

Compilers - Lexical & Syntax Analysis (Producing an AST) → Semantic Analysis (Producing a Symbol Table / Typed AST) → Runtime Memory Organization → IR Generation (Making an IR of assembly code. Optimizations) → Assembly Generation (writing the assembly string based on our IR).

1) Compilers - Lexical Analysis

Transforming a stream of characters into tokens, based on the formal description of tokens.

Lexical Analysis made of Regular Expressions.

1.1 Identifier Tokens (Part of Lexical Analysis):

LA uses fast string lookup with hash function.

1) **Keyword Identifiers:** Special reserved words in language e.g. 'class', 'while' → WHILE, etc. 2) **Non-Keyword Identifiers:** Programmer defined identifiers, such as variables. A generic token used that uses the provided string name. 'var1' → IDENT('var1')

Literals: Tokens are constant values embedded in the program (INTEGER(13), FLOAT(17.03), STRING("hi"))

Other tokens: Operators, whitespace, comments, pre-processing directives and Macros (think back to WACC).

1.2 Regular Expressions Rules

Can't be recursive. Regexes match strings.

Rule → Description - Example:

1) a - matches a symbol - x matches 'x' only.

2) \symbol - Escapes a regex character. \r matches '\r' only

3) e - matches the empty string. e - matches "" only

4) R₁R₂ - matches adjacent as a combined rule - ab89 matches 'ab89' only

5) R₁|R₂ - alternation, match one regex or the other. abc1 matches 'abc' or '1'

6) (R) - group regexes together, (a|b)c matches 'ac' and 'bc'

7) R+ - matches one or more repetitions of R, (a|b|c)+ matches all non-empty strings of a,b & c (eg: 'abbba')

8) R* - Zero or more repetitions. R* equivalent to (R+)?

Precedence from highest to lowest: grouping, repetition, concatenation, alternation. We can derive these compound rules which are much more useful:

Rule → Description - Example:

1) R? - Zero or One occurrence - a? matches "a" and "a".

2) ~ wildcard, matches any possible string.

3) [abcd] - Character set, match any chars in set.

4) [0-9] or [a-zA-Z] - match a single number from 0 to 9, or any alphabet character.

5) \"'abc\" - match any character except those in the set

We can use these in the production rules:

SignedInt → (+|-)? Int

Keyword → 'if' | 'while' | 'do'

When one or more expression matches, we choose the longest matching character sequence. Else, regex rules are ordered, with earlier rules taking precedence.

Lexical Analysers using Regex like this are quite easy to write, but it becomes hard to change token rules. In practice we use **Lexical Analyser Generators**, which take programmer defined token structures and functions based on the formal language definition to produce a **tokenizer** (program input → AST).

To generate a Lexical Analyser, we first write our **Regular Expressions** and then use **Thompson's Construction** to convert to a **Non-deterministic Finite Automata**. Then use **Subset Construction** to convert to a **Deterministic Finite Automata**. We can do further optimizations to get a **Minimum State DFA** and from here we can get a **Transition Table** + **GetToken function** - by mapping our graph (our FSM) into a table, which is our **Lexical Analyser**.

1.3) Finite Automata (Finite State Machines)

Arrows denote transitions between states.

Start state has an unlabelled transition to it.

Accepting states are double circles.

Will stop when no transition can be made (as a result matches the longest string possible to a state).

1.3.1) M1) Thompson's Construction Diagram:

(R)

R+

R1|R2

R1R2

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

2.4) LR(1) Parsers

→ Zero or more derivations.

→ One or more derivations.

A → Non-terminal A is nullable.

A → AB|BC|AB and BC are alternatives of A.

M3) Computing the First Set:

1) Start at the NFA start state.

2) Compute the ε-Closure for each state. These are the set of all states that can be reached from the current state by any number of ε-transitions. Set the new name of each node to be its ε-Closure.

3) Now, from our grouped node, consider all A transitions. Take the epsilon closure of all nodes we can go to with A, as a node. Do this for B and so on.

4) For all the new nodes we constructed, carry out step 2 and 3. Note: Be very careful when adding transitions - remember we can loop to ourselves if we had a connection from one node in our epsilon closure to another! Also, any node that contains an accepting state in its epsilon closure is now an accepting state in our DFA. This accepting state can still have transitions.

