

Hoisting, Allocation and Code Generation

Di Zhao

zhaodi01@mail.ustc.edu.cn

Document for the hoisting and allocation phases.

1 Hoisting

Hoist all the function definitions to top-level. Source language : closure passing language.

(terms)	K	\rightarrow	letval $x = V$ in K
			let $x = \#i$ y in K
			letcont k env $x = K$ in K'
			k env x
			f env k x
			case x of $\overrightarrow{\text{ini}_j x_j \Rightarrow K_j}$
			letprim $x = \text{PrimOp } \vec{y}$ in K
			if x then K else K'
			letfix f env k $x = K$ in K'
(values)	V	\rightarrow	()
			true
			false
			i
			" s "
			(x_1, x_2, \dots, x_n)
			ini x
			$\lambda \text{env } k$ $x.K$
(primitive operations)	PrimOp	\rightarrow	$+$ $-$ $*$
			$>$ $<$ $=$
			andalso orelse not
			print int2string

Figure 1: Closure syntax

The target language in this section is called Flat. The syntax is shown in Figure 2.

In the Flag language, a program p consists of a main function (denoted as m) and a list of functions (\vec{f}).

(program)	p	\rightarrow	$m; \vec{f}$
(functions)	m, f	\rightarrow	$x \vec{y} \{ \vec{b}; e; \}$
(bindings)	b	\rightarrow	$x = v$
(values)	v	\rightarrow	$()$ \mid true \mid false \mid i \mid $"s"$ \mid $\pi_i x$ \mid (x_1, x_2, \dots, x_n) \mid $\text{in}_i x$ \mid $\text{PrimOp } \vec{x}$
(primitive operations)	PrimOp	\rightarrow	\mid $+$ \mid $-$ \mid $*$ \mid $>$ \mid $<$ \mid $=$ \mid andalso \mid orelse \mid not \mid print \mid int2string
(transfers)	e	\rightarrow	$x \vec{y}$ \mid if x then e_1 else e_2 \mid case x of $\text{in}_{i_j} x_j \Rightarrow e_j$

Figure 2: Flat syntax

In a function definition $x \vec{y} \{ \vec{b}; e; \}$, x is the function name, \vec{y} is the list of arguments. The function body consists of a list of bindings \vec{b} and one control transferring expression e .

A binding binds a value v to a name x . Here the values to be named include empty value, constant integers, constant strings, projection operation, tuples, tagged values and primary operations.

There're three kinds of control transferring expressions (denoted as e).

The hoisting rules are illustrated in Figure 3.

Updates:

15-7-3: Change the form of **case** in **flat.sig**, **flat.sml**, and changed the respective hoisting rule in **hoist.sml** to enable multiple cases .

15-9-5: Added boolean values and corresponding operations. Changed **if0** expression into **if** expression.

2 Allocation

After the allocation process, the allocation and initialization of tuples and tagged values will be explicit. Besides, the arguments of a single function are packed

$$\begin{array}{c}
\frac{\vdash K \rightsquigarrow (\vec{f}_1, \vec{b}_1, e_1) \quad \vdash K' \rightsquigarrow (\vec{f}_2, \vec{b}_2, e_2)}{\vdash \text{letval } x = \lambda env \ k \ z. K \text{ in } K' \rightsquigarrow (\vec{f}_1 :: \vec{f}_2, :: f_x, \vec{b}_2, e_2)} \text{ (H-LETVALFUNC)} \\
\text{where } f_x \text{ is function: } x [env, k, z] \{ \vec{b}_1; e_1; \} \\
\\
\frac{\vdash K \rightsquigarrow (\vec{f}, \vec{b}, e)}{\vdash \text{letval } x = V \text{ in } K \rightsquigarrow (\vec{f}, (x = V) :: \vec{b}, e)} \text{ (H-LETVAL)} \\
\\
\frac{\vdash K \rightsquigarrow (\vec{f}, \vec{b}, e)}{\vdash \text{let } x = \pi_i \ y \text{ in } K \rightsquigarrow (\vec{f}, (x = \pi_i \ y) :: \vec{b}, e)} \text{ (H-LET)} \\
\\
\frac{\vdash K \rightsquigarrow (\vec{f}_1, \vec{b}_1, e_1) \quad \vdash K' \rightsquigarrow (\vec{f}_2, \vec{b}_2, e_2)}{\vdash \text{lecont } k \ env \ x = K \text{ in } K' \rightsquigarrow (\vec{f}_1 :: \vec{f}_2, :: f_k, \vec{b}_2, e_2)} \text{ (H-LETCONT)} \\
\text{where } f_k \text{ is function: } k [env, x] \{ \vec{b}_1; e_1; \} \\
\\
\frac{}{\vdash k \ env \ x \rightsquigarrow ([\], [\], k [env, x])} \text{ (H-CONTAPP)} \\
\\
\frac{}{\vdash f \ env \ k \ x \rightsquigarrow ([\], [\], f [env, k, x])} \text{ (H-FUNCAPP)} \\
\\
\frac{\vdash K_j \rightsquigarrow (\vec{f}_j, \vec{b}_j, e_j)}{\vdash \text{case } x \text{ of } \overline{\text{ini}_j \ x_j \Rightarrow K_j} \rightsquigarrow (\vec{f}_1 :: \dots :: \vec{f}_n, \vec{b}_1 :: \dots :: \vec{b}_n, \text{case } x \text{ of } \overline{\text{ini}_j \ x_j \Rightarrow e_j})} \text{ (H-CASE)} \\
\text{where } j = 1, \dots, n \\
\\
\frac{\vdash K \rightsquigarrow (\vec{f}, \vec{b}, e)}{\vdash \text{letprim } x = \text{PrimOp } \vec{y} \text{ in } K \rightsquigarrow (\vec{f}, (x, \text{PrimOp } \vec{y}) :: \vec{b}, e)} \text{ (H-LETPRIM)} \\
\\
\frac{\vdash K \rightsquigarrow (\vec{f}_1, \vec{b}_1, e_1) \quad \vdash K' \rightsquigarrow (\vec{f}_2, \vec{b}_2, e_2)}{\vdash \text{if } x \text{ then } K \text{ else } K' \rightsquigarrow (\vec{f}_1 :: \vec{f}_2, \vec{b}_1 :: \vec{b}_2, \text{if } x \text{ then } e_1 \text{ else } e_2)} \text{ (H-IF)} \\
\\
\frac{\vdash K \rightsquigarrow (\vec{f}_1, \vec{b}_1, e_1) \quad \vdash K' \rightsquigarrow (\vec{f}_2, \vec{b}_2, e_2)}{\vdash \text{letfix } f \ env \ k \ x = K \text{ in } K' \rightsquigarrow (\vec{f}_1 :: \vec{f}_2, :: f_f, \vec{b}_2, e_2)} \text{ (H-LETFIX)} \\
\text{where } f_f \text{ is function: } f [env, k, x] \{ \vec{b}_1; e_1; \}
\end{array}$$

