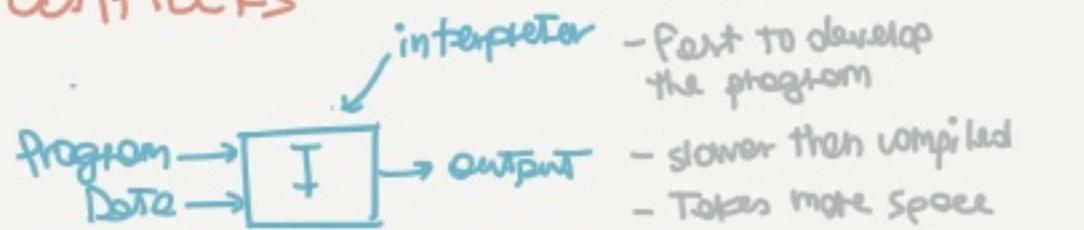
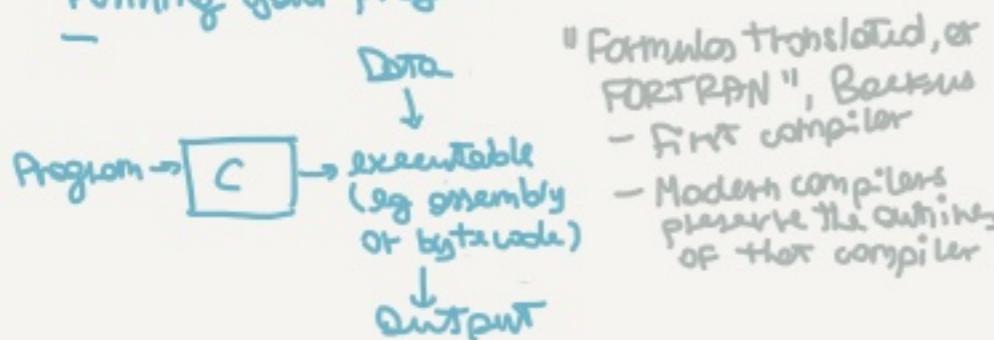


COMPILERS

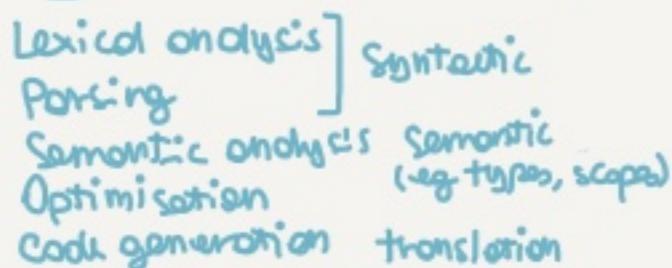


It doesn't do any processing of the program before it executes.
It immediately begins running on input

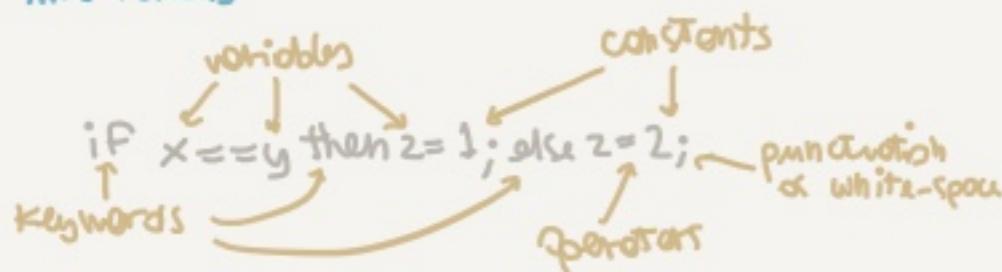
We can say it is on-line, meaning the work it does is all part of running your program



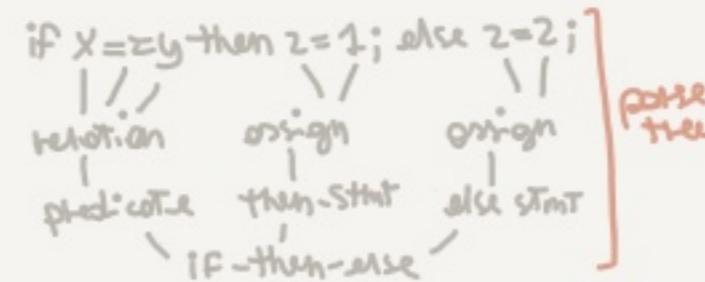
A compiler is off-line. It pre-processes the program first that produces an exec that can be run on many different inputs without having to recompile or do any other processing of the program



Lexical analysis, divides the program into tokens



Parsing, understand the structure of the sentence



Then we want to understand the meaning, and this is hard!

Compilers can only do very limited semantic analysis. It generally tries to catch inconsistencies

Programming languages define rules to avoid ambiguities

{ j0u=1
{ j0uk=3
print j0u }
3

it will print 3
as the outer definition is hidden

Optimisation, to make use of some resource

- Run faster
 - less memory
 - Power
 - Network usage
- X = Y ≠ 0
↓
X = 0
...
- This is NOT a correct rule.
It's not always obvious when to do optimisation.

This is valid for integers only but not for floating points
(NAN ≠ 0 = NAN)

Code Generation
(or code gen.)

It usually produces assembly code. But in general is a translation into some other language

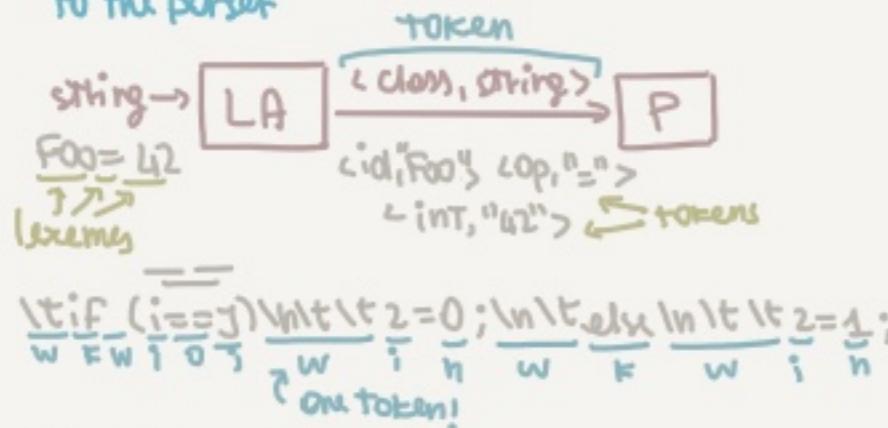
Lexical Analyses

```
if(i==j)
z=0;    => \t if(i==j) \n \t \t z=0; \n \t
else
z=1;    ... else \n \t \t z=1;
        ^ what the lexer receives
```

It divides the string into tokens and classifies them into token classes

- identifiers, string of letters or digits starting with a letter
- integer, a non-empty string of digits
- keywords, reserved words such as "if", "begin"
- whitespace, a non-empty sequence of blanks, newlines, and tabs

It will then communicate the tokens to the parser



operator |
 whitespace |
 keywords :
 Identifiers :
 Numbers =

classes by themselves

Lexical analysis examples

FORTRAN: whitespace is insignificant
 DO 5 I=1,25 \rightarrow FORTRAN loop VA R1 is the same as VTR2
 DO 5 I=1.25 \leftarrow variable assign.

To understand the role of DO we look ahead until there is no ambiguity (. or ; in that case)

Lookheads complicates the implementation of lexical analysis; we want to minimize them when building a compiler

When partitioning the string from left to right, lookahead may be required to decide where one token ends and the next begins

if (i == j)
 is 'i' a variable
 is this an assignment
 else
 z=1;

Regular languages

which set of strings belongs to each token is specified by a regular language

To define a regular language we use regular expressions.

$C = \{c\}$ language with only one character
 $\Sigma = \{\}$ empty string. It is NOT an empty language

$A + B = \{a|a \in A\} \cup \{b|b \in B\}$ Union

$AB = \{ab|a \in A \wedge b \in B\}$ concat.

$A^* = \bigcup_{i \geq 0} A^i$ $A^i = A \dots A$

$A^0 = \Sigma$ all strings

$\Sigma = \{0,1\}$ $1^* = \bigcup_{i \geq 0} 1^i = 1 + \dots + 1$

$(1+0)^* = \{ab|a \in 1+0 \wedge b \in 1\}$
 $= \{11, 01\}$

$0^k + 1^k = \{0^i 1^j 2^m 3^n | i \geq 0\}$
 $\hookrightarrow (0+1)^* = \bigcup_{i \geq 0} (0+1)^i = "", 0+1, 00+1, 0+10+1, \dots$

$R = \emptyset$
 $\cup \Sigma \subseteq \Sigma$

$| R+R$ grammar

$| R^k$ for e RE

also denoted

$\hookrightarrow \Sigma^*$

Formal languages

A formal language has an alphabet Σ

A language over Σ is a set of strings of characters drawn from Σ

If $\Sigma = \text{ASCII}$ Language = C programs

A meaning function maps syntax to semantics

$L(a) = M$
 $\hookrightarrow \text{reg exp}$ $\hookrightarrow \text{set of strings}$

We use meaning function to

- makes clear what is syntax, what is semantics
- allow us to consider notation as a separate issue
- Because expressions and meaning are not in one to one correspondence

Meaning is many to one, never one to many.

The meaning of expressions in our programming language must be unambiguous.

Lexical specifications

keyword: "if" or "else" or "then" or...

CONCAT. \hookrightarrow 'if' + 'else' + ...

integer: non-empty string of digits

digit = '0' + '1' ... '9' digit digit*
 \hookrightarrow most tools allow homing "digit")
 \hookrightarrow digits

identifier: strings of letters or digits starting with a letter

letter = 'a' + 'b' + 'c' ... 'z' + 'A' ... 'Z'
 $\rightarrow [a-zA-Z]$ character range

letter (letter + digit)*

whitespace: non-empty sequence of blanks, newline and tabs
 $\hookrightarrow (" " + "\n" + "\t")^*$

e.g. PASCAL

digit = '0' + ... '9' that means optional

digits = digit +

opt-fraction = ('.' digits) + { } $\rightarrow (.\{digits\})?$

opt-exponent = ('E' ('+' '-') digits) + { }

num = digits opt-fraction opt-exponent

$L('c') = \{ 'c' \}$
 $L(A+B) = L(A) + L(B)$

Lexical specification

A^* , $A B^*$
 $A \setminus B$, $A + B$ any character
 $[a_1 a_2 \dots a_n]$, $[a_1 - a_n]$ except that range
 complement of $[a_1 - a_n]$ $[^a - ^n]$

1 number=digit
 keyword= ...
 identifier=letter(lowercase or uppercase)
 openpar='(' + digit)*
 ...

first step of lexical Spec. is to write a RE for each syntactic category in the language

3 Let input be $x_1 \dots x_n$
 for $1 \leq i \leq n$ check
 $x_1 \dots x_i \in L(R)$
 take the input, for every prefix of that input we check whether it is in the language of R

5 Remove $x_1 \dots x_i$ from input and go to step 3

Remove that prefix and look for the next one, until the string is empty

In the case where

- $x_1 \dots x_i \in L(R)$
- $x_1 \dots x_i \in L(R_1)$
- $x_1 \dots x_i \in L(R_2)$
- if $\in L(\text{Keywords})$
- if $\in L(\text{Identifiers})$