Context Free Grammars are transformed to NFAs. CFGs take the form rule → (rule|token)*

2.5) Parsing - LR Parsing

2.1) Chomsky Hierarchy Given:

R - non-terminal

t - sequence of tokens

a, b, c - sequences of terminals and non-terminals

α, β, γ - context-free

c - context-sensitive

rec - enumerable

2.2) Parser for Context Free Grammars

Context-free grammars parsed with complexity O(n³).

LL and LR grammars subsets of CFG, parsed in O(n), which is why we use them. LL(n) is a subset of LR(n). Must be generated. LR Parsers are bottom up, building the AST from the leaves to the root. LL Parser are top down, building the AST from the root to leaves. LL/LR(k) denotes an LL/LR parser with k token lookahead.

2.3) LR(0) Parsers, 2.3.1) LR(0) Items

LR(0) parsers don't use the current token to perform a reduction. • represents the current position of the parser.

Initial Item - we haven't parsed any part of a rule yet.

Complete (Reduce) Item - we've fully matched a rule and can be reduced.

The rule X → AB has 3 LR(0) items, the initial item •AB, A•B, and the reduce item AB•. We add a **start rule**: E' → E \$ - where \$ is end of input. If omitted, its implied.

2.3.2) NFA from LR(0) Items - we can create transitions between LR(0) Items 1) Given an item X → A•BC, X → A•B•C if we have a rule for B.

2) If B is non terminal (no rule matches), for each rule B → D, we add a new ε transition (we can't reduce yet). (X → A•BC) → (X → A•B•D) This means if we're about to process a block that's a rule in a DFA state, we have to add the productions of this rule to our DFA state. To convert to the DFA from here, we do subset construction, but it's simple enough right away.

2.3.3) DFA to LR(0) Item Table

The parsing table describes the rules to apply, and states to move to when a given item is encountered. It is generated from a DFA using the rules:

1) For each Terminal translation X → Y

Add P(X, T) = SY (shift Y, we cannot reduce yet)

2) For each Non-terminal translation X → Y

X → Y• Add P(X, N) = gY (Goto Y)

3) For each State X containing item R' → ... •

Add P(X, \$) = a (accept) (end of input)

4) For each State X containing item R → ... •

Add P(X, T) = n (reduce) (use rule n to reduce)

An empty cell indicates an error.

We have a table of our states as rows, and each possible input as our columns (including \$, and Goto for movements of S / whatever grammar rule we're parsing).

2.3.4) LR Parser Algorithm

For action, we write T [state, lookahead] - what we do. We start with 0 (state 0) on the stack. \$ ends our tokens

Accepting states are double circles.

Will stop when no transition can be made (as a result matches the longest string possible to a state).

1.3.1) M1) Thompson's Construction Diagram:

(R)

R+

R1|R2

R1R2

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

2.4) LR(1) Parsers

→ Zero or more derivations.

→ One or more derivations.

A → Non-terminal A is nullable.

A → AB|BC|AB and BC are alternatives of A.

M3) Computing the First Set:

1) Start at the NFA start state.

2) Compute the ε-Closure for each state. These are the set of all states that can be reached from the current state by any number of ε-transitions. Set the new name of each node to be its ε-Closure.

3) Now, from our grouped node, consider all A transitions. Take the epsilon closure of all nodes we can go to with A, as a node. Do this for B and so on.

4) For all the new nodes we constructed, carry out step 2 and 3. Note: Be very careful when adding transitions - remember we can loop to ourselves if we had a connection from one node in our epsilon closure to another! Also, any node that contains an accepting state in its epsilon closure is now an accepting state in our DFA. This accepting state can still have transitions.

Context Free Grammars are transformed to NFAs. CFGs take the form rule → (rule|token)*

2.5) Parsing - LR Parsing

2.1) Chomsky Hierarchy Given:

R - non-terminal

t - sequence of tokens

a, b, c - sequences of terminals and non-terminals

α, β, γ - context-free

c - context-sensitive

rec - enumerable

2.2) Parser for Context Free Grammars

Context-free grammars parsed with complexity O(n³).

LL and LR grammars subsets of CFG, parsed in O(n), which is why we use them. LL(n) is a subset of LR(n). Must be generated. LR Parsers are bottom up, building the AST from the leaves to the root. LL Parser are top down, building the AST from the root to leaves. LL/LR(k) denotes an LL/LR parser with k token lookahead.