Figure 3: Hoisting rules for Closure.t

into one single argument, so that it would be easier to scan over the arguments when performing garbage collection. The target language in this phase is called: Machine. The syntax of Machine language is shown in Figure 4.

Given the Machine syntax, we can decide the conversion rules intuitively, as shown in Figure 5, 6, 7.

(program)	P	\rightarrow	$m; \vec{f}$
(functions)	m, f	\rightarrow	$x\ y\ \{ \overrightarrow{\text{binding};\ b}; \}$
(bindings)	binding	\rightarrow	$x = v$
			$ \text{dst}[i] = \text{src}$
(values)	v	\rightarrow	null
			$ $ true
			$ $ false
			$ i$
			$ \text{"s"}$
			$ x[i]$
			$ \text{allocTuple}(i)$
			$ \text{allocTag}(i)$
			$ \text{PrimOp } \vec{x}$
(primitive operations)	PrimOp	\rightarrow	$+ \mid - \mid *$
			$ > \mid < \mid =$
			$ \text{andalso} \mid \text{orelse} \mid \text{not}$
			$ \text{print} \mid \text{int2string}$
(blocks)	b	\rightarrow	$\overrightarrow{\text{binding};\ e}$
(transfers)	e	\rightarrow	$x\ y$
			$ \text{if } x \text{ then } \overrightarrow{b_1} \text{ else } \overrightarrow{b_2}$
			$ \text{case } x \text{ of } \overrightarrow{\text{in}_i x_j \Rightarrow b_j}$

Figure 4: Machine syntax

As for bindings, because tuples and tagged values are structured data elements, we need to allocate space for them and perform some initializations. In this way, both of them require special treatment:

- For $x = (x_1, \dots, x_n)$, allocate a tuple x with n components and initialize component $x[i]$ with x_i .
- For $x = \text{in}_i y$, allocate structure x with tag value i and initialize component $x[1]$ with y .

Arguments of a function call now need to be packed into a tuple, and the tuple needs to be initialized by these arguments. In this way, control transferring expressions will be translated into a *block* which consists of a list of bindings and an expression.

For the same reason above, in a function definition, the original arguments need to be fetched out from the sole parameter of the function. Namely we need to add some bindings ahead of the original function body.

Pay attention that as the main function m takes no arguments, no additional bindings for argument folding is needed.