Choose the one listed first. They have priority ordering

IF no rule matches

$x_1 \dots x_i \notin L(R)$
 Don't let this happen!
 ERROR = all strings not in the lexical spec
 put it last in priority
 (so that it can overlap other RE's in the language)

Finite automate QUICK RECAP!

is a 5-tuple with

- Σ alphabet
- S , states
- s_0 , start state
- $F \subseteq S$ final states
- set of transitions

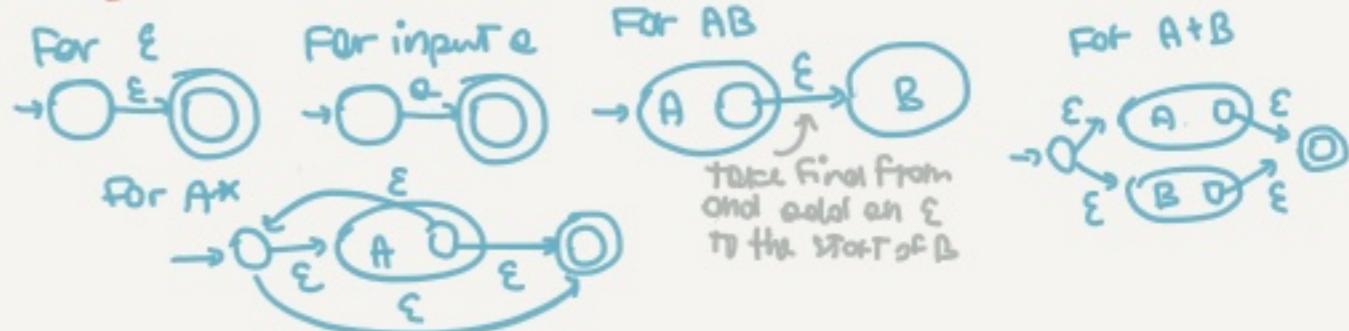
$s_1 \xrightarrow{a} s_2$ In s_1 with input a go to s_2

If after reading the string you're in a final state, accept, the string is in the language, else reject.
 If the machine gets stuck (no transitions on the input), also reject.

NFA's have ϵ -moves, meaning they can either change state without consuming any input, or to not move. More transitions per state makes the automata non-deterministic.

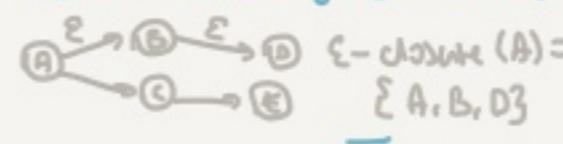
An NFA then, can wind up in any number of states. We accept if an accept state is reached after reading the string. An NFA is not more powerful than a DFA.

Regex TO NFA



NFA TO DFA

ϵ -closure, starting from a state, it is all the states we can reach by only following ϵ -moves.



An NFA can be in many states at the same time. But how many?

N states
 $1 \leq N \leq 2^n - 1$
 possible subsets of n states

NFA

states S
 start $s \in S$
 final $F \subseteq S$

$\alpha(X) = \{y | x \in X \wedge x \xrightarrow{\alpha} y\}$
 for a set of states X , show the all reachable states via α .

DFA

states, subsets of S
 start, ϵ -closure(s)

final, $\{x | x \in F \neq \emptyset\}$
 $x \xrightarrow{\alpha} y$ if
 $y = \epsilon\text{-closure}(\alpha(x))$

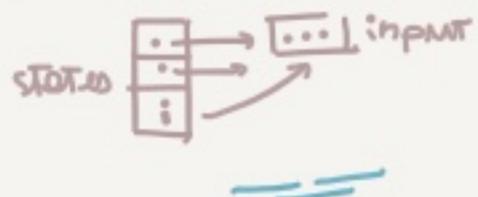
Implementing DFA

- A DFA can be implemented as a two-dimensional array
- one dimension is the state
- other is input symbol
- for every transition $s_i \rightarrow s_k$ define $T[i, a] = k$

input symbols

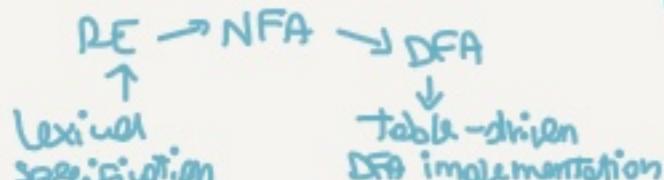
s_0, s_1, \dots, s_k	// store input and state index 0
$i=0;$	$s_0 = 0;$
$\text{while } (\text{input}[i])$	$s_0 = A[s_0, \text{input}[i+1]]$

With that representation there are lots of repetitions
To avoid that we could...



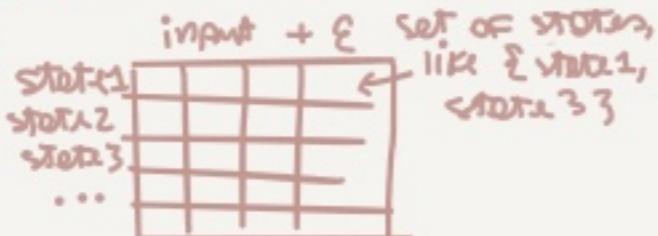
There's a pointer to the input rows.
That way you can share duplicated rows.
Pointers to other tables will need lookup a bit slower.

Lexical analyser construction



What if we want from NFA directly to table,
maybe because our specification would make
the DFA table huge.
We can do that by...

Now the inner loop is
gonna be much more
expensive.



Parsing

Regular languages can't express all languages

i.e. $\{(i)^i \mid i \geq 0\}$ language of balanced parenthesis

It can only count how many k is the number of states

It cannot count to an arbitrary i , with a finite set of states

A parser takes a sequence of tokens from the lexer and it produces a parse tree of the program

CoC, if $x=y$ then 1 else 2 fi
 ↓
 IF ID=ID THEN INT ELSE INT FI
 ↓
 IF-THEN-ELSE — INT
 ↓
 ID = ID — INT — parser input
 ↓
 parser output

In some compilers, the parser is powerful enough to exert lexical analysis in addition to parsing

Context-free grammars

We need a way to describe valid strings of tokens and a method to distinguish valid ones

Programming languages have a hierarchical structure.

CoC, - Expr ← can be if while...
 - if Expr then Expr else Expr fi

These expressions are recursively composed of other expressions

Context-free grammars are used to describe such recursive structures.

A CFG consists of

- set of terminals T
- set of non-terminals N
- start symbol S
- set of productions

$X \rightarrow y_1 \dots y_n$ ← production

$X \in N, y_i \in N \cup T \cup \{E\}$

—

e.g. $S \rightarrow (S)$

$S \rightarrow LS$
 $S \rightarrow ((S))$

- Begin with S , top symbol initially
- Replace any non-terminal X by the right hand side of some production
- Repeat until there are no terminals

Let G be a CFG, $L(G)$ is the language of G

—

- Terminals are so called because there are no rules for replacing them. Once generated, they are permanent

- Terminals ought to be the tokens

e.g. CoC

$\text{Expr} \rightarrow \text{if Expr then Expr else Expr fi}$
 $\text{while Expr loop Expr pool id ...}$

$E \rightarrow E+E \mid E * E \mid (E) \mid id$

- A CFG G just tells us whether the string is in $L(G)$ or not

- We want to handle those errors gracefully

Derivation

$S \rightarrow \dots \rightarrow \dots \xleftarrow{\text{series of productions}}$

can be drawn as a tree

$$\begin{array}{l} E \rightarrow E+E | E * E | (E) | id \\ E \rightarrow E+E \rightarrow \\ E+E+E \rightarrow \\ id+E+E \rightarrow \\ id+id+E \rightarrow \\ id+id+id \end{array}$$

Parse trees have terminals at the leaves and non-terminals at the interior nodes.

As in order traversal of the leaves is the original input

It shows associations of operations.
In the example the "*" binds more tightly than the +, as it is at a lower level, we do * first before +.

The example shows a left-most derivation where at each step, we replace the left-most non-terminal

Right-most and left-most derivations have the same parse tree

Ambiguity

$$\begin{array}{l} E \rightarrow E+E | E * E | (E) | id \\ id * id + id \end{array}$$

$$\begin{array}{l} E \rightarrow E+E \rightarrow \\ id \rightarrow id \\ E \rightarrow E+E \rightarrow \\ E * E \rightarrow \\ E * E+E \rightarrow \\ id * id + id \end{array}$$

$$\begin{array}{l} E \rightarrow E+E \rightarrow \\ id \rightarrow id \\ E \rightarrow E * E \rightarrow \\ id \rightarrow id \\ E \rightarrow E+E \rightarrow \\ E * E \rightarrow \\ E * E+E \rightarrow \\ id * id + id \end{array}$$

There are two parse trees for this grammar

Each of these parse trees have many derivations.
This is not about derivations.

This grammar is ambiguous; it has more than one parse trees for some string.

Ambiguity is bad, leaves meaning of some programs ill-defined

To eliminate ambiguity, write a new grammar that is unambiguous

$$\begin{array}{l} E \rightarrow E+E | E * E | (E) | id \\ E \rightarrow id + E | id | (E) * E | (E) \end{array}$$

Enforces precedence of * over +!

—
if E_1 then if E_2 then E_3 else E_4



"use" must match with the closest unmatched "then"

$$\begin{array}{l} E \rightarrow MIF / \text{"all "then" are matched} \\ \quad UIF / \text{"some "then" are matched"} \end{array}$$

MIF \rightarrow IF E then MIF else MIF | OTHER

UIF \rightarrow IF E then E | if E then MIF else UIF

It's impossible to automatically correct an ambiguous grammar to an unambiguous one

Used with care, ambiguities can simplify the grammar

- Allows more natural definitions
- We need some disambiguation method

Most practical parsing tools use this approach, the tool provides disambiguation declarations

Most tools allow precedence and associativity declarations

$$\begin{array}{l} E \rightarrow E+E | int \\ \text{left associativity} \end{array}$$

we declare + to be left associative

$$\begin{array}{l} E \rightarrow E+E | E * E | int \\ \text{precedence} \end{array}$$

% left + "Not too mixed
in bison"

the order gives precedence to higher precedence

Top-down parsing

In compiler

- translates valid programs
- detect non-valid programs

Many kinds of errors, lexical, syntax, semantic or context errors

An error handler should...

- Report errors accurately and clearly
- Recover from an error quickly
- Not slow down compilation of valid code

panic mode is the simplest, most popular method

When an error is detected

- Discard tokens until one with a clear role is found
- Continue from there

Look for synchronizing tokens, to identify which one.

e.g. $(1++2)+3 \leftarrow$ ++ not in the language, just skip it

$$E \rightarrow \dots | \text{error} \int \text{int} | (\text{error})$$

Bison uses the special terminal error to describe how much input to skip

They inserted last so that they'll run only if the previous ones didn't match



•••

To start the parser

- Initialise next to point to the first token
- Invoke $E()$, to match everything derivable from the `start` symbol

Recurive descent algorithm limitations

$E \rightarrow \text{int}$	$E \rightarrow \text{int} * \text{int}$	$E \rightarrow T T + E$
$E_1 \downarrow$	$E_2 \downarrow$ only you had int, but except and repeat the rest of the input	$T \rightarrow \text{int} \text{int} * T$
$\downarrow T$	$\downarrow T$	$I(E)$
$\downarrow \text{int}$	$\downarrow \text{int}$	

While there is backtracking for a non-terminal, there is no backtracking once we found a production that succeeds for a non-terminal.