2.3) LR(0) Parsers, 2.3.1) LR(0) Items

LR(0) parsers don't use the current token to perform a reduction. • represents the current position of the parser.

Initial Item - we haven't parsed any part of a rule yet.

Complete (Reduce) Item - we've fully matched a rule and can be reduced.

The rule X → AB has 3 LR(0) items, the initial item •AB, A•B, and the reduce item AB•. We add a **start rule**: E' → E \$ - where \$ is end of input. If omitted, its implied.

2.3.2) NFA from LR(0) Items - we can create transitions between LR(0) Items 1) Given an item X → A•BC, X → A•B•C if we have a rule for B.

2) If B is non terminal (no rule matches), for each rule B → D, we add a new ε transition (we can't reduce yet). (X → A•BC) → (X → A•B•D) This means if we're about to process a block that's a rule in a DFA state, we have to add the productions of this rule to our DFA state. To convert to the DFA from here, we do subset construction, but it's simple enough right away.

2.3.3) DFA to LR(0) Item Table

The parsing table describes the rules to apply, and states to move to when a given item is encountered. It is generated from a DFA using the rules:

1) For each Terminal translation X → Y

Add P(X, T) = SY (shift Y, we cannot reduce yet)

2) For each Non-terminal translation X → Y

X → Y• Add P(X, N) = gY (Goto Y)

3) For each State X containing item R' → ... •

Add P(X, \$) = a (accept) (end of input)

4) For each State X containing item R → ... •

Add P(X, T) = n (reduce) (use rule n to reduce)

An empty cell indicates an error.

We have a table of our states as rows, and each possible input as our columns (including \$, and Goto for movements of S / whatever grammar rule we're parsing).

2.3.4) LR Parser Algorithm

For action, we write T [state, lookahead] - what we do. We start with 0 (state 0) on the stack. \$ ends our tokens

Accepting states are double circles.

Will stop when no transition can be made (as a result matches the longest string possible to a state).

1.3.1) M1) Thompson's Construction Diagram:

(R)

R+

R1|R2

R1R2

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

(R)*

2.4) LR(1) Parsers

→ Zero or more derivations.

→ One or more derivations.

A → Non-terminal A is nullable.

A → AB|BC|AB and BC are alternatives of A.

M3) Computing the First Set:

1) Start at the NFA start state.

2) Compute the ε-Closure for each state. These are the set of all states that can be reached from the current state by any number of ε-transitions. Set the new name of each node to be its ε-Closure.

3) Now, from our grouped node, consider all A transitions. Take the epsilon closure of all nodes we can go to with A, as a node. Do this for B and so on.

4) For all the new nodes we constructed, carry out step 2 and 3. Note: Be very careful when adding transitions - remember we can loop to ourselves if we had a connection from one node in our epsilon closure to another! Also, any node that contains an accepting state in its epsilon closure is now an accepting state in our DFA. This accepting state can still have transitions.

Context Free Grammars are transformed to NFAs. CFGs take the form rule → (rule|token)*

2.5) Parsing - LR Parsing

2.1) Chomsky Hierarchy Given:

R - non-terminal

t - sequence of tokens

a, b, c - sequences of terminals and non-terminals

α, β, γ - context-free

c - context-sensitive

Compilers - Backend

5.1.1) Backus Naur Form

A CFG specifying the syntactic structure of a language
A CFG is a set of **Productions**, associated with a set of **tokens** (terminals), non-terminals (**rule**) & start symbol
Each production is of the form:

non-terminal \rightarrow String of terminals & non-terminals

1) Productions: A way to expand a non-terminal symbol into a string of terminals & non-terminals

2) Terminals: Symbols that can't be further expanded (tokens generated from Lex. Analysis)

3) Non-Terminal: Symbols that can be expanded further - outlined in a Production.

3) Parse Trees show how a string is derived from the start symbol.

5.1.1.1) Associativity Associativity can be enforced by using left or right recursive productions:

term \rightarrow const | ident
expr \rightarrow expr - term
expr \rightarrow term - expr
expr \rightarrow term

5.1.2.2) Precedence

To enforce precedence, we can consider **levels**. We factor those of highest precedence to the lowest level.

expr \rightarrow exp + term | exp - term | term
term \rightarrow exp * factor | term / factor | factor
factor \rightarrow const | ident

To prove something is a member of our grammar we can construct the derivation or a Parse Tree.

