



CS 412 Introduction to Compilers

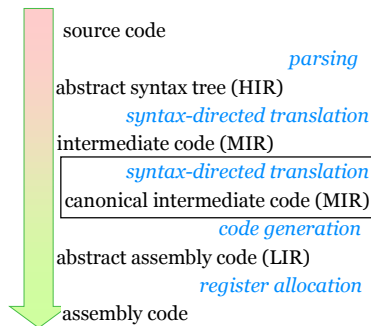
Andrew Myers
Cornell University

Lecture 15: Canonical IR
26 Feb 01

Administration

- HW3 due Friday
- Prelim 1 next Tuesday evening (7:30-9:30PM)
 - location TBA
 - covers topics up through this lecture
 - lexical, syntactic analysis
 - type checking and static semantics
 - syntax-directed translation and IR

Where we are



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

One SEQ node

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



- Function body 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

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: $x = a[(i = i+1)] \Rightarrow i=i+1; x=a[i];$

Lecture 15 CS 412/413 Spring '01 -- Andrew Myers

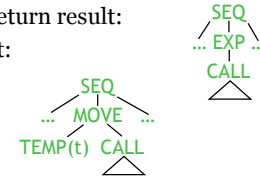
7

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:



Lecture 15 CS 412/413 Spring '01 -- Andrew Myers

8

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, l₁, l₂)
 LABEL(l)

Code is a just sequence of these statements
 \Rightarrow **Straightforward translation to assembly**

Lecture 15 CS 412/413 Spring '01 -- Andrew Myers

9

Simplifying a function body

- Last time: translate a function definition $f(a_1, \dots, a_n) = e$

as $SEQ(\text{MOVE}(\text{TEMP}(\text{RV}), \mathcal{L}\llbracket e \rrbracket), \text{LABEL}(\text{epilogue}))$

- Canonical form: SEQ node with all of $\llbracket SEQ(\text{MOVE}(\text{TEMP}(\text{RV}), \mathcal{L}\llbracket e \rrbracket), \text{LABEL}(\text{epilogue})) \rrbracket$ as children.



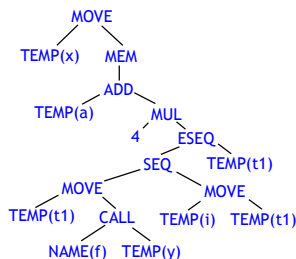
Lecture 15 CS 412/413 Spring '01 -- Andrew Myers

10

Example

$x = a[(i = f(y))]$

parsing,
intermediate code
generation

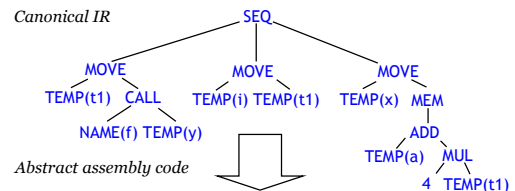


Lecture 15 CS 412/413 Spring '01 -- Andrew Myers

11

Canonical MIR to LIR

Canonical IR



Abstract assembly code

push y
call f
mov t1, rv

mov i, t1

mov x, [a + 4*t1]

Lecture 15 CS 412/413 Spring '01 -- Andrew Myers

12

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

Lecture 15 CS 412/413 Spring '01 -- Andrew Myers

13

ESEQ rewrite rules

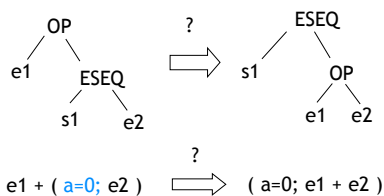
- Example transformations:
 - $ESEQ(s1, ESEQ(s2, e)) \Rightarrow ESEQ(SEQ(s1, s2), e)$
 - $MOVE(ESEQ(s1, e), dest) \Rightarrow SEQ(s1, MOVE(e, dest))$
 - $OP(ESEQ(s1, e1), e2) \Rightarrow ESEQ(s1, OP(e1, e2))$
 - $OP(e1, ESEQ(s1, e2)) \Rightarrow ?$

Lecture 15 CS 412/413 Spring '01 -- Andrew Myers

14

Rewriting expressions

- $OP(e1, ESEQ(s1, e2))$



Lecture 15 CS 412/413 Spring '01 -- Andrew Myers

15

Implementation options

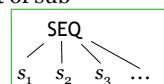
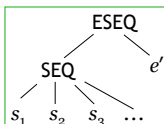
- Option 1: Walk over tree looking for places to apply rewrite rules
 - “bad” nodes (ESEQ, CALL) percolate upward, eventually disappear
 - Problem: may need to restart tree traversal at every rewrite
- Option 2: Rewrite whole IR tree in one pass to build canonical IR tree
 - Syntax-directed translation!