General descent algorithms have this kind of backtracking

This algorithm however

- is easy to implement by hand
- sufficient for grammars where for only non-terminal, at most one production can succeed
- grammars can be rewritten to work with this algorithm.

Left recursion

$S \rightarrow S_1$ $\text{bad } S_1() \{\text{return } S_1();\}$
 $S \rightarrow S_1$ $\{\text{return } S_2();\}$

$S_1()$ goes in an infinite loop

A left recursive grammar has a non-terminal $S \rightarrow^+ S d$ for some d

$$S \rightarrow Sd | \epsilon$$

Pearson parsing wants to see the first part of the input first.
In left recursions we produce the string right to left.

We can fix it by using right recursion

$$S \rightarrow \beta S'$$

$$S' \rightarrow d S' | \epsilon$$

$$S \rightarrow Sd_1 | \dots | Sd_n | \beta_1 | \dots | \beta_m$$

$$S \rightarrow \beta_1 S' | \dots | \beta_m S'$$

$$S' \rightarrow d_1 S' | d_n S' | \epsilon$$

$$S \rightarrow A d_1 S$$

$$A \rightarrow S \beta$$

still left recursive grammar!

Bottom-up Parsing

Predictive parsing

Like recursive descent but parser can "predict" which production to use

- By looking at the next few tokens

- No backtracking

lookahead

Predictive parsers accept LL(1) grammars

left to right
head input

left-most
derivation

↑ tokens
of lookahead
(usually 1)

In LL(1), at each step, only one production is guaranteed to be correct.

$E \rightarrow T + E | T$ - T has two productions starting with int
 $T \rightarrow \text{int} | \text{int} * T$ - E " with T .
 $| (E)$

This grammar cannot be predicted. We need to left-factor the grammar

$$E \rightarrow T + E | T$$

$$T \rightarrow \text{int} | \text{int} * T$$

$$| (E)$$

$$E \rightarrow T x$$

$$x \rightarrow + E | \epsilon$$

$$T \rightarrow \text{int} y | (E)$$

$$y \rightarrow * T | \epsilon$$

Method similar to recursive descent

- For the leftmost non-terminal S
- We look at the next input Token a
- Look up in the $LL(1)$ table the production to use for $[S, a]$

A stack of records to record the frontier

- Non-terminals that have yet to be expanded
- Terminals that have yet to be matched against input
- Leftmost pending terminal or non-terminal is at the top of the stack

reject on reaching error state

accept on end of input and empty stack.

stack bottom

initialise stack = $\langle \$\$ \rangle$
repeat

curr stack of
 $\langle x, \text{last} \rangle$:
 if $T[x, \text{next}] = y_1 \dots y_n$
 then stack $\leftarrow y_1 \dots y_n \text{ last}$
 $\langle t, \text{last} \rangle$:
 if $T[z, \text{next}] = y$
 then stack $\leftarrow \langle y, \text{last} \rangle$
 else error();
 until stack == $\langle \rangle$

First sets

consider $A \rightarrow \alpha$, and token t

$T[A, t] = d$ in two cases:

if $d \rightarrow^* t\beta$

- d can derive t in the first production
- we say $t \in \text{FIRST}(d)$

if $A \rightarrow d$ and $d \rightarrow^* \epsilon$ and $S \rightarrow^* \beta A t \delta$

- Useful if stack has A , input is t , and A cannot derive t
- In this case the only option is to get of A (by deriving ϵ)
Can work only if t can follow A in at least one derivation
- we say $t \in \text{follow}(A)$
-

$$\text{FIRST}(x) = \{t \mid x \rightarrow^* t\alpha\} \cup \{\epsilon \mid x \rightarrow^* \epsilon\}$$

Algorithm

1 $\text{first}(t) = \{t\}$

2 $\epsilon \in \text{first}(x)$

- if $x \rightarrow \epsilon$

- if $x \rightarrow A_1 \dots A_n$ and
 $\epsilon \in \text{first}(A_j)$ for $1 \leq j \leq n$

3 $\text{first}(d) \subseteq \text{first}(x)$ if $x \rightarrow A_1 \dots A_n$ and
- and $\epsilon \in \text{first}(A_j)$ for $1 \leq j \leq n$

$$E \rightarrow TX \quad X \rightarrow +E \mid \epsilon \quad T \rightarrow (E) \mid \text{int} \quad Y \rightarrow *T \mid \epsilon$$

$$\text{FIRST}(E) \supseteq \text{FIRST}(T)$$

$$\text{FIRST}(T) = \{(\), \text{int}\}$$

$$\text{FIRST}(E) = \{(\), \text{int}\}$$

$$\text{FIRST}(X) = \{+\}, \{\epsilon\}$$

$$\text{FIRST}(Y) = \{*, \epsilon\}$$

$$\text{FIRST}(+) = \{+\}$$

$$\text{FIRST}(\#) = \{\#\}$$

$$\vdots \quad \{(\}$$

$$\{ \}$$

$$\{ \text{int} \}$$

Follow sets

$$\text{follow}(x) = \{t \mid S \xrightarrow{*} \beta x t \delta\}$$

t is $\text{follow}(x)$ if there is some derivation where t appears right after x

If $X \rightarrow AB$ then

$$\text{FIRST}(B) \subseteq \text{follow}(A) \quad X \rightarrow AB \xrightarrow{*} AtB$$

$$\text{follow}(x) \subseteq \text{follow}(B) \quad S \rightarrow \dots xt \dots \rightarrow ABt$$

if $B \xrightarrow{*} \epsilon$ then

$$\text{follow}(x) \subseteq \text{follow}(A) \quad S \rightarrow \dots xt \dots \rightarrow ABt \xrightarrow{*} At$$

If S is the start symbol then $\$ \in \text{follow}(S)$

end of the
input symbol

Algorithm

1. $\$ \in \text{follow}(S)$

2. $\text{FIRST}(\beta) - \{\epsilon\} \subseteq \text{follow}(x)$

- For each production $A \rightarrow \alpha x \beta$

3. $\text{follow}(A) \subseteq \text{follow}(x)$

- For each production $A \rightarrow \alpha x \beta$ where
 $\epsilon \in \text{FIRST}(\beta)$

by definition

$$\text{follow}(\$) = \{\$\}$$

$$\text{follow}(X) = \{\$\}$$

$$\text{follow}(T) = \{+, \$\}$$

$$\text{follow}(Y) = \{*, \$\}$$

$$\text{follow}('(') = \{(\}, \text{int}\}$$

$$\text{follow}(')') = \{+, \$, '\}$$

$$\text{follow}('+') = \{(\}, \text{int}\}$$

$$\text{follow}('*') = \{(\}, \text{int}\}$$

$$\text{follow}(\text{int}) = \{\#, +, \$, '\}$$

$$\begin{array}{ll} E \rightarrow TX & X \rightarrow +E \mid \epsilon \\ T \rightarrow (E) \mid \text{int} & Y \rightarrow *T \mid \epsilon \end{array}$$

equel $\left\{ \begin{array}{l} \text{follow}(X) \subseteq \text{follow}(E) \text{ because } X \rightarrow +E \\ \text{follow}(E) \subseteq \text{follow}(X) \text{ because } E \rightarrow TX \\ \text{follow}(E) \subseteq \text{follow}(T) \quad E \rightarrow TX \rightarrow T \end{array} \right.$

equel $\left\{ \begin{array}{l} \text{follow}(Y) \subseteq \text{follow}(T) \quad Y \rightarrow *T \\ \text{follow}(T) \subseteq \text{follow}(Y) \quad T \rightarrow \text{int} Y \end{array} \right.$

LL1 Parsing Table

Construct parsing table T for CFG G

For each production $A \rightarrow \alpha$ in G do

- For each terminal $t \in \text{First}(\alpha)$ do

$$T[A, t] = \alpha$$

- If $\epsilon \in \text{First}(\alpha)$, for each $t \in \text{Follow}(A)$ do

$$T[A, t] = \alpha$$

- If $\epsilon \in \text{First}(\alpha)$ and $\$ \in \text{Follow}(A)$ do

$$T[A, \$] = \alpha$$

$E \rightarrow TX$ $X \rightarrow +E | \epsilon$ we need to know $\text{Follow}(X)$
 $T \rightarrow (E) | \text{int} Y$ $Y \rightarrow *T | \epsilon$ in order to know when to use $X \rightarrow \epsilon$

	(+	*	int	\$
E	TX			TX	
T	(E)			int Y	
X		$\epsilon + E$			ϵ
Y			$*T$		

blank entries
are parsing errors

$$S \rightarrow Sa | b$$

$$\text{first}(S) = \{b\}$$

$$\text{follow}(S) = \{\$, a\}$$

$$a \quad b \quad \$$$

S	b	\$
---	---	----

multiply defined!

If any entry is multiply defined, then G is not LL1

- not left factored

- left recursive

- ambiguous

- other grammars are not LL1

Most programming language CFGs are not LL1

Bottom-up parsing

Bottom up parsing is more general than top-down parsing

- Just as efficient
- Builds on ideas in top-down parsing

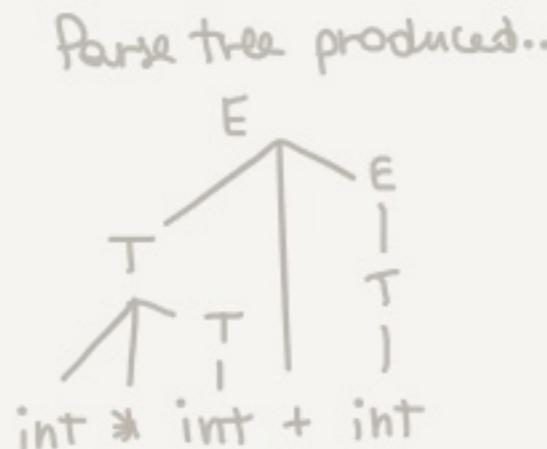
Bottom-up is the preferred method!

- For our example, we no longer need left factored grammars

$$\begin{aligned} E &\rightarrow T + E \mid T \\ T &\rightarrow \text{int}^* T \mid \text{int} \mid (E) \end{aligned}$$

Bottom-up parsing reduces a string to the start symbol by inverting productions

input string
↓
 $\text{int}^* \text{ int} + \text{int}$
 $\text{int}^* T + \text{int}$
 $T + \text{int}$
 $T + T$
 $T + E$
 $E \curvearrowleft \text{start symbol}$



A bottom-up parser traces a rightmost derivation in reverse

IMPORTANT FACT #1

Shift Reduce parsing

- Let $d\beta w$ be a step of a bottom-up parse
- Assume the next reduction is by $X \rightarrow \beta$
- Then w is a string of terminals

$dXw \rightarrow d\beta w$ is a step in a right-most derivation

Split string into two substrings

- Right substring is the unexamined input
 - Left substring has terminals and non-terminals
 - Split marked by a |
- —

Bottom-up parsing uses only two kinds of actions. Those are shift and reduce

- Shift move 1 one place to the right
- $$ABC|xyz \Rightarrow ABCx|yz$$

- Reduce is to apply on inverse production of the right end of the left string

$$\text{if } A \rightarrow xy \\ Cbxy|ijk \Rightarrow CbA|ijk$$

— —

$\text{int}^* \text{ int} + \text{int}$ shift
 $\text{int}^* \text{ int} + \text{int}$...
 $\text{int}^* \text{ int} + \text{int}$...
 $\text{int}^* T + \text{int}$ reduce
...

