

**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 based of **Regular Expressions**.

**1.1) Identifiers (Tokens of Partial 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) **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 defined 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) R+ - matches adjacent as a combined rule - ab89 matches 'ab89' only

5) R|R - alternation, match one regex or the other. ab|1 matches "abc" or "1"

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

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

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) [^abcd] - match any character except those in the set

We can use these in **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 (or 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+**

**R\***

**R1|R2**

**R1 R2**

**R1 R2**

**R1 R2**

**R1 R2**

**R1 R2**

**R1 R2**

**R1 R2**

**R1 R2**

**R1 R2**

**R1 R2**

**1) NFA: Space:** O(n), **Time:** O(n \* len x)

**2) DFA: Space:** O(2<sup>len</sup>), **Time:** O(n \* len x)

**1.4) M2) NFAs to DFA (subset construction)**

1) Start at the NFA start state.

2) Compute the ε-closure for each state.

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 to our epsilon closure to another!! Also, any node that contains an accepting state in its DFA. This accepting state can still have transitions.**

**2) Parsing - LR Parsing**

**2.1) LR(0) Parsers. 2.2.1) LR(0) Items**

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

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.

**2.2.2) NFA from LR(0) Items**

We can create transitions between LR(0) Items

1) Given an item X → A•BC →<sub>a</sub> we get X → AB•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 → AB•C →<sub>a</sub> (B →•D) **This means if we're about to process a token that's a rule in a NFA state, we have to add the productions of this rule to our NFA state.** To convert to the DFA from here, we do subset construction, but its simple enough right away.

**2.3.3) DFA to LR(0) Parsing Table**

X →•Y Add P[X, T] = sy (shift Y, we cannot reduce yet)

X →•Y Add P[X, T] = sy (shift Y, we cannot reduce yet)

X →•Y Add P[X, T] = sy (shift Y, we cannot reduce yet)

X →•Y Add P[X, T] = sy (shift Y, we cannot reduce yet)

X →•Y Add P[X, T] = sy (shift Y, we cannot reduce yet)

X →•Y Add P[X, T] = sy (shift Y, we cannot reduce yet)

X →•Y Add P[X, T] = sy (shift Y, we cannot reduce yet)

X →•Y Add P[X, T] = sy (shift Y, we cannot reduce yet)

X →•Y Add P[X, T] = sy (shift Y, we cannot reduce yet)

X →•Y Add P[X, T] = sy (shift Y, we cannot reduce yet)

X →•Y Add P[X, T] = sy (shift Y, we cannot reduce yet)

X →•Y Add P[X, T] = sy (shift Y, we cannot reduce yet)

X →•Y Add P[X, T] = sy (shift Y, we cannot reduce yet)

X →•Y Add P[X, T] = sy (shift Y, we cannot reduce yet)

X →•Y Add P[X, T] = sy (shift Y, we cannot reduce yet)

X →•Y Add P[X, T] = sy (shift Y, we cannot reduce yet)

X →•Y Add P[X, T] = sy (shift Y, we cannot reduce yet)

X →•Y Add P[X, T] = sy (shift Y, we cannot reduce yet)

X →•Y Add P[X, T] = sy (shift Y, we cannot reduce yet)

X →•Y Add P[X, T] = sy (shift Y, we cannot reduce yet)

X →•Y Add P[X, T] = sy (shift Y, we cannot reduce yet)

X →•Y Add P[X, T] = sy (shift Y, we cannot reduce yet)

X →•Y Add P[X, T] = sy (shift Y, we cannot reduce yet)

X →•Y Add P[X, T] = sy (shift Y, we cannot reduce yet)

X →•Y Add P[X, T] = sy (shift Y, we cannot reduce yet)

X →•Y Add P[X, T] = sy (shift Y, we cannot reduce yet)

X →•Y Add P[X, T] = sy (shift Y, we cannot reduce yet)

X →•Y Add P[X, T] = sy (shift Y, we cannot reduce yet)

**2.4) LR(1) Parsers**

→ Zero or more derivations.