Updates:

$$\begin{aligned}
\mathcal{A}_b : \text{string} * \text{Flat.binding} &\rightarrow \text{Machine.binding list} \\
\mathcal{A}_b(x = ()) &= [x = \text{null}] \\
\mathcal{A}_b(x = \text{true}) &= [x = \text{true}] \\
\mathcal{A}_b(x = \text{false}) &= [x = \text{false}] \\
\mathcal{A}_b(x = i) &= [x = i] \\
\mathcal{A}_b(x = "s") &= [x = "s"] \\
\mathcal{A}_b(x = \pi_i y) &= [x = y[i]] \\
\mathcal{A}_b(x = (x_1, x_2, \dots, x_n)) &= [x = \text{allocTuple}(n), \\
&\quad x[1] = x_1, \\
&\quad \dots \\
&\quad x[n] = x_n] \\
\mathcal{A}_b(x = \text{in}_i y) &= [x = \text{allocTag}(i), \\
&\quad x[1] = y] \\
\mathcal{A}_b(x = \text{PrimOp } \vec{y}) &= [x = \text{PrimOp } \vec{y}]
\end{aligned}$$

Figure 5: Conversion from Flat to Machine for bindings

$$\begin{aligned}
\mathcal{A}_e : \text{Flat.e} &\rightarrow \text{Machine.block} \\
\mathcal{A}_e(x \vec{y}) &= [z = \text{allocTuple}(n), \\
&\quad z[1] = y_1, \\
&\quad \dots \\
&\quad z[n] = y_n]; \\
&\quad x \ z \\
&\quad (\text{given that } \vec{y} = [y_1, \dots, y_n]) \\
\mathcal{A}_e(\text{if } x \text{ then } e_1 \text{ else } e_2) &= []; \\
&\quad \text{if } x \text{ then } \mathcal{A}_e(e_1) \text{ else } \mathcal{A}_e(e_2) \\
\mathcal{A}_e(\text{case } x \text{ of } \overrightarrow{\text{ini}_j x_j \Rightarrow e_j}) &= []; \\
&\quad \text{case } x \text{ of } \overrightarrow{\text{ini}_j x_j \Rightarrow \mathcal{A}_e(e_j)}
\end{aligned}$$

Figure 6: Conversion from Flat to Machine for expressions

- 15-7-4:** Change the form of `case` in `machine.sig`, `machine.sml`, and changed the respective allocation rule in `codegen.sml` to enable multiple cases .
- 15-9-5:** Added boolean values and corresponding operations. Changed `if0` expression into `if` expression.

$$\mathcal{A}_f(x \vec{y} \{ \vec{b}; e; \}) = x(z) \{$$

$$\begin{array}{l} \overrightarrow{y_i = z[i]} :: \\ \overrightarrow{\mathcal{A}_b(b)}; \\ \mathcal{A}_e(e); \end{array}$$

$$\}$$

(given that $\vec{y} = [y_1, \dots, y_n]$, $i = 1, \dots, n$)

Figure 7: Conversion from Flat to Machine for functions

3 Code Generation

Code generation procedure: encode the functions to output the Machine syntax tree into a `.c` file which can be compiled by a C compiler and then be executed.

As is shown Figure 4, the Machine syntax is basically C style syntax. Consequently we can output C code for most of the syntax terms directly without making much modifications. However, there are a few aspects to which you may pay attention:

- Type and type-casting. As C is strong typed, you'll have to declare variables and insert type-castings with explicit types. For example, you can declare all the variables with `int*` type; when performing a function call, you may cast the function name into `void * (*)()` type first.
- Bindings of tuples and tagged values are encoded as function calls to `allocTuple` and `allocTag` respectively. You'll need to implement these two functions in the runtime system.
- A case expression: $\text{case } x \text{ of } \overrightarrow{ini_j x_j \Rightarrow b_j}$ can be encoded as an switch structure. The variables x_j should be initialized (by what value?) before the switch structure.

To support the output code, we need to provide a runtime system, including the definition of function `allocTuple` and `allocTag`. Therefore, we need to decide the memory map for these structures.

As we rely on garbage collection to retrieve memory space, some additional information needs to be attached, such as the length of a tuple, or to distinguish between integers, tuples and tagged values. We will discuss about this with more details in the next lab. For now we just use a simple strategy shown in Figure 8 and 9.

As is shown in Figure 8, a tuple with n components will occupy $n+1$ memory cells. The first cell is to store the length of the tuple, the rest n cells store the component values in order. In this lab, each cell is constraint to be 4 Bytes. In this way, `allocTuple(n)` will allocate a memory space of $4 * (n + 1)$ Bytes and

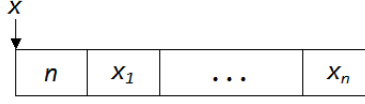


Figure 8: memory map for tuple $x = (x_1, \dots, x_n)$

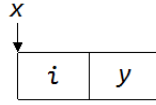


Figure 9: memory map for tagged value $x = \text{in}_i y$

store the value n in the first cell. Finally it returns a pointer to the first cell. The value of the components will be loaded later via this pointer.

Note that we will never visit a tuple with the length of 0, so we can simply output " $\mathbf{x}=\mathbf{0}$ " instead of " $\mathbf{x}=\text{allocTuple}(0)$ ".

As shown in Figure 9, a tagged value $\text{in}_i y$ takes two memory cells. The first cell is to store the tag value i (1 or 2), while the second one for the value to be tagged. In this way, `allocTag(i)` will allocate a memory space of 8kb and store the tag value i in the first cell. And finally it will return a pointer to the first cell. The value y will be loaded later via this pointer.