...
 $T + \text{int}$ reduce
 $T + \text{int}$ shift
 $T + \text{int}$...
 $T + E$ reduce
 $T + T$ reduce

• • •

Left string can be implemented by a stack

- Top of the stack is 1

Shift pushes a terminal onto the stack

Reduce pops symbols from top of the stack and pushes one terminal on the stack

— =

If it is legal to both shift and reduce, we say there's a shift-reduce conflict

If it is legal to reduce by two productions, then there's a reduce-reduce conflict

^TBAD

Handlers

How to decide when to shift/reduce?

$$E \rightarrow T + E \mid T \quad T \rightarrow \text{int}^* \mid \text{int} \mid (\text{E})$$

...
int | * int + int

we could decide to reduce,
that would be a mistake!

We want to reduce only if the result can still be reduced to the start symbol

$$\text{Assume } S \xrightarrow{*} dXw \xrightarrow{\text{reduction}} d\beta w$$

Then $d\beta$ is a handle of $d\beta w$.
A handle is a reduction that also allows further reductions back to the start symbol

We only want to reduce handles!

IMPORTANT FACT #2

In shift-reduce parsing handles appear only at the top of the stack!

Recognising handles

There are no efficient algorithms to recognise handles. But there are good heuristics for guessing them and for some CFG they always guess correctly

— =

d is a viable prefix if there is an w ^{rest of input} such that $d|w$ is a valid configuration _{stack} of a shift-reduce parser

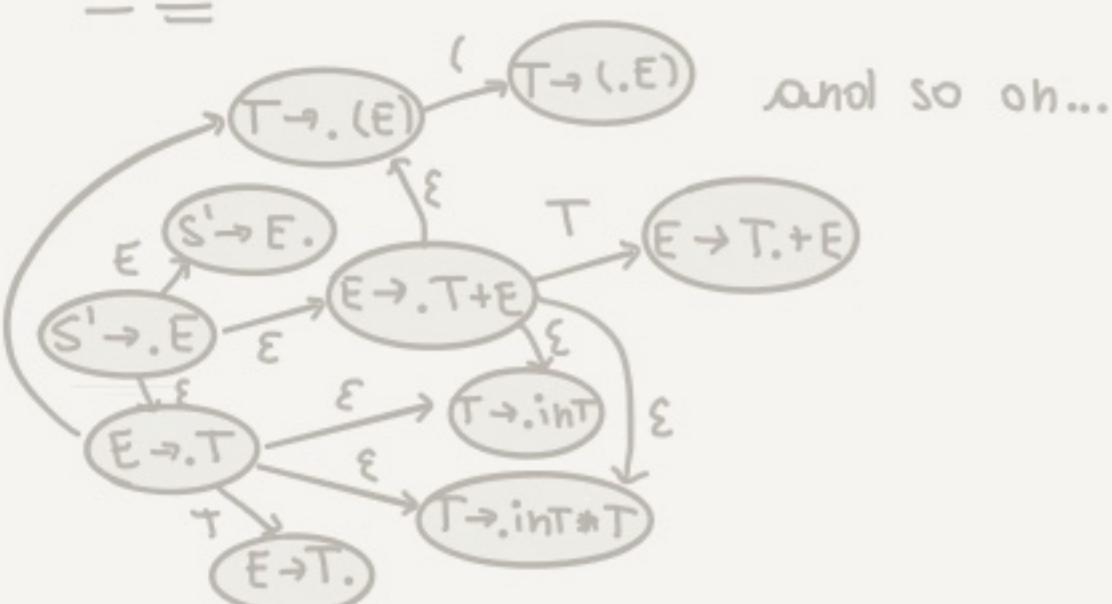
- A viable prefix does not extend past the right end of a handle
- It is viable because it is a prefix of the handle
- As long as the parser has viable prefixes on the stack, no error has been encountered

For any grammar, the set of viable prefixes is a regular language

IMPORTANT FACT #3

Recognizing viable prefixes

- 1 Add a dummy production $S' \rightarrow S$ to G
- 2 The NFA states are the items of G
 - The NFA reads the stack and it is either going to say 'yes' it's valid, or no.
- 3 For item $E \rightarrow \alpha \cdot X \beta$ add transition
 $E \rightarrow \alpha X \cdot \beta \xrightarrow{X} E \rightarrow \alpha \cdot X \beta$
seen α on stack
- 4 For item $E \rightarrow \alpha \cdot X \beta$ and production $X \rightarrow \gamma$ add
 $E \rightarrow \alpha \cdot X \beta \xrightarrow{\gamma} X \rightarrow \gamma$
non-terminal
- 5 Every state is an accepting state
- 6 Start state is $S \rightarrow S'$



Valid items

Using the standard subset of states construction, we can build a DFA that is equivalent to the previous NFA

Each state is a set of items

These states are called canonical collection of LR(0) items

— =

Item $X \rightarrow \beta \cdot \gamma$ is valid for a viable prefix $\alpha \beta$ if
 $S' \xrightarrow{*} \alpha X w \rightarrow \alpha \beta \gamma w$
 by a rightmost derivation

After parsing $\alpha \beta$, the valid items are the possible tops of the stack items

SLR parsing

LR(0) parsing

- stack contains α
- next input is t
- DFA on input t terminates in state s

Reduce by $X \rightarrow \beta$ if

- s contains item $X \rightarrow \beta$.

Shift if

- s contains item $X \rightarrow \beta \cdot t w$
- equivalent to saying s has a transition labeled t



• • •
LR(0) has a reduce-reduce conflict if
- Any state of the DFA has two reductions

LR(0) has a shift-reduce conflict if
- Any state has a reduce item and a shift item

—
SLR stands for simple LR

SLR improves LR(0) by adding some heuristics that will refine when we shift over when we reduce, so that fewer states have conflicts

We add one condition to the reduction case in LR(0)

...
- $t \in \text{follow}(x) \leftarrow$ additional rule
—

If there are conflicts under these rules, then the grammar is not SLR

These rules amount to a heuristic for detecting the handles

Algorithm
—

1 Let M be a DFA for viable prefixes of G

2 Let $|x_1 \dots x_n \$|$ be initial configuration

3 Repeat until configuration is $S | \$$

- Let $\alpha | w$ be current configuration
- Run M on current stack α
- If M accepts α with items I , let ' a ' be the input
 - Shift if $x \rightarrow \beta, a \gamma \in I$
 - Reduce if $x \rightarrow \beta, \epsilon \in I$ and $a \in \text{follow}(x)$
 - Report parsing error if neither applies

In $SLR(k)$, k is the amount of lookahead, in practice $k=1$

SLR improvements

Running the viable prefix automaton at the top of the stack is wasteful

- Most work is repeated

To avoid repeating, we want to remember the state of the automaton at each prefix

Change stack to contain pairs $(\text{symbol}, \text{state})$

The state is the result of running the DFA on all symbols to its left

At the bottom of the stack we will have $(\text{only}, \text{start})$

- 'only' is only dummy symbol
- 'start' is the start state

—

Define a table goto that maps a symbol to another state

This is just the transition functions of the DFA written as an array

—

Our algorithm moves are

- Shift; push (a, x) on the stack
 - a is the current input
 - x is a DFA state
- Reduce $x \rightarrow \alpha$ (just as before)
- Accept and Error moves

—

Define an action table which tells us which move to make in every possible state

The action table is indexed by the state and the next input symbol. It returns one of the moves above

• • •

• • •

For each state s_i and terminal a

- If s_i has item $X \rightarrow d, a \beta$ and $\text{goto}[i, a] = j$
then $\text{action}[i, a] = \text{shift } j$
 - If s_i has item $X \rightarrow d$ and $a \in \text{low}(X)$ and
 $X \neq S'$ then $\text{action}[i, a] = \text{red } \perp X \rightarrow d$
 - If s_i has item $S' \rightarrow S$, then $\text{action}[i, \$] = \text{accept}$
 - Otherwise $\text{action}[i, a] = \text{error}$
- =

Full SLR algorithm

Let $I = w\$$ be the initial input

Let $j = 0$

Let DFA state 1 have item $S' \rightarrow .S$

Let $\text{stack} = \langle \text{dummy}, 1 \rangle$

repeat

```

    case action[top-state(stack), I[j]] of
        shift k: push < I[j+1], k >
        reduce  $X \rightarrow A$ :
            pop | A | pairs,
            push <  $X, \text{goto}[\text{top-state}(\text{stack}), A]$  >
        accept: halt normally
        error: halt and report error
    endcase
endrepeat

```

Note that this algorithm only uses the DFA states and the input. Stack symbols are never used

We still need stack symbols for later compiler stages

— =

Widely used bottom-up parsing algorithms are based on a more powerful class of grammars, the LR(1) where

- Lookahead is built into items
- An item is a pair LR(0) item and lookahead
- More accurate than just using follow sets

Semantic Analysis

Lexical analysis

- Detects inputs with illegal tokens

Parsing

- Detects inputs with ill-formed parse trees

Semantic analysis

- Catches all remaining errors

— =
Checks of many kinds... code checks:

- All identifiers are declared
- Types
- Inheritance and relationships
- Reserved identifiers not misused
- Classes and methods defined only once
- Others...