→ One or more derivations.

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

**M3) Computing the First Set:**

In each alternative we iterate through each terminal/non-terminal B, if e is first(B) then we include first(B) \ {ε} in the first set and move on the next in the sequence

**Example: 1) Factor** → ('(' Expr ')' | id Term2 → '\*' Factor Term2 | ε

3) Term2 → Factor Term2 | ε

4) Expr2 → '+' Term Expr2 | ε

5) Expr → Term Expr2 | ε

6) Expr → Expr \$

FIRST(Factor) = {'(', id}

FIRST(Term2) = {'\*', ε}

FIRST(Term) = FIRST(Factor) = {'(', id}

**M4) Computing the Follow Set:**

The follow set for a non-terminal rule is the set of all tokens (terminals) that could immediately follow, and \$ if the end of the input could follow. Hence we must check each rule which contains the non-terminal.

Given rule B → C•AD we have follow(A) = first(D) \ {ε} ∪ {follow(B) if D →•\* ε

C can be empty here, since it is at the start of the rule, it has no bearing on the follow set. If A can end the input, we include \$ in follow set.

**Example for the rules above:**

Follow(Expr) = {'\$', ε}

Follow(Term) = FIRST(Expr2) + Follow(Expr)

+ Follow(Expr2) (because Expr2 can be a null string and thus we need the follow of Expr2) = {'+', '\$', ε}

Follow(Term2) = Follow(Term) + Follow(Term2) = {'\*', '\$', ε}

Follow(Factor) = Follow(Term) + Follow(Factor) = {'(', id, '\$', ε}

Follow(Expr2) = Follow(Expr) + Follow(Expr2) = {'+', '\$', ε}

Follow(Expr) = Follow(Expr) + Follow(Expr) = {'\$', ε}

Follow(Expr) = Follow(Expr) + Follow(Expr) = {'\$', ε}

Follow(Expr) = Follow(Expr) + Follow(Expr) = {'\$', ε}

Follow(Expr) = Follow(Expr) + Follow(Expr) = {'\$', ε}

Follow(Expr) = Follow(Expr) + Follow(Expr) = {'\$', ε}

Follow(Expr) = Follow(Expr) + Follow(Expr) = {'\$', ε}

Follow(Expr) = Follow(Expr) + Follow(Expr) = {'\$', ε}

Follow(Expr) = Follow(Expr) + Follow(Expr) = {'\$', ε}

Follow(Expr) = Follow(Expr) + Follow(Expr) = {'\$', ε}

Follow(Expr) = Follow(Expr) + Follow(Expr) = {'\$', ε}

Follow(Expr) = Follow(Expr) + Follow(Expr) = {'\$', ε}

Follow(Expr) = Follow(Expr) + Follow(Expr) = {'\$', ε}

Follow(Expr) = Follow(Expr) + Follow(Expr) = {'\$', ε}

Follow(Expr) = Follow(Expr) + Follow(Expr) = {'\$', ε}

Follow(Expr) = Follow(Expr) + Follow(Expr) = {'\$', ε}

Follow(Expr) = Follow(Expr) + Follow(Expr) = {'\$', ε}

Follow(Expr) = Follow(Expr) + Follow(Expr) = {'\$', ε}

Follow(Expr) = Follow(Expr) + Follow(Expr) = {'\$', ε}

Follow(Expr) = Follow(Expr) + Follow(Expr) = {'\$', ε}

Follow(Expr) = Follow(Expr) + Follow(Expr) = {'\$', ε}

Follow(Expr) = Follow(Expr) + Follow(Expr) = {'\$', ε}

**For Reduce-Reduce conflicts**, caused by two rules having the same RHS, we add precedence to some rules (e.g. earliest).

