

CS 412/413

Introduction to Compilers and Translators Cornell University Andrew Myers

Lecture 14: Canonical Intermediate Code
23 Feb 00

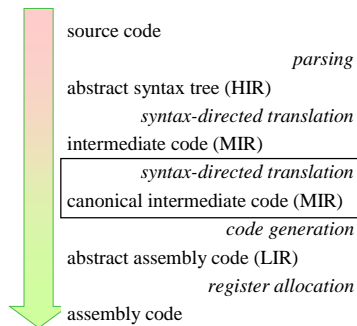
Administration

- HW3 due Friday
- Prelim 1 in one week
 - covers topics up through this lecture
 - lexical, syntactic analysis
 - type checking and static semantics
 - syntax-directed translation and IR
- Prelim review Monday evening 7-9 PM
 - location TBA

CS 412/413 Spring '00 Lecture 14 Andrew Myers

2

Where we are

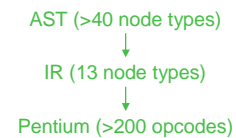


CS 412/413 Spring '00 Lecture 14 Andrew Myers

3

What makes a good IR?

- Easy to translate from AST
- Easy to translate to assembly
- Narrow interface: small number of node types (instructions)
 - Easy to optimize
 - Easy to retarget



CS 412/413 Spring '00 Lecture 14 Andrew Myers

4

Canonical form

- Intermediate code has general tree form
 - easy to generate from AST, but...
- Hard to translate directly to assembly
 - assembly code is a sequence of statements
 - Intermediate code (IR) has nodes corresponding to assembly statements deep in expression trees
- Canonical form: all statements brought up to top level of tree — generate assembly directly

CS 412/413 Spring '00 Lecture 14 Andrew Myers

5

One SEQ node

- In canonical form, only one SEQ node at the very top of tree



- Function is just a list of statements:
SEQ($S_1, S_2, S_3, S_4, S_5, \dots$)
- Can translate to assembly by translating each S_i to assembly statement(s) and concatenating

CS 412/413 Spring '00 Lecture 14 Andrew Myers

6

Canonical form

- Idea: rewrite IR to get rid of constructs incompatible with assembly code
 - arbitrarily deep expression trees — deal with this later as part of *instruction tiling*
 - ESEQ & CALL nodes — rewrite tree so no ESEQ nodes, CALLs moved to top
 - CJUMP is two-way jump rather than fall-through — convert to one-way jumps

no ESEQ nodes

- ESEQ nodes put a statement node underneath an expression:



$S \llbracket x = a[(i = i+1)] \rrbracket = ?$

- Problem: statement can have arbitrary number of side effects; assembly can't

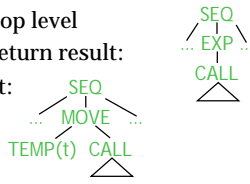
- Canonical form: **no ESEQ nodes**
similar to: $x = a[(i = i+1)] \Rightarrow i=i+1; x=a[i];$

Top-level CALL statements

- CALL nodes have arbitrary side effects
- CALL node deep in expression tree will break translation to assembly

Example: $x = f(g(x) + h(y))$

- Solution: move to top level
- Call that discards return result:
- Call that uses result:



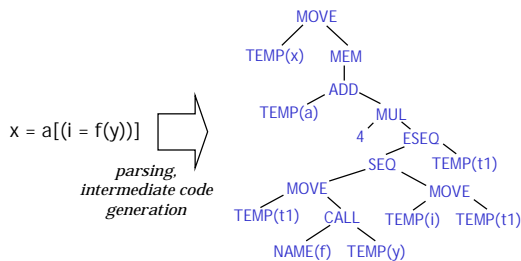
Canonical form

- Canonical tree has top-level SEQ node with following kinds of children:

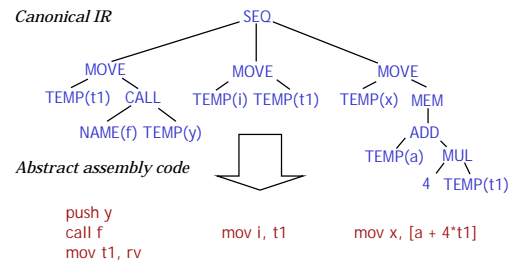
$MOVE(dest, e)$
 $MOVE(TEMP(t), CALL(...))$
 $EXP(CALL(...))$
 $JUMP(e)$
 $CJUMP(e, I_1, I_2)$
 $LABEL(l)$

\Rightarrow Straightforward translation to assembly

Example



Canonical IR to assembly



ESEQ rewriting

- Want to move ESEQ nodes up to top of tree where they can become SEQ nodes
- Idea: define syntax-directed rules that take an IR tree and move ESEQ nodes to top.
- **Goal:** avoid ripping apart expressions more than necessary -- leads to better code because expression patterns can be recognized and mapped to instruction set