The requirements depend on the language
Scope

Matching identifier declarations with uses

let $y: \text{String} \leftarrow \text{"abc"}$ in $y+3$
 y cannot be added?

let $y: \text{Int}$ in $x+3 \Leftarrow$ no definition
of x

The scope of an identifier is the portion of a program in which the identifier is accessible

The same identifier may refer to different things in different parts of the program

- Different scopes for the same name can't overlap

An identifier can have restricted scope

Most programming languages today have static scope

- The scope of a variable depends only on the program text, not run-time behaviour

A Few languages are dynamically scoped

Like LISP, SNOBOL

- Scope of a variable depends on the execution behaviour of the program

```
let x:Int <- 0 in  
  {  
    same scope  
    let x:int <- 1 in  
      x;  
    }  
  --=
```

A variable binds to the definition that is most closely enclosing it of the same name

In cool...

class definitions:

- cannot be nested
- are globally visible

A class name can be used before it is defined

Attribute names are global within the class in which they are defined

Methods need not be defined in the class in which it is used, but in some parent class

A method can be overridden

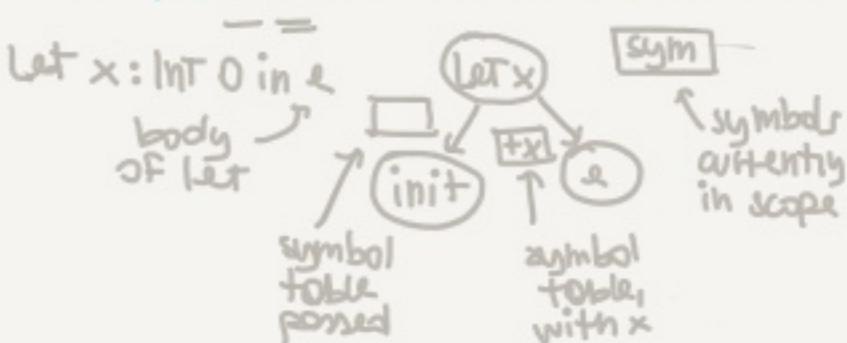
Symbol table

Much of semantic analysis can be seen as a recursive descent of an abstract syntax tree

Before; Process on AST node h
Recuter; Process the children of h
After; Finish processing node h

Recursive ✓

When performing semantic analysis on a portion of the AST, we need to know which identifiers are defined



Before processing λ , add x to symbol table overriding any existing definition

Recursive

After processing λ , remove x from table and restore old definition of x

Symbol table tracks the current binding of identifiers

For a simple symbol table we can use a stack

- `add-symbol(x)`, push x and associated info
- `find-symbol(x)`, return first x occurrence, or NULL
- `remove-symbol()`, pop

When using a stack, declarations are perfectly nested

Consider `f(x:Int, x:Int) {}`

that is not legal, functions introduce multiple names at once and the walk will not detect any error

We change stack operations slightly

- `enter-scope()`, start a new nested scope
- `exit-scope()`, exit current scope
- `check-scope(x)`, true if x defined in current scope
- `find-symbol(x)`, find current x or NULL
- `add-symbol(x)`, add x to the table

This is a stack of scopes!

Class names can be used before being defined

- Pass through code before and gather all class names
- Pass again and do the check

Semantic analysis requires multiple passes

Types

Notion varies from language to language

A set of values

A set of operations on those values

\equiv

Doesn't make sense to add pointer and integer

Does make sense to add int with int

Both have the same assembly implementation

A language's type system specifies which operations are valid for which types

Type checking ensures operations are used only with the correct types

\equiv

Statically typed, all or almost all checking of types is done as part of compilation

Dynamically typed, almost all checking of types is done as part of program execution

Untyped, No type checking

Static typing

- Static checking catches many program errors at compile time
- Avoids overhead of runtime type checks

Dynamic typing

- Static types are restrictive
- Rapid prototyping difficult within a static type system

cool types are

- class names
- SELF-TYPE

The user declares types for identifiers

The compiler infers types for expressions

\equiv

Type checking is the process of checking a fully typed program

Type inference is the process of filling in missing types

Those terms are different, but they're usually used interchangeably

Type checking

We've seen two examples of formal notation for specifying parts of a compiler. RE's and CFB's

The formalism for type checking is logical rules and rules of inference

\equiv

- Inference rules have the form:
IF Hypothesis is true, then Conclusion is true

- Type checking computes by reasoning
IF E_1 and E_2 have certain types, then E_3 has certain type

- Rules of inference are a compact notation for "if-then" statements

Start with simple symbols and gradually add features

Building blocks

- symbol \wedge is "and"
- symbol \Rightarrow is "if-then"
- $x:T$ is " x has type T "
 \equiv

IF e_1 has type int and e_2 has type int, then $e_1 + e_2$ has type int

$(e_1 : \text{Int} \wedge e_2 : \text{Int}) \Rightarrow e_1 + e_2 : \text{Int}$

\equiv

The statement above is a special case of an inference rule

\downarrow
hypothesis $1 \dots$ hypothesis $n \Rightarrow$ conclusion

By tradition, inference rules are written $\frac{\text{hypothesis}_1 \dots \text{hypothesis}_n}{\text{conclusion}}$

called turnstiles, means
"It is provable that..."

$\frac{i \text{ is integer literal}}{+ i : \text{Int}}$ [Int]

$\frac{+ e_1 : \text{Int} \quad + e_2 : \text{Int}}{+ e_1 + e_2 : \text{Int}}$ [Add]

Subtyping

Define a relation \leq on classes

- $X \leq X$
- $X \leq Y$ if X inherits from Y
- $X \leq Z$ if $X \leq Y$ and $Y \leq Z$

$$O(x) = T_0$$

$$O[e_1] = T_1$$

$$T_1 \leq T_0$$

$$\frac{}{O[x \leftarrow e_1 : T_1]} \text{ [Assign]}$$

\equiv

if e_0 then e_1 else e_2

We don't know which branch is taken.
The best we can do is the smallest
Supertype larger than type of e_1 or e_2

$\text{lub}(X, Y)$, the least upperbound of X and
 Y is Z if

- $X \leq Z \wedge Y \leq Z$ (Z is an upper bound)
- $X \leq Z' \wedge Y \leq Z' \Rightarrow Z \leq Z'$
(Z is least among upper bounds)

The least upper bound is the least common
ancestor in the inheritance tree

$$O[e_0] = \text{Bool}$$

$$O[e_1] = T_1$$

$$O[e_2] = T_2 \quad [\text{IF-Then-Else}]$$

$$\frac{}{O[\text{if } e_0 \text{ then } e_1 \text{ else } e_2]} = \text{lub}(T_1, T_2)$$

A method 'foo' and an object 'foo'
can coexist in the same namespace

In the type rules, this is reflected
by a separate mapping ' M ' for
method signatures

$$M(C, F) = (T_1, \dots, T_h, T_{h+1})^{\text{Result}}$$

Means in class C there's
a method F

$$F(x_1 : T_1, \dots, x_n : T_h) : T_{h+1}$$

\equiv

Add method environment to
all rules

$$\frac{}{O, M[e_1] = \text{Int} \quad O, M[e_2] = \text{Int}} \text{ [Add]}$$

$$\frac{}{O, M[e_1 + e_2] = \text{Int}} \text{ [Add]}$$

The full type environment for COOL

- A mapping O giving types to
object id's
- A mapping M giving types to methods
- The current class C

The full form of a sentence in
the cool sentence logic is

$$O, M, C \vdash e : T$$

Implementing type checking

COOL type checking can be implemented
in a single traversal over the AST

Top-down phase, type environment is
passed down the tree

bottom-up phase, types are passed back up

$$\frac{}{O, M, C \vdash e_1 : \text{Int} \quad O, M, C \vdash e_2 : \text{Int}} \text{ [Add]}$$

$$\frac{}{O, M, C \vdash e_1 + e_2 : \text{Int}} \text{ [Add]}$$

$$\frac{}{\text{TypeCheck}(\text{Environment}, e_1 + e_2) = \{ O, M, C \}}$$

$$T_1 = \text{TypeCheck}(\text{Environment}, e_1);$$

$$T_2 = \text{TypeCheck}(\text{Environment}, e_2);$$

$$\text{Check } T_1 == T_2 == \text{Int};$$

$$\text{return Int;} // \text{check succeeds}$$

\equiv

$$O, M, C \vdash e_0 : T_0$$

$$\frac{O[T/x], M, C \vdash e_1 : T_1}{T_0 \leq T_1} \quad [\text{LET-Init}]$$

$$\frac{}{O, M, C \vdash \text{let } x : T \leftarrow e_0 \text{ in } e_1 : T_1}$$

$$\frac{}{\text{TypeCheck}(\text{Environment}, \text{let } x : T \leftarrow e_0 \text{ in } e_1) = \{ T_0 = \text{TypeCheck}(\text{Environment}, e_0); T_1 = \text{TypeCheck}(\text{Environment}, \text{add}(x : T), e_1); \}}$$

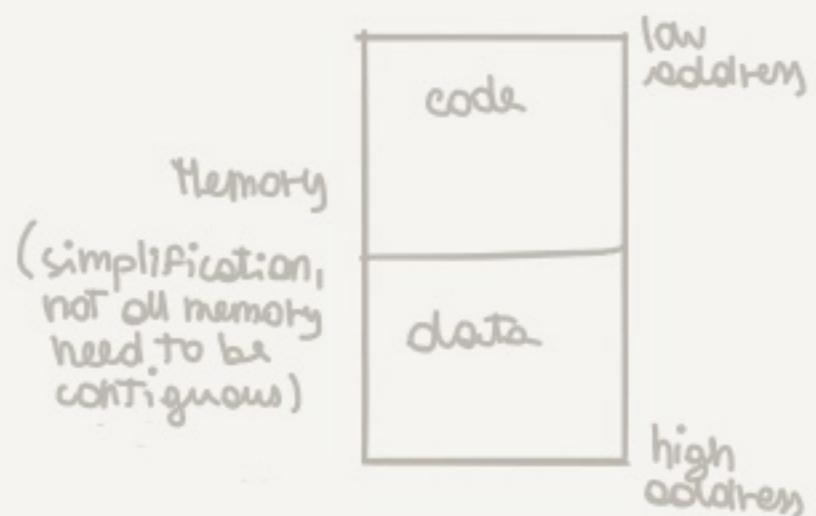
$$\frac{}{T_1 = \text{TypeCheck}(\text{Environment}, \text{add}(x : T), e_1); \text{Check subtype}(T_0, T_1); \text{return } T_1; }$$

Runtime organization

Before discussing code generation, we need to understand what we're trying to generate

When a program is invoked

- The OS allocates space for the program
- The code is loaded into that space
- The OS jumps to the entry point (ie main)



Compiler is responsible for generating code, but also has to decide where the layout of the data is going to be and then generate code that correctly manipulates that data

Activation

Complications in code generation come from trying to be fast as well as correct

Assumptions

- Execution is sequential (...in a well defined order)
- After a procedure is called, control returns to the point immediately after

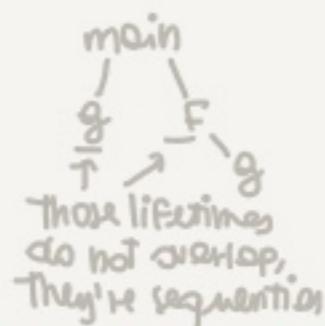
We'll talk on invocation of procedure P, on activation of P

The lifetime of P is all the steps in P and all the steps of the procedures that P calls

The lifetime of a variable x is the portion of the execution in which x is defined

When P calls Q, Q returns before P
Activation lifetimes can be depicted as a tree

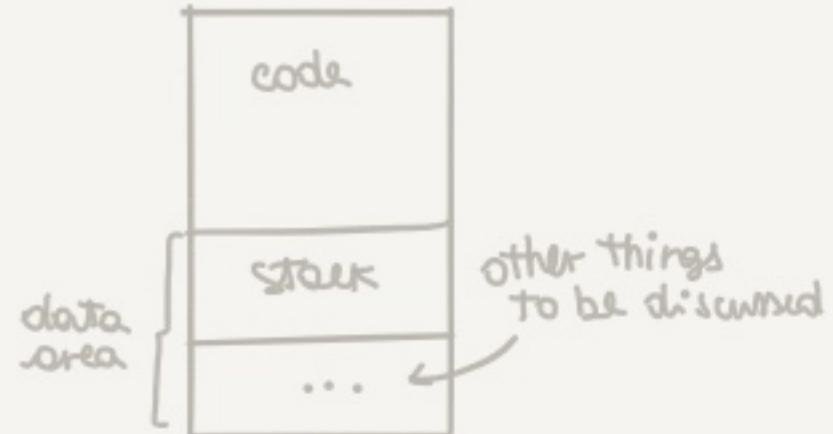
```
f(): Int { g();  
main(): Int  
{ g(); f(); };  
...}
```



Since activations are properly nested, we can use a stack to implement the currently active activations.