**2.6) Parse Trees**

**Leaf nodes built on a shift operations**

**Non-Leaf nodes created on a reduce.**

To make AST we do a pass over Parse Tree, or have AST construction rules to build AST while parsing. To make these properly, just look at our input and figure out how we could construct it (what rules were taken and uses).

**2.7) LR Parsing**

**Top-down** - recursive descent or DFA. Hand-coded or generated.

**LL(k) Grammar** - a grammar is LL(k) if a k token lookahead is sufficient to determine which alternative of a rule to use when parsing.

LL(1) uses the current token ONLY. LL(0) could not exist as it doesn't use ANY TOKENS. **Leaf Nodes are constructed from the root.**

**2.7.1) LL(1) Grammar**

A grammar is LL(1) if for a rule A → α | β (A is non-terminal):

first(α) ∩ first(β) = ∅

ε ∈ first(α) → (first(β) ∩ follow(A) = ∅)

ε ∈ first(β) → (first(α) ∩ follow(A) = ∅)

AKA - if α and β don't share first tokens, or if ε is in the first of A or B, then the first of B/A follow A don't share anything, makes sense!

**2.7.2) Extended Backus Naur Form (EBNF)**

Used to write CFGs. **{α} = 0 or more of α**

**[α] = 0 or 1. (α) = groups elems together.**

We left factors with alternatives with intersecting first sets: Expr → Term '+' Expr | Term becomes Expr → Term {'+' Expr}

We can remove left recursion: Sequence → Sequence' | Statement | Statement check (IF, PRIF, BEGIN;), syncret, "Error..."

Sequence → Sequence' | Statement | Statement check (IF, PRIF, BEGIN;), syncret, "Error..."

Sequence → Sequence' | Statement | Statement check (IF, PRIF, BEGIN;), syncret, "Error..."

Sequence → Sequence' | Statement | Statement check (IF, PRIF, BEGIN;), syncret, "Error..."

Sequence → Sequence' | Statement | Statement check (IF, PRIF, BEGIN;), syncret, "Error..."

Sequence → Sequence' | Statement | Statement check (IF, PRIF, BEGIN;), syncret, "Error..."

Sequence → Sequence' | Statement | Statement check (IF, PRIF, BEGIN;), syncret, "Error..."

Sequence → Sequence' | Statement | Statement check (IF, PRIF, BEGIN;), syncret, "Error..."

Sequence → Sequence' | Statement | Statement check (IF, PRIF, BEGIN;), syncret, "Error..."

Sequence → Sequence' | Statement | Statement check (IF, PRIF, BEGIN;), syncret, "Error..."

Sequence → Sequence' | Statement | Statement check (IF, PRIF, BEGIN;), syncret, "Error..."

Sequence → Sequence' | Statement | Statement check (IF, PRIF, BEGIN;), syncret, "Error..."

Sequence → Sequence' | Statement | Statement check (IF, PRIF, BEGIN;), syncret, "Error..."

Sequence → Sequence' | Statement | Statement check (IF, PRIF, BEGIN;), syncret, "Error..."

Sequence → Sequence' | Statement | Statement check (IF, PRIF, BEGIN;), syncret, "Error..."

Sequence → Sequence' | Statement | Statement check (IF, PRIF, BEGIN;), syncret, "Error..."

Sequence → Sequence' | Statement | Statement check (IF, PRIF, BEGIN;), syncret, "Error..."

Sequence → Sequence' | Statement | Statement check (IF, PRIF, BEGIN;), syncret, "Error..."

Sequence → Sequence' | Statement | Statement check (IF, PRIF, BEGIN;), syncret, "Error..."

Sequence → Sequence' | Statement | Statement check (IF, PRIF, BEGIN;), syncret, "Error..."

Sequence → Sequence' | Statement | Statement check (IF, PRIF, BEGIN;), syncret, "Error..."

Sequence → Sequence' | Statement | Statement check (IF, PRIF, BEGIN;), syncret, "Error..."

Sequence → Sequence' | Statement | Statement check (IF, PRIF, BEGIN;), syncret, "Error..."

