

Back-End Code Generation

- Given a list of tree fragments, how to generate the corresponding assembly code?

```

datatype frag
= PROC of {name : Tree.label,      function name
          body : Tree.stm,         function body/tree
          frame : Frame.Frame}     stack frame layout
| DATA of string
    
```

- Main challenges: certain aspects of tree statements and expressions do not correspond exactly with machine languages:

of temp. registers on real machines are limited

real machine's conditional-JUMP statement takes only one label

high-level constructs **ESSEQ** and **CALL** — side-effects

Itree Stmts and Exprs

- tree statements **stm** and tree expressions **exp**

```

datatype stm = SEQ of stm * stm
            | LABEL of label
            | JUMP of exp
            | CJUMP of test * label * label
            | MOVE of exp * exp
            | EXP of exp

and exp = BINOP of binop * exp * exp
        | CVTOP of cvtop * exp * size * size
        | MEM of exp * size
        | TEMP of temp
        | NAME of label
        | CONST of int
        | CONSTF of real
        | ESEQ of stm * exp
        | CALL of exp * exp list
    
```

Side-Effects

- Side-effects means updating the contents of a memory cell or a temporary register. What are tree expressions that might cause side effects? **ESEQ** and **CALL** nodes

- ESEQ**(**s**, **e**) where **s** is a list of statements that may contain **MOVE** statement

The natural way to generate assembly code for **ESEQ**(**s**, **e**)

```

Instructions to compute e1 into r1;
Instructions to compute e2 into r2;
r1 ← r1 op r2
    
```

But it won't work for this:

```
ESEQ(PROC, CONST a, ESEQ(MOVE(CONST a, r), r)
```

- CALL**(**s**, **r1**) by default puts the result in the return-result register.

```
BINOP(CALL, CALL(r, s), CALL(r, s))
```

Summary: IR → Machine Code

- Step #1 : Transform the tree code into a list of canonical trees
 - eliminate **SEQ** and **ESEQ** nodes
 - the arguments of a **CALL** node should never be other **CALL** nodes — the parent of each **CALL** node should either be **EXP**(**...**) or **MOVE**(**TEMP** **t**, **...**)
- Step #2 : Perform various code optimizations on canonical trees
- Step #3 : Rearrange the canonical trees (into traces) so that every **CJUMP**(**cond**, **L₁**, **L₂**) is immediately followed by **LABEL**(**L₂**)
- Step #4 : Instruction Selection — generate the pseudo-assembly code from the canonical trees in the step #3.
- Step #5 : Perform register allocations on pseudo-assembly code

Canonical Trees

- A canonical tree is a simple tree statement in the following form (it is really a restricted-kind of tree statement):

```

treeType stmt - LABEL nC label
| JUMP nC exp
| CJUMP nC cond * label * label
| MOVE nC exp * exp
| EXP nC exp

and exp - ESEQ nC ESeq * exp * exp
| CSEQ nC nCSeq * exp * stmt * stmt
| SEQ nC exp * stmt
| SEQ nC Seq
| LABEL nC label
| CONST nC int
| CONST nC real
| CALL nC exp * exp list

```

Restrictions:

- no SEQ statements, no ESEQ expressions.
- each CALL node doesn't contain other CALL nodes as subtrees

Canonicalizer

- The body of each PROC fragment is translated into an ordered list of canonical trees

```
stmt: Tree.stm list
```

- Step 1: transformation on CALL nodes.

```

CALL(---) ----->
    ESEQ(MOVE(TEMP C,CALL(---)),TEMP C)

```

- Step 2: elimination of ESEQ nodes. (see Appel pp 174-179)

lift them higher and higher until they become SEQ nodes ...

Rearranging itree statements

- Goal: rearrange the list of canonical trees so that every `JUMP (cond, L_1 , L_2)` is immediately followed by its false branch `LABEL (L_2)`.
- Step #1: take a list of canonical trees and form them into basic blocks

A basic block is a sequence of statements that is always entered at the beginning and exited at the end:

- the first statement is a LABEL
- the last statement is a JUMP or CJUMP
- there are no other LABELS, JUMPS, or CJUMPS in between

basic blocks are often used to analyze a program's control flow

- Step #2: re-order the list of basic blocks into traces

Canonical Trees => Basic Blocks

- Input: a sequence of statements (i.e., canonical trees — the body of a function); Output: a set of basic blocks

- Algorithm:

if a new LABEL is found, end the current block and start a new block;

if a JUMP or CJUMP is found, end the current block;

if it results a block not ending with a JUMP or CJUMP, then a JUMP to the next block's label is appended to the block;

if it results a block without a LABEL at the beginning, invent a new LABEL and stick it there;

invent a new label done for the beginning of the epilogue;

put JUMP(NAME done) at the end of the last basic block.

Basic Blocks => Traces

- Control Flow Graph (CFG): basic blocks as the nodes, pairs (a,b) as the edges if block a ends with a `goto` or `goto` statement to block b.
- Basic blocks can be arranged in any order, but we want:
 - each `goto` is followed by its false label
 - each `goto` should be followed by its target label whenever possible
- A trace is a path in the CFG — it characterizes some fragment of a real program execution.
- Algorithm for gathering traces: just do the depth-first traversal of the CFG — (can also take advantage of branch prediction information)

Traces => List of Statements

- Flatten the traces back to an ordered list of statements (canonical trees):
 - any `goto` followed by its false label: do nothing;
 - any `goto` followed by its true label: switch its true and false label, and `negate` the condition;
 - remove `goto` (1) if it is followed by its target 1;
 - any `goto` (cond, L_1 , L_2) followed by neither label: invent a new false label L_3 , rewrite it into:

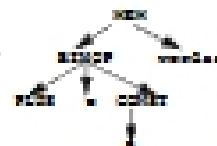

```
goto (cond,  $L_1$ ,  $L_3$ )
label  $L_3$ 
goto (cond,  $L_2$ )
```
- We are now ready to do instruction selection: generate assembly code for your favourite target machine.

Instruction Selection

- Input: an ordered list of canonical trees;
- Output: the pseudo-assembly code (without register assignments)
- Algorithm: translating each canonical tree into an assembly code sequence, and then concatenate all sequences together.
 - Main Problem: how to map the canonical tree to the assembly code?
- Each machine instruction can be expressed as a tree pattern — a fragment of the canonical tree:

Load the value at addr $a + c$ in the memory!

Each machine instruction may correspond to several layer of tree expressions



Instruction Selection via Tiling

- Express each machine instruction as a tree pattern.
- Given a canonical tree, the instruction selection is just to tile the tree using various tree patterns (for all possible machine instructions) — cover the canonical tree using nonoverlapping tiles.
- Optimum Tiling: — one whose tiles sum to the lowest possible value (suppose we give each machine instruction a cost)
- Optimal Tiling: — one where no two adjacent tiles can be combined into a single tile of lower cost
- Even optimum tiling is also optimal, but no vice versa.
- Algorithm: maximum munch finds the optimal tiling; dynamic programming finds the optimum tiling.