5.1.3.3) Production Choice

We may have a grammar where we cannot determine which production for a non-terminal token to use based on the first symbol.

stat \rightarrow 'loop' statlist 'until' expr
stat \rightarrow 'loop' statlist 'while' expr
stat \rightarrow 'loop' statlist 'forever'

When we have taken 'loop' we cannot determine which production to use. Methods to deal with this:

1) Delay the choice

Delay creating this tree (from stat) until it is known which production matches. It is still possible to create the statlist inside while doing so.

2) Modify the grammar

Change the grammar to factor out the difference.

stat \rightarrow loop' statlist loopstat
loopstat \rightarrow 'until' expr
loopstat \rightarrow 'while' expr
loopstat \rightarrow 'forever'

However there are more difficult problems, which can be more easily fixed with bottom-up parsing.

Top down parsing is done by Rec. Desc. Parsers, which can deal with left recursion.

Bottom-up Parsing

The grammar's productions are used right \rightarrow left.

Input is compared against the right hand side to produce a non-terminal on the left.

Parsing is complete when the whole input is replaced by the start symbol. Bottom up parsers are difficult to implement, so parser generators are recommended.

5.2) Visitor Pattern

The visitor pattern is a design pattern that is commonly used in object-oriented programming to separate algorithms from the objects they operate on. The idea behind the pattern is to create a separate class (the visitor) that can traverse a complex object structure and perform operations on its elements.

In the context of compilers, the visitor pattern can be used to implement the different phases of a compiler as separate visitor classes. For example, a lexer could be implemented as a visitor that traverses the source code and generates a stream of tokens, while a parser could be implemented as another visitor that takes this stream of tokens and generates an abstract syntax tree (AST).

6) Code Generation

1) The language:

data Stat = Assign Name Exp | Seq Stat Stat | ForLoop Name Exp Exp Stat

data Exp = BinOp Op Exp Exp | UnOp Op Exp | Ident Name | Const Int

data Op = Plus | Minus | Tim

name Name = [Char]

2) Assembly Instructions:

data Instruction = Add | Sub | Mul | Div | Push | Pop | Load Name | LoadImm Int | Store Name | Jump Label | JTrue Label | JFalse Label | Define Label

3) Assembly Pseudocode:

DD / MINUS / MUL / DIV:

T := store[SP]
SP := SP + 4
T := store[SP] [+*+//] T
store[SP] := T
PUSHIMM:
SP := SP - 4
store[SP] := operand (IR)
PUSHABS:
T := store[operand (IR)]
SP := SP - 4
store[SP] := T
POP:
T := store[SP]
SP := SP + 4
store[operand (IR)] := T
COMPEQ:
T := store[SP]
SP := SP + 4
T := store[SP] - T
store[SP] := T=0 ? 1 : 0
JTRUE / JFALSE:
T := store[SP]
SP := SP + 4
PC := T=1 / 0 ? operand (IR) : PC

4) Translate Functions for Exp and Stat

transExp::Exp \rightarrow [Instruction]
transExp (BinOp op e1 e2) =
transExp e1 ++ transExp e2 ++ [case op of
Plus \rightarrow Add
Minus \rightarrow Sub
Times \rightarrow Mul
Divide \rightarrow Div]
transExp (UnOp Minus e)
= transExp e ++ [PushImm (-1), Mul]
transExp (UnOp _) = error "(transExp)
Only '-' unary operator supported"
transExp (Ident id) = [PushAbs id]
transExp (Const n) = [PushImm n]

transStat::Stat \rightarrow [Instruction]
transStat (Assign id exp) = transExp exp ++
[Pop id]
transStat (Seq s1 s2) =
transStat s1 ++ transStat s2
transStat (ForLoop x e1 e2 body)
= transExp e1 ++ [Pop x] ++ [Define "loop"]
++ transExp e2 ++ [CompEq::JTrue "break"]
++ transStat body ++ [PushImm 1, Add, Pop x,
Jump "loop", Define "break"]

7) Improving our Assembly

1) Using Immediate Instructions

movl \$3, %eax
imull \$3, %eax
addl \$4, %eax
Rather than moving into registers first and then doing work.

We simply add new AddImm, SubImm, etc. data types. The translateFunctions for these are very simple.