Sequence → Sequence' | Statement | Statement check (IF, PRIF, BEGIN;), syncret, "Error..."

Sequence → Sequence' | Statement | Statement check (IF, PRIF, BEGIN;), syncret, "Error..."

Sequence → Sequence' | Statement | Statement check (IF, PRIF, BEGIN;), syncret, "Error..."

Sequence → Sequence' | Statement | Statement check (IF, PRIF, BEGIN;), syncret, "Error..."

**As an example of Left Recursion Removal:**

Expr → Expr {'+' | '-' } Term | Term

Term → Term {'\*' | '/' } F factor

Factor → '(' Expr ')' | Int

Parse tree may no longer represent the association, hence we will need to ensure the arithmetic is still associative when we construct the AST:

Expr → Term {'+' | '-' } Term

Term → Factor {'\*' | '/' } Factor

Factor → '(' Expr ')' | Int

**2.7.6) Error Recovery**

We need to emit good error messages, and skip as little code as possible to recover. **Panic Mode**

**Recovery:** Each parse function has a syncret of tokens. When an error occurs the parser skips forward until it encounters one of these tokens.

**2.7.6.1) Panic Mode Error Recovery**

**Panic Mode Recovery:** Each parse function has a syncret of tokens (usually the follow set - provided as extra args to the function).

When an error occurs the parser skips forward until it encounters one of these tokens (thus repairing the state).

We improve our check function:

def check (expectset, syncret, error):

if next\_token() not in expectset:

add error (next\_token(),

parser\_pos(),

expectset, error)

skipto(expectset + syncret)

skipto simply pop. token(s) while our next token isn't in syncret or EOF.

We use check at the top of our higher parsers to make sure things are ok. We use match in our rules parsers to match:

def Statement(syncret):

check (IF, PRIF, BEGIN;), syncret, "Error..."

if token == IF: ifStatement(syncret)

if token == IF: ifStatement(syncret)

if token == IF: ifStatement(syncret)

if token == IF: ifStatement(syncret)

if token == IF: ifStatement(syncret)

if token == IF: ifStatement(syncret)

if token == IF: ifStatement(syncret)

if token == IF: ifStatement(syncret)

if token == IF: ifStatement(syncret)

if token == IF: ifStatement(syncret)

if token == IF: ifStatement(syncret)

if token == IF: ifStatement(syncret)

if token == IF: ifStatement(syncret)

if token == IF: ifStatement(syncret)

if token == IF: ifStatement(syncret)

if token == IF: ifStatement(syncret)

if token == IF: ifStatement(syncret)

if token == IF: ifStatement(syncret)

if token == IF: ifStatement(syncret)

if token == IF: ifStatement(syncret)

if token == IF: ifStatement(syncret)

**4.1) Structs and Records**

1) Allocated in contiguous block of memory.

2) Alignment is used to space fields out correctly.

3) C orders fields by their code position, other compilers optimize (its important for us in C to optimize this)

4) Tuples are anonymous structs.



## 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 Define a base term.  
expr  $\rightarrow$  expr - term Left associative  
expr  $\rightarrow$  term - expr Right associative  
expr  $\rightarrow$  term

### 5.1.1.2) Precedence

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

exp  $\rightarrow$  exp + term | exp - term | term  
term  $\rightarrow$  term \* 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.

### 5.1.1.3) 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't 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 Op = Plus | Minus | Tim | type Name = [Char]

**2) Assembly Instructions:** data Instruction = Add | Sub | Mul | Div | PushImm Int | PushAbs Name | Pop Name | CompEq | 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,True "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) Combine the register and accumulator approach: Use registers as normal, and when we run out of registers take the Accumulator machine approach.