CS 412/413 Spring '00 Lecture 14 Andrew Myers

13

ESEQ rewrite rules

- **Example transformations:**

$$\text{ESEQ}(s1, \text{ESEQ}(s2, e)) \Rightarrow \text{ESEQ}(\text{SEQ}(s1, s2), e)$$

$$\text{MOVE}(\text{ESEQ}(s1, e), \text{dest}) \Rightarrow \text{SEQ}(s1, \text{MOVE}(e, \text{dest}))$$

$$\text{OP}(\text{ESEQ}(s1, e1), e2) \Rightarrow \text{ESEQ}(s1, \text{OP}(e1, e2))$$

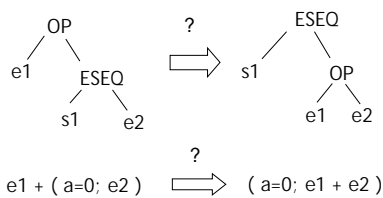
$$\text{OP}(e1, \text{ESEQ}(s1, e2)) \Rightarrow ?$$

CS 412/413 Spring '00 Lecture 14 Andrew Myers

14

Rewriting expressions

- $\text{OP}(e1, \text{ESEQ}(s1, e2))$

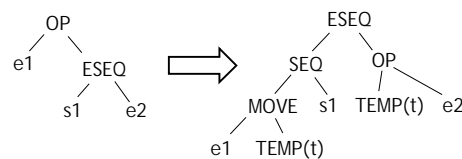


CS 412/413 Spring '00 Lecture 14 Andrew Myers

15

Introducing temporaries

- If $e1$ does not *commute* with $s1$
 – i.e., $(s1; e1; e2) \neq (e1; s1; e2)$
 must save value of $e1$ in temporary



CS 412/413 Spring '00 Lecture 14 Andrew Myers

16

Implementation options

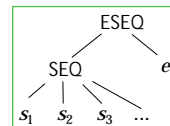
- Option 1: Walk over tree looking for places to apply rewrite rules (can use visitors)
 - “bad” nodes (ESEQ, CALL) percolate upward
 - Problem: need to restart tree traversal at every rewrite
- Option 2: Rewrite whole IR tree in one pass to build canonical IR tree
 - Syntax-directed translation!

CS 412/413 Spring '00 Lecture 14 Andrew Myers

17

General case

- When we move all ESEQ nodes to top, arbitrary expression node looks like:
- ESEQ transformation takes arbitrary expression node, returns list of sub-statements s_j to be executed (all side effects of e) plus final expression e **free of side effects**.
- Arbitrary statement node becomes a new SEQ node with no ESEQ nodes (or list of sub-statements s_j)



CS 412/413 Spring '00 Lecture 14 Andrew Myers

18

IR Simplification Interface

```
class CanonExpr {
  IRStmt[ ] pre_stmts;
  IRExpr expr;
}

abstract class IRExpr { CanonExpr simplify(); }
abstract class IRStmt { IRStmt[ ] simplify(); }
```

CS 412/413 Spring '00 Lecture 14 Andrew Myers

19

Simplification

Two translation functions:

- $T \llbracket e \rrbracket$ gives a list of canonical statements s_i and a canonical expression e' where executing the s_i and **then** evaluating e' has same side effects and value as e (= IRExpr.simplify)

$$T \llbracket e \rrbracket = (s_1, \dots, s_n) ; e'$$
- $T \llbracket s \rrbracket$ gives a list of canonical statements s_i such that executing the s_i has same side effects as s (= IRStmt.simplify)

$$T \llbracket s \rrbracket = (s_1, \dots, s_n)$$

CS 412/413 Spring '00 Lecture 14 Andrew Myers

20

Simplifying a function body

- Last time: translate a function definition $f(a_1, \dots, a_n) = e$ as $SEQ(\text{MOVE}(\text{TEMP}(RV), E \llbracket e \rrbracket), \text{LABEL}(\text{epilogue}))$
- Canonical form: SEQ node with all of $T \llbracket SEQ(\text{MOVE}(\text{TEMP}(RV), E \llbracket e \rrbracket), \text{LABEL}(\text{epilogue})) \rrbracket$ as children.