For the above example
The stack will at any time, contain the active procedures

push main
push g
pop g
push F
push ...
push g
pop g
push F
push main
pop F
push ...



Activation records (also called frames)

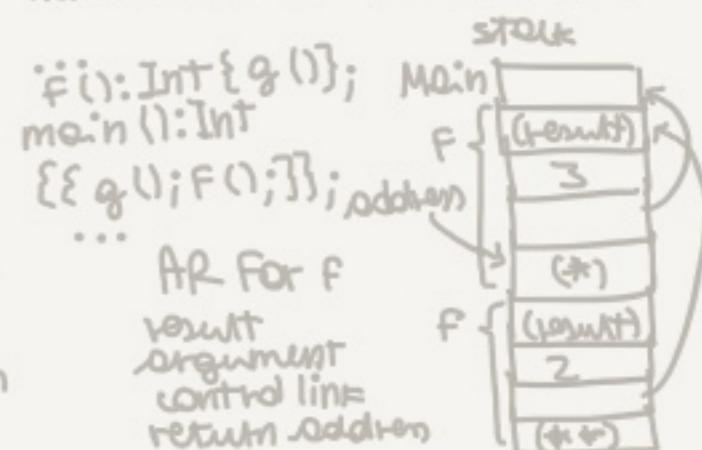
An activation record is all the information that's needed to manage the execution of one procedure activation

If procedure F calls procedure G, then G's activation record contains a mix of info about F and G

F calls G

F is suspended until G completes, at which point F resumes

G's activation record contains information needed to complete the execution of G and resume F



Compilers hold as much of the frame as possible in registers

- Especially method result and arguments
- =

The compiler has to determine at compile-time, the layout of activation records and generate code that correctly encodes locations in the activation records

AR layout and code generator must be designed together

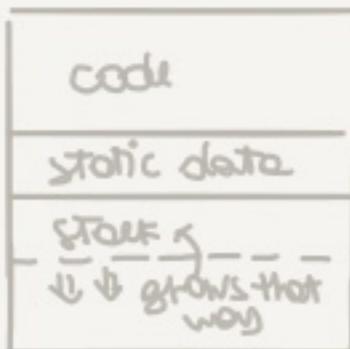
Globals and heaps

All references to a global variable point to the same object

- We can't store a global in an activation record

Globals are assigned a fixed address once

- Said to be "statically allocated" (at compile time)



A value that survives the procedure that creates it cannot be kept in the AR

method foo() { return new Bar; }

'Bar' must survive the deallocation of Foo's AR

Language with dynamically allocated data use a heap to store dynamic data

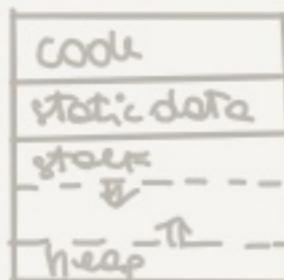
Different kinds of data that a language implementation has to deal with

- The code area; fixed size and read only (in many languages)
 - The static area contains data with fixed addresses (global data); fixed size, may be readable or writable
 - Stack contains an AR for each currently active procedure; each AR has fixed size and contains all the local information needed to execute a particular activation
 - Heap contains all other data; In C, heap is managed by malloc and free
- =

Both heap and stack can grow!

We don't want them to grow into each other

Solution, start heap and stack at opposite ends of memory and let them grow into each other



If these pointers ever become equal, the program is out of memory and the runtime system will take some sort of action

Alignment

Most modern machines are 32 or 64 bits

- 4 or 8 bytes in a word
 - Machines are either byte or word addressable
- =

Data is word aligned if it begins at a word boundary

Most machines have some alignment restrictions

- Or performance penalties for poor alignment
- = unused]

"Hello"

He llo w

text data item goes to the next word

Stack machine

Only storage is a stack

$$r = F(a_1, \dots, a_n):$$

- pops n operands
- Computes F using the operands
- push the result + on the stack

Contents prior to r are preserved

Addresses are not explicitly stated

- Always refer to the top of the stack

More compact programs

(...then register machine)

Java bytecode uses stack evaluation

—

n-register stack machine

- keep to n stack locations in registers

1-register stack machine

- Also called accumulator

Takes ADD, in a stack machine

- Pop two arguments
- Add them
- Put result back

In a 1-register stack machine

$$\rightarrow acc \leftarrow acc + top_of_stack$$

accumulator register

—

Consider an expression $op(e_1, \dots, e_n)$

- e_1, \dots, e_n are subexpressions

For each $e_i (0 \leq i < n)$

- compute $e_i \leftarrow$ result in acc
- push result

$e_n \rightarrow$ just evaluate (don't push result)

Pop n-1 values from the stack, compute op

Store result in accumulator

After evaluating e , the accumulator holds the value of e and the stack is unchanged

—

$3 + (7 + 5)$

code	acc	stack
acc ← 3	3	linit
push acc	3	3, linit
acc ← 7	7	3, linit
push acc	7	7, 3, linit
acc ← 5	5	7, 3, linit
acc ← acc +	12	7, 3, linit
...top_of_stack		
pop	12	3, linit
acc ← acc +		
...top_of_stack	15	3, linit
pop	15	linit

Code generation

We focus on generating code for a stack machine with accumulator
(simplest code gen strategy)

We use the MIPS processor to run the resulting code on a real machine

Simulate stack machine instructions using MIPS (simulator)

The accumulator is kept in MIPS register \$a0

The stack is kept in memory

- Grows towards lower addresses

The address of the next stack location is kept in MIPS register \$sp (stack pointer)

- Top of the stack is \$sp+4

MIPS architecture

- Prototypical RISC
- Most operations use registers for operands and results

—

lw reg1 offset(reg2)

- Load 32 bit word from reg2+offset into reg1
addi reg1 reg2 reg3; reg1 ← reg2 + reg3

sw reg1 offset(reg2)

- Store 32 bit word in reg1 at address reg2+offset
addiu reg1 reg2 imm; reg1 ← reg2 + imm

li reg imm; reg ← imm

—

acc ← 7
push acc

acc ← 5

acc ← acc +

top_of_stack

pop

li \$a0 7
sw \$a0 0(\$sp)

addiu \$sp \$sp -4

Temporary Register

li \$a0 5

lw \$t1 4(\$sp)

add \$a0 \$a0 \$t1

addiu \$sp \$sp 4

Code generation I

Language with integers and operations

$P \rightarrow D; P1 D$

$D \rightarrow \text{def id (ARGS)} = E;$

$\text{ARGS} \rightarrow \text{id}, \text{ARGS} | \text{id}$

$E \rightarrow \text{int} | \text{id} | \text{if } E_1 = E_2 \text{ then } E_3 \text{ else } E_4$
 $| E_1 + E_2 | E_1 - E_2 | \text{id} (E_1, \dots, E_n)$

- The first function definition f is the entry point

—

For each expression e we generate MIPS code that

- Compute e and store result in $\$a0$
- Preserve $\$sp$ and stack contents

cgen(e) generates code for e

—

To evaluate an expression, which is just a constant, load that constant into the accumulator

cgen(i) = li $\$a0, i$; compile time

buftime

cgen($e_1 + e_2$) =
cgen(e_1)

print sw $\$a0, 0(\$sp)$ put e_1 on the stack
... addiu $\$sp, \$sp - 4$ bump the stack pointer
cgen(e_2)
... lw $\$t_1, 4(\$sp)$ load e_1 into t_1 from stack
... addiu $\$a0, \$t_1, \$a0$
... addiu $\$sp, \$sp, 4$ pop stack

Assembly is printed into a file

Put the result of e_2 directly in t_1

cgen($e_1 + e_2$) =
cgen(e_1) WRONG
move $\$t_1, \$a0$ Values are overridden
cgen(e_2)
add $\$a0, \$t_1, \$a0$ in t_1

—

The code for $+$ is a template with holes for evaluating e_1 and e_2

stack machine code gen is recursive

—

bne $\$reg_1, \$reg_2, label$
- Branch if $\$reg_1 = \reg_2
b $label$
- Jump to label

cgen (if $e_1 = e_2$ then e_3 else e_4) =
cgen(e_1)
addiu $\$sp, \$sp - 4$
cgen(e_2)
lw $\$t_1, 4(\$sp)$
addiu $\$sp, \$sp, 4$
beq $\$a0, \$t_1, true_branch$

false_branch:

cqz (e_4)

b end_if

true_branch:

cgen(e_3)

end_if:

Code generation II

Code for function calls and function definitions depends on the layout of the AR

A very simple AR suffices for this language

- The result is always in the accumulator

No need to store result in AR

- AR holds actual parameters

For $f(x_1, \dots, x_n)$ push x_n, \dots, x_1

These are the only variables in this language

—

Stack pointer is preserved across function calls

- No need for a control link in AR

We need return address in AR

A pointer to the current activation is useful

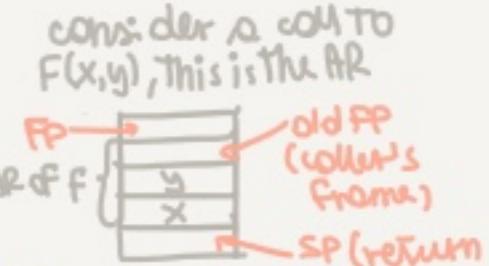
- This pointer is in register $\$fp$ (frame pointer)

—

So we need an AR with

- frame pointer
- actual parameters
- return address

—



The calling sequence is the sequence of instructions to set up a function invocation

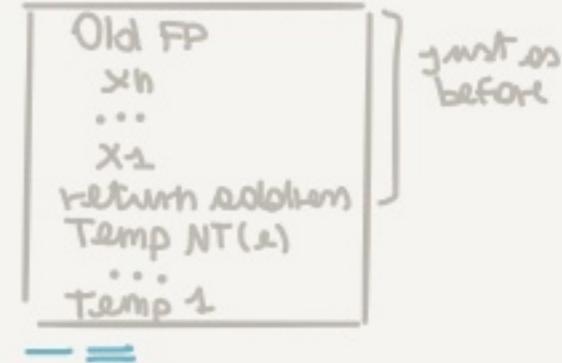
—

Jal label; jump to label and save address of next instruction in $\$ra$ (return address)



For a function definition
 $f(x_1 \dots x_n) = \dots$, the AR is $2 + n + NT(\ell)$
 elements

- Return addresses
- Frame pointer
- n arguments
- $NT(\ell)$ locations for intermediate results



- Solution, we use the frame pointer
- Always points to the return address of the stack
 - Since it doesn't move, it can be used to find the variables

Let x_i be the i th ($i=1 \dots n$) argument of the function

$$\text{cgen}(x) = \text{lw } \$a0 \ 4*i(\$fp)$$

Production compilers do different things

- Emphasize on keeping values on registers
- Intermediate results are laid out in the AR, not pushed and popped from stack