2) Dealing with Bounded Numbers of Registers

1) Accumulator machines have 1 register. We store the accumulated value in there.

2) Register Allocation for Function Calls

Need to know where parameters are when passed. (f(x) + 1) + (1 + (a + y)) Which side of the + should be evaluated first depends on the context (e.g registers that need to be saved at the call site, and registers used by the callee, calling convention).

3.2.1) Infeasible Control Path Problem

Control Flow Graphs capture control flow inside functions and methods but not between them. We could have infeasible paths as follows: (a,b,f invalid, a,d,e,c invalid due to how return works)

Call site

Jump & Link F
<Next Instruction>
...

Jump & Link F
<Next Instruction>
...

return

Body of F

saveRegs unusedRs = [Mov (Reg x) Push | x <- unusedRs]
where usedRs = allRegs \ unusedRs

restoreRegs unusedRs = [Mov Pop (Reg x) | x <- revUsedRs]
where revUsedRs = reverse (allRegs \ unusedRs)

4.2.2) Caller and Callee Saving

We must enforce a **calling convention** to ensure registers are not **clobbered** (e.g non-argument or return registers are changed).

1) Caller Saved: save registers used by the caller in case the callee clobbers them.
(e.g: Save used registers by caller that callee also uses, call method, restore used registers by caller that callee also uses. **Problem:** We have to know what registers the callee uses.)

2) Callee Saved: Save registers that the callee uses only. (e.g: In the called method, save registers that the callee uses if they are also used by the caller - at the end of our method restore them. **Problem:** We have to know what registers the caller uses.)

7.2.3) IA 32 Calling Convention solves this

Given an expression E1 op E2 always evaluate the subexpression that uses most registers 1st %eax %ecx %edx %esi %edi %esp %ebp

7.3) Register Allocation by Graph Colouring

1) Use simple traversal to generate intermediate code Temporary values are always saved in a named location. (e.g t0...). This way we can consider all values including intermediate ones.

2) Construct an Inference Graph each node is a temporary location, each edge connects simultaneously live locations. Registers that need to simultaneously store values must be different colours (different registers).

3) Attempt To Colour Nodes: If colouring is not possible spilling occurs.

(a) Find an edge, to remove it either split the live range (e.g temporarily put to memory).

(b) Redo the analysis to determine if the graph can now be coloured.

When choosing which values to spill it is important to consider how often a variable is used, e.g avoid spilling from innermost loop. While NP hard, heuristics exist.

8) Optimization

High level optimizations use **high-level info** encoded in the program: (types, func analysis), e.g: Function Inlining.

Low level optimizations use **low-level info** (instruction types, the ISA, the order of instructions in the IR, etc) to optimise the output. e.g: Instruction Scheduling.

8.1.1) Peephole Optimization

Scan through the assembly in order, looking for obvious cases to optimise.

1) Can catch some of the worst cases (e.g store followed by load of the same location).

2) Very easy to implement (at smallest just consider two adjacent instructions).

3) The sparse ordering problem in what order should the optimisations be applied to get the best result?

8.2) Lowering Representation

Taking high-level features and converting them into lower-level representations. For example taking arrays and converting them into pointer arithmetic/address calculation. When lowering you lose high-level information (e.g that values are part of an array), but can optimise the lower level representation (optimise address calculations). We usually start with high level IRs, analyse, optimize, then move to lower IRs, optimizing based on the info we have on each level.

8.3) Other Optimizations

1) Induction Variable - a variable which increases / decreases by a (loop invariant) constant on each iteration.

2) Strength Reduction - an optimization where we calculate induction variables by breaking it down into a single addition rather than a compound expr which might have multiplication - which is expensive. (In general, replace a complex expression with a simpler one).

3) Control variable selection - replace loop control variable (as in the l in "for l in range") with an induction variable in our loops instead, and then rework the bounds check to work with the values of this induction variable (so we have less increments / variables).

4) Dead Code Elimination code that does not produce a use result can be eliminated. many other optimisations result in dead code (e.g inlining a function where not all the function's returned values or optional arguments are used.)