CS 412/413 Spring '00 Lecture 14 Andrew Myers

21

Rules

- Simplify arbitrary expression e :

$$T \llbracket e \rrbracket = (s_1, s_2, s_3, \dots) ; e'$$
- Goal: define $T \llbracket e \rrbracket$ and $T \llbracket s \rrbracket$ for all 13 node types
- 3 trival cases:

$$T \llbracket \text{CONST}(t) \rrbracket = () ; \text{CONST}(t)$$

$$T \llbracket \text{NAME}(n) \rrbracket = () ; \text{NAME}(n)$$

$$T \llbracket \text{TEMP}(t) \rrbracket = () ; \text{TEMP}(t)$$
- These expressions are already free of side effects: already in canonical form

CS 412/413 Spring '00 Lecture 14 Andrew Myers

22

JUMP, CJUMP, MEM

- $JUMP(e)$, $CJUMP(e, I_1, I_2)$, $MEM(e)$
- Need to make sure e is canonical
- $T \llbracket JUMP(e) \rrbracket = (s_1, \dots, s_n, JUMP(e'))$ if $T \llbracket e \rrbracket = (s_1, \dots, s_n) ; e'$
- Similarly for CJUMP
- Can write as inference rule:

$$\frac{T \llbracket e \rrbracket = (s_1, \dots, s_n) ; e'}{T \llbracket MEM(e) \rrbracket = (s_1, \dots, s_n) ; MEM(e')}$$

CS 412/413 Spring '00 Lecture 14 Andrew Myers

23

ESEQ

- How to simplify an expression $ESEQ(s, e)$?

$$\frac{T \llbracket e \rrbracket = (s_1, \dots, s_n) ; e'}{T \llbracket ESEQ(s, e) \rrbracket = (s, s_1, \dots, s_n) ; e'}$$
- What is wrong with this rule?

CS 412/413 Spring '00 Lecture 14 Andrew Myers

24

Correct ESEQ rule

$$\frac{\begin{array}{l} \mathbb{T} \llbracket e \rrbracket = (s_1, \dots, s_n) ; e' \\ \mathbb{T} \llbracket s \rrbracket = (s'_1, \dots, s'_m) \end{array}}{\mathbb{T} \llbracket \text{ESEQ}(s, e) \rrbracket = (s'_1, \dots, s'_m, s_1, \dots, s_n) ; e'}$$

- Assuming $\mathbb{T} \llbracket e \rrbracket, \mathbb{T} \llbracket s \rrbracket$ produce canonical statements s_i, s'_j and canonical expression e , $\mathbb{T} \llbracket \text{ESEQ}(s, e) \rrbracket$ works properly.

Translating binary operators

$$\frac{\begin{array}{l} \mathbb{T} \llbracket e_1 \rrbracket = (s_1, \dots, s_n) ; e'_1 \\ \mathbb{T} \llbracket e_2 \rrbracket = (s'_1, \dots, s'_n) ; e'_2 \end{array}}{\mathbb{T} \llbracket \text{OP}(e_1, e_2) \rrbracket = (s_1, \dots, s_m, \text{MOVE}(\text{TEMP}(t), e'_1), s'_1, \dots, s'_n) ; \text{OP}(\text{TEMP}(t), e'_2)}$$

- Idea: save value of e_1 in a temporary before executing all the side effects of e_2

SEQ nodes

- How to get rid of SEQ nodes: concatenate canonical versions of all sub-statements

$$\frac{\begin{array}{l} \mathbb{T} \llbracket s_1 \rrbracket = (s'_1, \dots, s'_m) \\ \mathbb{T} \llbracket s_2 \rrbracket = (s''_1, \dots, s''_n) \end{array}}{\mathbb{T} \llbracket \text{SEQ}(s_1, s_2) \rrbracket = (s'_1, \dots, s'_m, s''_1, \dots, s''_n)}$$

CALL nodes

- CALL nodes call a function which may have side effects
- Overwrites return value register at least; can't be operand at assembly-code level
- Therefore, CALL nodes must move to top

$$\frac{\begin{array}{l} \mathbb{T} \llbracket e_f \rrbracket = (s_1, \dots, s_m) ; e'_f \\ \mathbb{T} \llbracket e_1 \rrbracket = (s'_1, \dots, s'_n) ; e'_1 \end{array}}{\mathbb{T} \llbracket \text{CALL}(e_f, e_1) \rrbracket = (s_1, \dots, s_m, s'_1, \dots, s'_n, \text{MOVE}(\text{TEMP}(t), \text{CALL}(e'_f, e'_1)); \text{TEMP}(t))}$$

EXP

- $\text{EXP}(e)$ evaluates e for its side effects, discards value
- Simplified IR does same:

$$\frac{\mathbb{T} \llbracket e \rrbracket = (s_1, \dots, s_n) ; e'}{\mathbb{T} \llbracket \text{EXP}(e) \rrbracket = (s_1, \dots, s_n)}$$

Canonical intermediate code

- Syntax-directed translation function $\mathbb{T} \llbracket \cdot \rrbracket$ simplifies an IR tree into canonical form
- Yields recursive impl. of $\text{IRstmt.simplify}, \text{IRExpr.simplify}$
- Canonical form: IR is a sequence of simple IR statements ready for translation to assembly