Temporaries

Keep temporaries in the AR
 IDEA: ... rather than pushing and popping

The code generation must assign a fixed location in the AR for each temporary

Let $NT(\ell)$ be the # of temporaries needed to evaluate ℓ

$$NT(\ell_1 + \ell_2)$$

- Needs at least as many temporaries as $NT(\ell_1)$
- Needs at least as many temporaries as $NT(\ell_2) + 1$

$$NT(\ell_1 + \ell_2) = \max(NT(\ell_1), 1 + NT(\ell_2))$$

hold the result of ℓ_1

$$NT(\text{if } \ell_1 = \ell_2 \text{ then } s_3 \text{ else } s_4) = \max(NT(\ell_1), 1 + NT(\ell_2), NT(\ell_3), NT(\ell_4))$$

$$NT(\text{id}(\ell_1, \dots, \ell_n)) = \max(NT(\ell_1), \dots, NT(\ell_n))$$

no additional ones for result because result is saved in the AR

$$NT(\text{int}) = 0 \quad NT(\text{id}) = 0$$

variable reference

caller side

```

cgen(f(e1...en)) =
  sw $fp 0($sp) save fp in stack
  addiu $sp $sp-4 (start building AR)
  cgen(en)
  sw $a0 0($sp) push en result
  addiu $sp $sp-4
  ... for all arguments
  cgen(e1)
  sw $a0 0($sp)
  addiu $sp $sp-4
  jol f-entry
  jr reg; jump to address in reg
  
```

The AR so far
 is $4 * n + 4$ bytes
 from frame pointer

cgen(def(x1...xn)=e) =

F-entry:

```

move $fp $sp
sw $ra 0($sp)
addiu $sp $sp-4
cgen(e)
lw $ra 4($sp)
addiu $sp $sp 4*n+8
lw $fp 0($sp)
jr $ra
  
```

At this point, the AR
 is completely set up.
 Next we generate
 code for the function body

with n arguments
 + return address
 + old fp

Variable references are the last constraint

The variables of a function are just its parameters
 - They're all in the AR

But the stack grows when intermediate results
 are saved. The variables are not at a fixed
 offset from \$sp

Object Layout

```
Class A {
    a: Int ← 0;
    b: Int ← 1;
    f(): Int {a ← a+d};
}
;
```

```
Class C inherits A {
    c: Int ← 3;
    h(): Int {a ← a*c};
}
;
```

```
Class B inherits A {
    b: Int ← 2;
    f(): Int {a};
    g(): Int {a ← a-b};
}
;
```

— =
Objects are laid out in contiguous memory

Each attribute is stored at a fixed offset in the object

— =
The attribute is in the same place for every object of that class

When the method is invoked, the object is 'self' and the fields are the object's attributes (like 'this' in java)

Offset — =

0	class tag	header information
4	object size	
8	dispatch pointer	
12	attribute 1	
16	attribute 2	
...		

Each block is a word long.

The class tag is an integer
— Identifies the class of the object

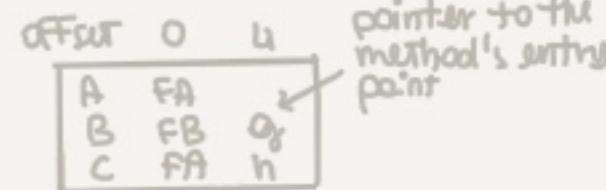
The object size is an integer
— Size of the object in words

Dispatch pointer is a pointer to a table of methods

Attributes in subsequent slots
— =

Given a layout for a class A, a layout for a subclss B can be defined by extending the layout of A

— Layout of A is unchanged (offset for A's attributes is the same for all subclasses)
— B is an extension



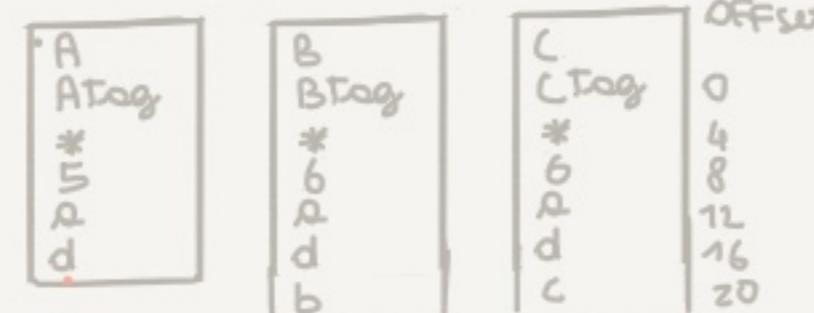
FB is method f overridden in B
The method f is always found at the same offset

The dispatch pointer in the object header points to the table of method entries for that class

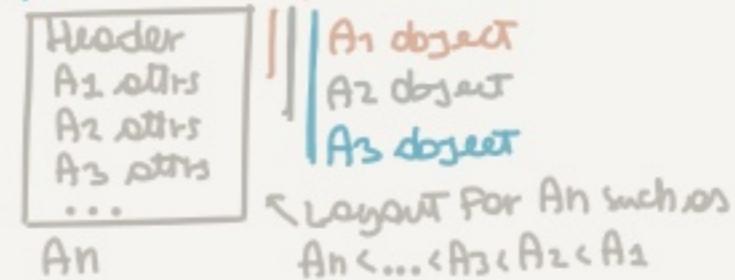
We use a separate table because that way it can be shared between all objects of that class
— =

To implement a dynamic dispatch e.g. f()
— evaluate 2, giving an object x
— call D[Of]

— D is the dispatch table for x
— In the cell, self is bound to x



If we have a chain of inheritance relationships
each prefix is a valid object



A dispatch table indexes a class' methods
— An array of method entry points
— A method F lives at fixed offset in a dispatch table for a class and all of its subclasses

Operational semantics

We must specify for every cool expression what happens when it's evaluated

- The "meaning" of an expression

— =
Quick recap

In defining a programming language

- The tokens → lexical analysis
- the grammar → syntactic analysis
- The typing rules → semantic analysis
- The evaluation rules → code generation and optimisation

We have specified evaluation rules indirectly

- Compilation of cool to a stack machine
- evaluation rules of the stack machine

But assembly language descriptions of language implementations have a lot of irrelevant detail

- Whether to use a stack machine
- Which way the stack grows
- How ints are represented
- ...

We need a complete description instead

- But not overly restrictive specification

There are many ways to specify semantics, we're going to use operational semantics

- Describes a program evaluation via execution rules
- Most useful for specifying implementations

Other ways to specify semantics are Denotational sem and Axiomatic sem

We're going to be using logical rules of inference, so int type checking

In type checking they had the form
 $\text{context } t \vdash e : C$

In a given context expression 'e' has type 'C'

For evaluation

$\text{context } t \vdash e : v$

In the given context expression 'e' evaluates to 'v'

$\text{context } t \vdash e_1 : S$

$\text{context } t \vdash e_2 : T$

$\text{context } t \vdash e_1 + e_2 : T$

— =

Consider the evaluation of $y \leftarrow x + 1$

To evaluate it have to know

- environment, where in memory a variable is stored
- store, what is in the memory

A variable environment maps variables to locations

$E = [a:l_1, b:l_2]$ list of variable and location pairs
keep track of which variables are in scope

— =
Store maps locations to values

$S = [l_1 \rightarrow 5, l_2 \rightarrow 7]$

$S' = S[l_2/l_1]$ takes S and updates the value at l_1 to be 12 by defining a new store S'

Stores are just functions

$x = (a_1=l_1, \dots, a_n=l_n)$ is a cool object where

- x is the class of the object
- a_i are the attributes (including inherited)
- l_i is the location of the value of a_i

— =

special cases (classes without attributes)

$\text{Int}(5); \text{Bool}(\text{true}); \text{String}(4, "cool")$
length

There is a special value void of type Object

- No operation can be performed on it
- Except for the test is void
- Concrete implementations might use NULL to represent void

— =
The evaluation judgement is

$s_0, E, S \vdash e : v, S'$

- s_0 is the self object
- E is the environment
- S is the store
- In that context an expression 'e' will ... produce a value 'v'
- ... produce a modified store S' only if e evaluates

Result of an evaluation is a value in a new store

The environment E and the self object don't change during the evaluation

Note: Contents of the self might change, but the layout of the object and its location doesn't change

Intermediate code

It is a language between the source language, and the target language.

- More details than source → e.g. you want to optimise register merge.
 - Fewer details than target ↑
- Above a particular instruction set of a machine. So it's easier to re-target to lots of different kinds of machines
- Some compilers have more than one interm. lang
- =

Intermediate language = high-level assembly

- Uses register names, an unlimited # (e.g. jump, labels)
- Uses control structures like assembly labels)
- Uses opcodes but some are higher level

Each instruction is of the form

$x := y \text{ op } z$ — y, z are either registers or constants
 $x := \text{op } z$

Really common form! → Three address code, it has three addresses in it:

- at most two arguments
- a destination

The expression $x + y * z$ is

$t_1 := y * z$ ← each subexpression
 $t_2 := x + t_1$ ↑ has a "name"

— =
 $\text{igen}(e, t)$; code to compute 'e' in register 't'

$\text{igen}(e_1 + e_2, t) =$
 $\quad \text{igen}(e_1, t_1)$
 $\quad \text{igen}(e_2, t_2)$
 $\quad t := t_1 + t_2$

Unlimited registers leads to very simple code generation of intermediate code

Optimisation

Lexical analysis → Parsing → Semantic analysis ...

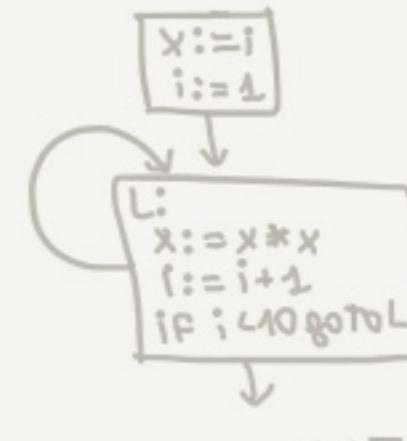
Optimisation → Code generation

Most complex part of the compiler

We optimise intermediate language

- Machine independent
- Exposes optimisation opportunities

$P \rightarrow S \ P \ I \ S$
 $S \rightarrow id := id \ op \ id \mid id := op \ id \mid$
 $id := id \mid push \ id \mid id := pop \mid$
 $if \ id \ relop \ id \ goto \ L \ L : \mid jump \ L$



- The body of a method can be represented as a control flow graph
- There is an initial node
- All return nodes have no edges coming out

optimisation improves resource utilisation

- execution time
 - code size
 - memory usage and so on
- ... without altering what the program produces

— =
(1) Local optimisation
Applies to basic blocks

(2) Global optimisation
Applies to control-flow graphs

(3) Inter-procedural optimisation
Applies across method boundaries

Most compilers do (1), many do (2), few do (3)

- =
- Some optimisations are hard to implement
 - Some are costly in compilation time
 - ... Or have low-payoff

- Many fancy optimisations have all three

Goal: Maximum benefit for minimum cost!