8.4) Data Flow Analysis for Live Ranges

Live Range - the range of instructions for which a temporary value must be maintained. A live range starts at a definition, and ends when either the variable is used, or immediately if the value is never used. Like with Graph Colouring we have a similar process:

1. Generate code using temporaries T0... instead of regs.

2. For each temporary T_i, find T_i's live range - set of instructions for which T_i must reside in a register.

3. If liveRange(T_i) intersects liveRange(T_j) then they must be allocated to different registers - **they interfere**.

4. Assemble the Register Inference Graph (RIG).

5. Colour the RIG. If successful replace temporaries with register and generate code. If Graph can't be recoloured, then find a temporary to spill, retry. Num colours = regs.

data CFG = ControlFlowGraph [CFGNode]

data CFGNode = Node ID Instruction [Register] [Id] type Id = Int

data Register = D Int | T Int

buildCFG :: [Instruction] -> CFG

List of Nodes in our graph. The nodes contain an ID, an instruction, the temporaries used by the instruction, and the ones defined by it, and the successors (the ids after the node - edges).

First, we build the CFG. Things to note

1) Succs refers to all the possible paths that can be taken from the current node - e.g

IR Code, CFG (line, instruction, uses, defs, succs):

Bra L2 1 Bra L2 [] [] [10]
L1:
cmp b a 2 cmp b a [b, a] [] [3]
bge L3 3 bge L3 [] [] [4, 8]
mul #7 a 4 mul #7 a [a] [a] [5]
mov a b 5 mov a b [a] [b] [6]
add #1 b 6 add #1 b [b] [b] [7]
bra L4 7 bra L4 [] [] [10]
L3:
mov b a 8 mov b a [b] [a] [9]
sub #1 a 9 sub #1 a [a] [a] [10]
L4:
L2:
Cmp b #10 10 cmp b #10 [b] [] [11]
Bltt L1 11 bltt L1 [] [] [2, 12]

1) Point: any location between adjacent nodes.

2) Path: a sequence of points traversing through CFG.

3) Live: A variable is live immediately after a node n if it is live before any of n's successors. A variable is live before a node n if it's used by n, or it's alive after n and **IS NOT** overwritten by n.

4) Live Out: If the temporary is live / used in any nodes afterwards: $LiveOut(n) = \bigcup_{a \in succ(n)} LiveIn(s)$

"Do live not keep the temporary alive after this node?"

5) Live In: If the temporary is used by our current node, or is liveout after the current node unless our current node defines that temporary:
 $LiveIn(n) = uses(n) \cup (LiveOut(n) - defines(n))$

Examples: b is liveout from live 1, as used by line 10. b is live from 2, as its used on the path following 8. b is **NOT** live out from node 4 as we do mov a b - overwriting our old b (so it's GONE). The new b we use on that path is a different one.

The whole point is we define a set for each node, $LiveIn(n)$ (temporaries alive immediately before n) and $LiveOut(n)$ (temporaries alive immediately after n).

Iterative code for Live Ranges

for node in CFG {
LiveIn(node) = {} ; LiveOut(node) = {}
}
repeat {
for each n in CFG {
LiveIn(n) = uses(n) U (LiveOut(n) - defs(n))
LiveOut(n) = $\bigcup_{s \in succ(n)} LiveIn(s)$
}
}
until LiveIn and LiveOut do not change

To improve this method, to update the nodes from last \rightarrow first as data propagates from back to front - as we use successors - so this is more efficient.

9) Loop Invariant Code Motion

An instruction is loop-invariant if its operands are only defined **outside of the loop**. Hence the value it defines is loop-invariant (same for every iteration) and hence it can be moved to move instruction outside the loop.

9.1) Finding Reaching Definitions (Forward DFA)

Formally we attempt to find definition nodes of the form:

$t_d := u_1 * u_2 \mid \frac{u_1}{c}$

Node ID Binary Op Copy Op Constant

Where t_d is the destination, and u_i is for temporary variables used. d is **loop-invariant** if every definition of a $u_i \in uses(d)$ that **reaches** d is outside the loop.

Reaching definitions: A definition d reaches p if there is a path $d \rightarrow p$ where d isn't killed.