Our only code change from the before is the BinOp case with two exprs in transExp (and we pass in the register we want to store on in transExp):

(final case)  
transExp (BinOp op e1 e2) r =  
| r == maxReg = transExp e2 r ++  
| [Push r] ++ transExp e1 r ++  
| [translateOpStack op r]  
| otherwise = transExp e1 r ++  
| transExp e2 (r + 1) ++  
| [translateOpOp r (r + 1)]

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

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

**3) Register Allocation**  
transOp :: Op  $\rightarrow$  (Int  $\rightarrow$  Instruction) %eax %edx %ecx

transOpImm Plus = AddImm  
transOpImm Minus = SubImm  
transOpImm Times = MultImm  
transOpImm Divide = DivImm

transOp :: Op  $\rightarrow$  Instruction is obvious

transExp :: Exp  $\rightarrow$  [Instruction]  
transExp (Const n) = [LoadImm n]  
transExp (Ident x) = [Load x]  
transExp (Unop Minus e) = transExp e ++ [MultImm (-1)]

Const Int  
data Op = Plus | Minus | Tim | type Name = [Char]

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

data Op = Plus | Minus | Tim | type Name = [Char]

data Instruction = Add | Sub | Mul | Div | PushImm Int | PushAbs Name | Pop Name | CompEq | Jump Label | JTrue Label | JFalse Label | Define Label

data Op = Plus | Minus | Tim | type Name = [Char]

data Instruction = Add | Sub | Mul | Div | PushImm Int | PushAbs Name | Pop Name | CompEq | Jump Label | JTrue Label | JFalse Label | Define Label

data Op = Plus | Minus | Tim | type Name = [Char]

data Instruction = Add | Sub | Mul | Div | PushImm Int | PushAbs Name | Pop Name | CompEq | Jump Label | JTrue Label | JFalse Label | Define Label

data Op = Plus | Minus | Tim | type Name = [Char]

data Instruction = Add | Sub | Mul | Div | PushImm Int | PushAbs Name | Pop Name | CompEq | Jump Label | JTrue Label | JFalse Label | Define Label

data Op = Plus | Minus | Tim | type Name = [Char]

data Instruction = Add | Sub | Mul | Div | PushImm Int | PushAbs Name | Pop Name | CompEq | Jump Label | JTrue Label | JFalse Label | Define Label

data Op = Plus | Minus | Tim | type Name = [Char]

data Instruction = Add | Sub | Mul | Div | PushImm Int | PushAbs Name | Pop Name | CompEq | Jump Label | JTrue Label | JFalse Label | Define Label

data Op = Plus | Minus | Tim | type Name = [Char]

data Instruction = Add | Sub | Mul | Div | PushImm Int | PushAbs Name | Pop Name | CompEq | Jump Label | JTrue Label | JFalse Label | Define Label

data Op = Plus | Minus | Tim | type Name = [Char]

## 7.1) Register Usage

**1) Sethi Ullman Weights**

Given an expression E1 op E2 always evaluate the subexpression that uses most registers 1st.

1) If E1 evaluated first, registers needed is max(E1, E2 + 1)

2) If E2 evaluated first, registers needed is max(E1 + 1, E2)

This is nicely encoded by the weight function weight :: Exp  $\rightarrow$  Int

weight (Const) = 1  
weight (Ident) = 1  
weight (Unop e) = weight e  
weight (BinOp Times (Const) e) = weight e  
weight (BinOp Plus (Const) e) = weight e  
weight (BinOp e (Const)) = weight e

-- Use maximum of either  
weight (BinOp e1 e2) = min e1.f e2.f  
where

e1.f = max (weight e1) (weight e1 + 1)  
e2.f = max (weight e2) (weight e1 + 1)

**This doesn't work for divide or subtraction** - as they can't be reordered, instead we use register targeting.

**2) Register Targetting**

We give transExp a list of all the registers that are free and by convention we want the result of the operation to be stored into the first register in the list:

transExp :: Exp  $\rightarrow$  [Register]  $\rightarrow$  [Instruction]  
transExp (Const n) (dst:rst) = [LoadImm dst n]  
transExp (Ident x) (dst:rst) = [Load dst x]