Lecture 15 CS 412/413 Spring '01 -- Andrew Myers

16

General case

- When we move all ESEQ nodes to top, arbitrary expression node e looks like:
- Transformation returns list of sub-statements s_i plus final expression e'
 - each s_i has at most one side-effect
 - e' is *free of side effects*.
- Arbitrary statement node becomes a new SEQ node with no ESEQ nodes (or list of sub-statements s_i)
 - each s_i has at most one side-effect



Lecture 15 CS 412/413 Spring '01 -- Andrew Myers

17

IR simplification Interface

```
class CanonicalExpr {
    IRStmt[ ] pre_stmts;
    IExpr expr; }
}
class CanonicalStmt {
    IRStmt[ ] stmts; }
}
abstract class IExpr { CanonicalExpr simplify(); }
abstract class IRStmt { IRStmt[ ] simplify( ); }
```

Lecture 15 CS 412/413 Spring '01 -- Andrew Myers

18

Simplification

Two translation functions:

- $\mathcal{T}[e]$ 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)

$$\mathcal{T}[e] = (s_1, \dots, s_n); e'$$

- $\mathcal{T}[s]$ gives a list of canonical statements s_i such that executing the s_i has same side effects as s (= IRStmt.simplify)

$$\mathcal{T}[s] = (s_1, \dots, s_n)$$

Rules

- Simplify arbitrary expression e :

$$\mathcal{T}[e] = (s_1, s_2, s_3, \dots); e'$$

- Goal: define $\mathcal{T}[e]$ and $\mathcal{T}[s]$ for all 13 node types

- 3 trivial cases:

$$\mathcal{T}[\text{CONST}(i)] = (); \text{CONST}(i)$$

$$\mathcal{T}[\text{NAME}(n)] = (); \text{NAME}(n)$$

$$\mathcal{T}[\text{TEMP}(t)] = (); \text{TEMP}(t)$$

- Already in canonical form!

JUMP, CJUMP, MEM

- JUMP(e), CJUMP(e, l_1, l_2), MEM(e)
– Need to make sure e is canonical

- $\mathcal{T}[\text{JUMP}(e)] = (s_1, \dots, s_n, \text{JUMP}(e'))$
if $\mathcal{T}[e] = (s_1, \dots, s_n); e'$

- Similarly for CJUMP

- Can write as inference rule:

$$\frac{\mathcal{T}[e] = (s_1, \dots, s_n); e'}{\mathcal{T}[\text{MEM}(e)] = (s_1, \dots, s_n); \text{MEM}(e')}$$

ESEQ

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

$$\frac{\mathcal{T}[e] = (s_1, \dots, s_n); e'}{\mathcal{T}[\text{ESEQ}(s, e)] = (s, s_1, \dots, s_n); e'}$$

?

Correct ESEQ rule

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

Assuming $\mathcal{T}[e], \mathcal{T}[s]$ produce canonical statements s_i, s'_j and canonical expression e , $\mathcal{T}[\text{ESEQ}(s, e)]$ works properly.

SEQ nodes

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

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

EXP

- EXP(e) evaluates e for its side effects, discards value
- Simplified IR does same:

$$\begin{aligned} \mathcal{T}[e] &= (s_1, \dots, s_n); e' \\ \mathcal{T}[\text{EXP}(e)] &= (s_1, \dots, s_n) \end{aligned}$$

Translating binary operators

$$\begin{aligned} \mathcal{T}[e_1] &= (s_1, \dots, s_m); e'_1 \\ \mathcal{T}[e_2] &= (s'_1, \dots, s'_n); e'_2 \\ \hline \mathcal{T}[\text{OP}(e_1, e_2)] &= (s_1, \dots, s_m, s'_1, \dots, s'_n); \text{OP}(e'_1, e'_2) \end{aligned}$$

- When does this rule work?
 - Note: OP allows either e_1 or e_2 to be evaluated first

Translating binary operators

- Previous rule works if e'_1 commutes with each of s'_i or e'_2 commutes with each of s_i
- Idea: save value of e_1 in a temporary before executing all the side effects of e_2

$$\begin{aligned} \mathcal{T}[e_1] &= (s_1, \dots, s_m); e'_1 \\ \mathcal{T}[e_2] &= (s'_1, \dots, s'_n); e'_2 \\ \hline \mathcal{T}[\text{OP}(e_1, e_2)] &= \\ & (s_1, \dots, s_m, \text{MOVE}(\text{TEMP}(t), e'_1), s'_1, \dots, s'_n); \text{OP}(\text{TEMP}(t), e'_2) \end{aligned}$$

- Works, but:
 - Introduces extra register t
 - No opportunity to do e'_1 and e'_2 in one instruction

Reordering

Statement s and expression e commute if executing s does not change result of e .
False if:

- s overwrites any TEMP used by e
- s overwrites any MEM location that might be the same location as (alias) a MEM in e
 - conservative assumption: all memory locations may alias one another
 - less conservative: use *alias analysis* to determine which memory locations may alias

CALL nodes

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

$$\begin{aligned} \mathcal{T}[e_f] &= (s_1, \dots, s_m); e'_f \\ \mathcal{T}[e_1] &= (s'_1, \dots, s'_n); e'_1 \\ \hline \mathcal{T}[\text{CALL}(e_f, e_1)] &= \\ & (s_1, \dots, s_m, s'_1, \dots, s'_n, \text{MOVE}(\text{TEMP}(t), \text{CALL}(e'_f, e'_1))); \text{TEMP}(t) \end{aligned}$$

- Assumes e'_f commutes with s'_1, \dots, s'_n
- Iota, C: operands may be reordered freely

Canonical intermediate code

- Syntax-directed translation function $\mathcal{T}[\cdot]$ simplifies an IR tree into canonical form
- Yields recursive implementation of `IRStmt.simplify`, `IRExpr.simplify`
- Canonical form: IR is a sequence of simple IR statements ready for translation to assembly