Gen(n) = {n} = set of defs generated by the node

Kill(n) = Set of all def of t except for n

ReachIn(n) = Set of definitions reaching up to n

ReachOut(n) = Set of definitions reaching after n

$ReachIn(n) = \bigcup_{p \in pred(n)} ReachOut(p)$

$ReachOut(n) = Gen(n) \cup (ReachIn(n) \setminus Kill(n))$

Informal Algorithm

1. Initialise $ReachIn(n)$ and $ReachOut(n)$ to {}

2. Iterate, updating $ReachIn(n)$ and $ReachOut(n)$ using definitions above, until convergence.

At each step, the sets increase in size

From our reaching definitions, we can reduce each set to the **relevant reaching definitions**, by considering only the reaches that are actually used by the instruction (for operands)

Lines of Code. Reaching definitions (RDs) Relevant RDs

1. x=1 [] []
2. w=100 [1] []
3. z=200 [3, 1, 2] []
4. x=x+1 [1, 2, 3, 4, 5] [1, 4]
5. y=w+z [2, 3, 4, 5] [2, 3]
6. if (x<10) [2, 3, 4, 5] [4]
go back to 4
else continue

We just found that all the definitions used by node 5 lie **outside the loop**!! We can hoist.

9.2) Identifying Loops

A loop in a control flow graph is a set of nodes S including a header node h, with the following properties: 1) From any node in S there is a path leading to h

2) There is a path from h to any node in S

3) There is no edge from any node outside S to any node in S other than h

So there's only way in, through a node in S, and h can lead to all of the nodes back in S

Dominators: A node d dominates a node n if every path from the CFG's start node to n must go through d. Every node dominates itself.

M8) Finding all the dominators of a node

1. Set all Dom sets to the set of all nodes. Set the start node's dom to be itself.

2. Apply the Doms rule.

$Doms(n) = \{n\} \cup \left(\bigcap_{p \in pred(n)} Doms(p) \right)$

As the start node will have a set of {start} this will propagate, reducing the sizes of the sets for other nodes.

3. Once the sets stop changing, we have our solution.

Back Edge: An edge in the CFG from n \rightarrow h where h dominates n is a back edge:

Start $\rightarrow \dots \rightarrow h \rightarrow \dots \rightarrow n \rightarrow \dots$

A dominator h

Natural Loop: The Natural Loop of a back edge (n, h) is the set of nodes S such that:

1) All nodes x \in S are dominated by h

2) For all nodes x \in S (except h), there is a path from x \rightarrow n that does not contain h. This represents a loop, with the header node h.

Multiple Loops can share the same header. But this is not obvious from the CFG. If we have a natural loop in another, then this is a nested loop.

Control Tree

We can construct a tree to show which loops are nested, what the headers and final nodes in each loop are. We then group our nodes with the closest enclosing circle they're in and draw the control tree (children of a node are the nodes in circles within that circle node group)

Pre Header: A node inserted immediately **before the header node** of a natural loop.

9.3) Hoisting Instructions

The Conditions for Hoisting are:

1. **All reaching definitions used by d occur outside the loop:** Use reaching definition analysis for this.

2. **Loop invariant node must dominate all loop exits** Use dominators analysis

3. **There can only be one definition of t** Count the definitions.

4. **t cannot be liveout from the loop's pre-header** Use live range analysis.

Process of hoisting loop-invariant instructions out of a loop is:

1. Compute dominance sets for each node.

2. Use dominance sets to identify natural loop and their headers.

3. Compute the reaching sets for nodes.

4. Use relevant reaching definitions to identify loop-invariant code.

5. Attempt the loop invariant code to a pre-header.

6. Check that the semantics of the program are not altered.

9.3.1) Static Single Assignment (SSA)

An IR which avoids side conditions by only allowing a single assigner temporary. Each **Reassignment of a variable is renamed**, this splits all live ranges. Each variable has only one reaching definition. The phi function is used for branching. A phi statement $\phi(a_1, a_2)$ means either a1 or a2 could be used.

Once in SSA form, we can reassess the requirements for hoisting:

1. All reaching definitions used by d occur outside the loop (Same as prior to SSA). Use reaching definitions analysis for this.

2. Loop invariant node must dominate all loop exits No longer an issue.

3. There can only be one definition of t guaranteed by SSA form.

4. t cannot be liveout from the loop's pre-header Cannot occur with SSA due to single assignment.

start = {A}
nodes = {F, G, K, L}

start = {B_h}
nodes = {C_n, D_n}

start = {B_n}
nodes = {A_n}

start = {H_n}
nodes = {I_n}

start = {J_n}

branch

a₁ = ...

a₂ = ...

t = $\phi(a_1, a_2)$

x = t

The Either chosen