The other cases are standard. Binop case:  
transExp :: Exp  $\rightarrow$  [Register]  $\rightarrow$  [Instruction]  
transExp (Const n) (dst:rst) = [LoadImm dst n]  
transExp (Ident x) (dst:rst) = [Load dst x]

transExp (BinOp op e1 e2) (dst) =  
transExp e2 (dst:rst) ++ [Push dst] ++  
transExp e1 (dst+1) ++ [transOpStack op dst]  
transExp (BinOp op e1 e2) (dst:next:rst) =  
| weight e1 > weight e2 =  
| transExp e1 (dst:next:rst) ++  
| transExp e2 (next:rst) ++  
| [transOp op dst next]  
| otherwise =  
| transExp e2 (next:dst:rst) ++  
| transExp e1 (dst:rst) ++  
| [transOp op dst next]

This works well, as the expression size accommodated by N registers, is 2<sup>n</sup>.

**7.2) Register Allocation for Function Calls**

Need to know where parameters are when called. ((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).

**7.2.1) 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) invalid, a,d,e,c invalid due to how return works)

**7.2.2) 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 (e.g. Saved 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**

**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 l0...). 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) 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) Phase 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 operation with a simpler one).

**3) Control variable selection** - replace loop control variable (as in the 'in' for '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 used 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 l0... instead of regs.

2. For each temporary T<sub>i</sub>, find T<sub>i</sub>'s live range - a set of instructions for which T<sub>i</sub> must reside in a register.

3. If liveRange(T<sub>i</sub>) intersects liveRange(T<sub>j</sub>) then they must be allocated to different registers - **they interfere**.

4. Assemble the Register Interference 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 CFGNode = ControlFlowGraph (CFGNode) data CFGNode = Node ID Instruction [Register] [Id] type ID = Int

data Register = D Int | T Int buildCFG :: [Instruction]  $\rightarrow$  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 [ Bra L2 [] [] [10] ]  
L1:

cmp b a 2 cmp b a [b, a] [] [3]  
bge l3 3 bge l3 [b] [] [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 [b] [] [10]

L3:  
mov b a 8 mov b a [b] [a] [9]  
sub #1 a 9 sub #1 a [a] [a] [10]

L4:  
cmp b #10 10 cmp b #10 [b] [] [11]  
blt L1 11 blt L1 [b] [] [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.

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

"Do we need to 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) = U<sub>s in succ(n)</sub> 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 all defined **outside of the loop**. Hence the value it defines is loop-invariant (same for every iteration) and hence it may be possible to move instruction outside the loop.

**9.1) Finding Reaching Definitions (Forward DFA)**

Formally we attempt to find definition nodes of the form:  $l_i = t_i = (u_1 \bullet u_2 \mid u_1 \mid c_i)$

Where  $u_i$  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

$ReachIn(n) \triangleq \bigcup_{p \in Pre(n)} ReachOut(p)$   
 $ReachOut(n) \triangleq Gen(n) \cup (ReachIn(n) \setminus Kill(n))$

**Informal Algorithm**  
Initialise  $ReachIn(n)$  and  $ReachOut(n)$  to {}  
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 instructions (for operands)

Lines of Code Reaching definitions (RDs) Relevant RDs  
1: x=1 [ ] [ ]  
2: w=100 [ ] [ ]  
3: z=z-200 [1, 2] [ ]  
4: x=x+1 [1, 2, 3, 4, 5] [1, 4]  
5: y=y+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 set to 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 {} that 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:

**Pre Header:**

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

**Start**  $\rightarrow \dots \rightarrow h$  dominates n  $\rightarrow \dots \rightarrow n$

back edge (n,h)

h dominates n

n

back edge (n,h)

h dominates n

n

back edge (n,h)

h dominates n

n

back edge (n,h)

h dominates n

n

back edge (n,h)

h dominates n