (x_1, x_2) \mathbf{in} (\bar{\lambda} v. (f v k) b) 4) 3) (\lambda a k'. \llbracket \#1 a \rrbracket k')) \\
&= (\bar{\lambda} f. (\bar{\lambda} x_1. (\bar{\lambda} x_2. \mathbf{let} b = (x_1, x_2) \mathbf{in} (\bar{\lambda} v. (f v k) b) 4) 3) (\lambda a k'. \llbracket a \rrbracket (\bar{\lambda} t. \mathbf{let} y = \#1 t \mathbf{in} k' t)))
\end{aligned}$$

Clearly, the translation generates a lot of administrative lambdas, which will be quite expensive if they are compiled into machine code. To make the code more efficient and compact, we will optimize it using some simple rewriting rules to eliminate administrative lambdas.

$\bar{\beta}$ -Reduction

We can eliminate unnecessary application to a variable, by copy propagation:

$$(\bar{\lambda} x. e) y \longrightarrow e\{y/x\}$$

Other unnecessary administrative lambdas can be converted into **lets**:

$$(\bar{\lambda} x. c)v \longrightarrow \mathbf{let} x = v \mathbf{in} c$$

We can also perform administrative η -reductions:

$$\bar{\lambda} x. k x \longrightarrow k$$

If we apply these rules to the expression above, we get

```

let f = (λk'a. let y = #1 a in k' y) in
  let x1 = 3 in
    let x2 = 4 in
      let x3 = (x1, x2) in
        f b k

```

This is starting to look a lot more like our target language.

The idea of separating administrative terms from real terms and performing a compile-time *partial evaluation* is powerful and can be used in many other contexts. Here, it allows us to write a very simple CPS conversion that treats all continuations uniformly, and perform a number of control optimizations.

Note that we may not be able to remove all administrative lambdas. Any that cannot be reduced using the rules above are converted into real lambdas.

4.1 Tail call optimization *

A tail call is a function call that determines the result of another function. A tail-recursive function is one whose recursive calls are all tail calls. Continuations make tail calls easy to optimize. For example, the following program has a tail call from f to g :

```

let g = λx. #1 x in
  let g = λx. g x
  in
    f(2)

```

The translation of the body of f is $(g (\overline{\lambda y. k' y}) x)$, which permits optimization by η -reduction to $(g k x)$. In this optimized code, g does not bother to return to f , but rather jumps directly back to f 's caller. This is an important optimization for functional programming languages, where tail-recursive calls that take up linear stack space are converted into loops that take up constant stack space.

5 Intermediate Language 1 \rightarrow Intermediate Language 2

The next step is the translation from Intermediate language 1 to Intermediate Language 2. In this intermediate language, all lambdas are at the top level, with no nesting:

$$\begin{aligned}
 P & ::= \text{let } x_f = \lambda k x_1 \dots x_n. c \text{ in } P \\
 & \quad | \text{let } x_c = \lambda x_1 \dots x_n. c \text{ in } P \\
 & \quad | c \\
 c & ::= \text{let } x = e \text{ in } c \mid x_0 x_1 \dots x_n \\
 e & ::= n \mid x \mid \mathbf{halt} \mid x_1 + x_2 \mid (x_1, x_2, \dots, x_n) \mid \#n x
 \end{aligned}$$

The translation requires the construction of closures that capture all the free variables of the lambda abstractions in intermediate language 1. We have covered closure conversion earlier; it too can be phrased as a translation that exploits compile-time partial evaluation.

6 Intermediate Language 2 \rightarrow Assembly

The translation is given below. Note: ra is the name of the dedicated register that holds the return address.

$$\begin{aligned}
 \mathcal{P}[[p]] &= \text{program for } p \\
 \mathcal{C}[[c]] &= \text{sequence of commands } c_1; c_2; \dots; c_n \\
 \\
 \mathcal{P}[[s]] &= \text{main} : \mathcal{C}[[c]]; \text{halt} : \\
 \mathcal{P}[[\text{let } x_f = \lambda k x_1 \dots x_n. c \text{ in } p]] &= x_f : \text{mov } k, ra; \\
 &\quad \text{mov } x_1, a_1; \\
 &\quad \vdots \\
 &\quad \text{mov } x_n, a_n; \\
 &\quad \mathcal{C}[[c]]; \\
 &\quad \mathcal{P}[[p]] \\
 \mathcal{P}[[\text{let } x_c = \lambda x_1 \dots x_n. c \text{ in } p]] &= x_c : \text{mov } x_1, a_1; \\
 &\quad \vdots \\
 &\quad \text{mov } x_n, a_n; \\
 &\quad \mathcal{C}[[c]]; \\
 &\quad \mathcal{P}[[p]] \\
 \mathcal{C}[[\text{let } x_1 = x_2 \text{ in } c]] &= \text{mov } x_1, x_2; \mathcal{C}[[c]] \\
 \mathcal{C}[[\text{let } x_1 = x_2 + x_3 \text{ in } c]] &= \text{add } x_1, x_2, x_3; \mathcal{C}[[c]] \\
 \mathcal{C}[[\text{let } x_0 = (x_1, x_2, \dots, x_n) \text{ in } c]] &= \text{cons } x_0, x_1, x_2, \dots, x_n; \mathcal{C}[[c]] \\
 \mathcal{C}[[\text{let } x_1 = \#n x_2 \text{ in } c]] &= \text{load } x_1, x_2[n]; \mathcal{C}[[c]] \\
 \mathcal{C}[[x_0 k x_1 \dots x_n]] &= \text{mov } ra, k; \\
 &\quad \text{mov } a_1, x_1; \\
 &\quad \vdots \\
 &\quad \text{mov } a_n, x_n; \\
 &\quad \text{jump } x_0
 \end{aligned}$$

At this point, we are still assuming an infinite supply of registers. We need to do register allocation and possibly spill registers to a stack to obtain working code.

While this translation is very simple, it is possible to do a better job of generating calling code. For example, we are doing a lot of register moves when calling functions and when starting the function body. These could be optimized.