A basic block is a maximal sequence of instructions

- No label, except first instruction
- No jumps, except last instruction

Ideas, execution is guaranteed to proceed from the first to the last block statement

1: $t := 2 * x$ (3) presents only after (2)
2: $w := t + x$ (3) can be $w := 3 * x$
3: if $w > 0$ goto L' - We could eliminate (2) after that optimisation if 't' is not used elsewhere

A control-flow graph is a directed graph of basic blocks. There's an edge from node A to B if the execution can pass from A's instruction to B's first instruction

Local optimisation

Optimise one basic block

can be deleted if $w := \text{rhs}$ and w is never used

$x := x + 0$ can be simplified

$x := x * 1$

$x := x * 8$ $x := x \ll 3$ \equiv $y := y * 2$ $y := y * y$
bit shifts are usually machine instruction faster than $*$

\equiv Algebraic simplifications

Compute operations on constants at compile time Constant folding

$x := 2 + 2$ If $2 < 0$ jump L // deleted
 $\Rightarrow x := 4$

REALLY COMMON!

Eliminate unreachable basic blocks

- Makes the code smaller

It is quite common

e.g. libraries. Not used methods can be removed to make the code smaller
debugging. Not running a piece of code while debugging
result of other optimisations. They may create unreachable blocks

\equiv

Some optimisations are simpler to express if each register occurs only once on the left-hand side of an assignment

Single assignment form

$x := z + y$ $b := z + y$

$a := x \Rightarrow a := b$

e.g. $x := 2 * x$ $x := 2 * b$ common subexpression

- if basic block in single a. form
then all assignments with the same right hand side compute the same value

- If $w := x$ appears in a block, then all subsequent w 's can be replaced with x - copy propagation

if $w := \text{rhs}$ and w is never used
then the statement can be deleted

\equiv
Typically one optimisation enables another
- repeat optimisation until no improvement is possible

$\begin{array}{ll} l := x * * x & l := x * x \\ b := 3 & b := 3 \\ c := x & c := x \\ d := c * c \xrightarrow{\text{algebraic}} & d := c * c \\ e := b * 2 & e := b \ll 1 \\ f := a + d & f := a + d \\ g = e * f & g = e * f \end{array}$

$\begin{array}{ll} \xrightarrow{\text{copy}} l := x * x & l := x * x \\ \xrightarrow{\text{propag.}} b := 3 & b := 3 \\ c := x & c := x \\ d := x * x & d := x * x \\ e := 3 \ll 1 & e := 6 \\ f := a + d & f := a + d \\ g = e * f & g = e * f \end{array}$

$\begin{array}{ll} \xrightarrow{\text{common subexp.}} l := x * x & l := x * x \\ c := x & c := x \\ d := l & \xrightarrow{\text{copy}} d := l \\ e := 6 & e := 6 \\ f := a + d & f := a + d \\ g = e * f & g = e * f \end{array}$

Peephole optimisation

Optimise assembly code directly

The "peephole" is a short sequence of

Replace the sequence with another equivalent one (but faster)

\equiv
Peephole optimisations are written as replacement rules $i_1, \dots, i_n \rightarrow j_1, \dots, j_m$

move \$a \$b, move \$b \$a \rightarrow move \$a \$b
add \$a \$a i, add \$a \$a j \rightarrow

add \$a \$a i+j

\equiv

Many optimisations from local optimisations can be used as peephole optimisations

Just like locals, they must be applied repeatedly

Code produced by "optimisers" is not optimal,
it is just improved

$\begin{array}{l} l := x * x \\ f := a + l \\ g = 6 * f \end{array}$

$\xrightarrow{\text{dual code elim.}}$
There are other optimisations that are harder to find

$\begin{array}{l} l := x * x \\ f := 2 * l \Rightarrow g := 12 * l \\ g := 12 * l \end{array}$

Dataflow analysis

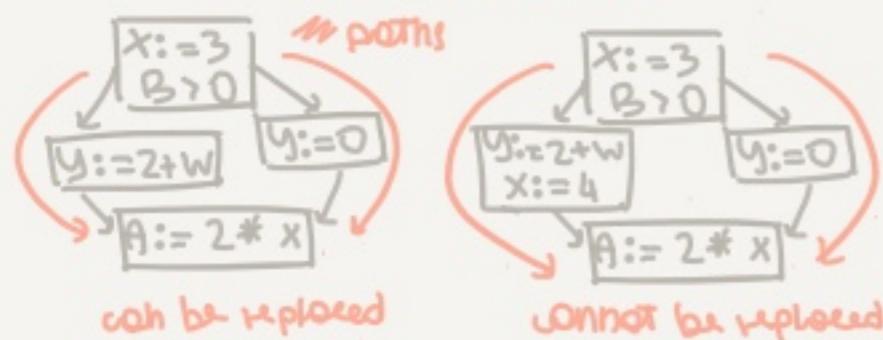
Local optimizations can be extended to an entire control-flow graph

Cybernetic decision analysis is a technique to solve problems with these characteristics

- The optimisation depends on knowing a property x at a particular point in the program
 - Proving x at any point requires knowledge of the entire program
 - It is ok to be conservative. Either x is true or we don't know (and don't do the optimisation)

Constant propagation

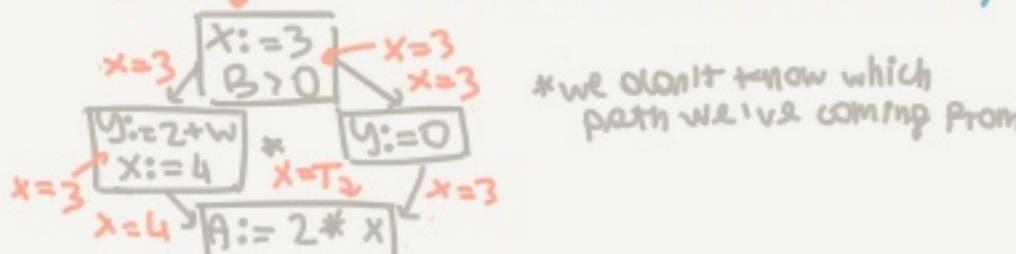
- To replace a var of x with a constant k if
 - On every path to the var of x, the last assignment to x is $x := k$



For statements with x , we associate

T "top"; x is not a constant (we can't figure out if x is a constant)
 C ; x is a constant

\downarrow "bottom"; this statement never executed
 $x = T \downarrow$ (we don't know if it executes)



The analysis of a complicated program can be expressed as a combination of simple rules relating to the change of information between adjacent statements

$C(x, s, \text{in})$ = value of x before s "C stands for
 $C(x, s, \text{out})$ = value of x after s "constant information"

A transfer function transfers information from one statement to another

Every statement has a set of predicates
— = $(p_1 \dots p_n)$

Rule 1 $p_1 p_2 p_3 p_4$ if $C(p_i, x, \text{out}) = T$
 $x \in T$ \downarrow for any i ,
 s \boxed{s} then $C(x, s, \text{in}) = T$

Rule 2  If $c(p_i, x, \alpha_{\text{mt}}) = c \notin$
 $c(p_j, x, \alpha_{\text{mt}}) = d \in$
 then $c(\cdot, \cdot, \cdot)$ not equal

Rule 3 then $C(s, x, \text{in}) = T$

$x=1 \downarrow \quad \downarrow x=c \quad \xrightarrow{x=c}$

$x \text{ is } c \quad \boxed{s}$

if $C(p_i, x, \text{out}) = c \text{ or } 1$
for all i
then $C(s, x, \text{in}) = c$

Rule 4 then $C(s, x, \text{out}) = 1$

$x=1 \downarrow \quad y=1 \downarrow \quad z=1 \downarrow$

\boxed{S}

If $C(p_i, x, \text{out}) = 1$
for all i :
then $C(s, x, \text{out}) = 1$

Rules 1-4 are for out-in statements
next we will see for in-out

Rule 5 $\downarrow x = 1$

$C(x := c, x, \text{ONT}) = 1$
 $\text{if } C(s, x, \text{in}) = 1$

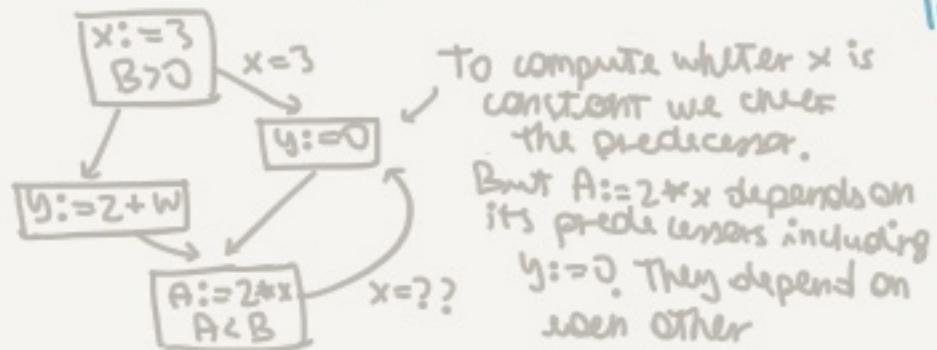
Rule 7 \downarrow $C(x := f(...), x, \text{out}) = T$
 Rule 5 has precedence over rule 7

Rule 8 $\boxed{y := \dots}$ over rule 7

Algorithm

1. For every entry s to the program,
SET $C(s, x, \text{in}) = T$
 2. Set $C(s, x, \text{in}) = C(s, x, \text{out}) = L$
everywhere else
 3. Repeat until all points satisfy 1-8
 - Pick s not satisfying 1-8 and
update using appropriate rule

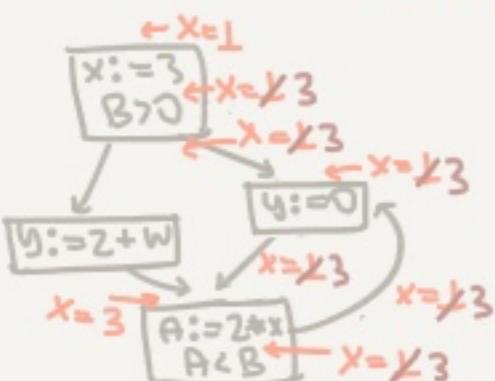
Analysis of loops



Break the cycle by starting with an initial guess \perp .

This allows us to make progress.

By following the algorithm previously described (which also assigns initial guesses):



Ordering

\perp, \sqsubset, T are called **abstract values**

We can simplify things by ordering the abstract values

$$\perp \sqsubset \sqsubset T$$

We can now define lub (least upper bound) operation. lub takes the smallest element that is bigger than anything in the least upper bound. Then we have incomparable (both c)

Rules 1-4 can be written as

$$C(s, x, \text{in}) = \text{lub} \{ C(p, x, \text{out}) \mid p \text{ predecessor of } s \}$$

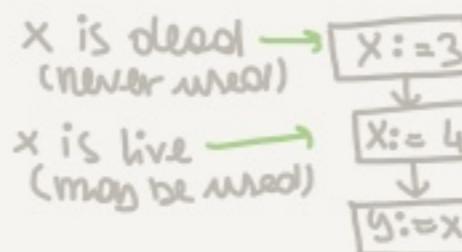
lub explains why the algorithm terminates

- Values start at \perp and only increase
- \perp can change to a constant, and even to T
- Thus $C(s, x, -)$ can change at most twice

Constant propagation then is linear in program size

$$\begin{aligned} \text{steps} &= C(\dots) \text{ values computed} \approx 2 \\ &= \# \text{ of program statements} \approx 4 \end{aligned}$$

Liveness analysis



A variable is live at statement s if

- There exist a statement s_i that uses x
- There is a path from s to s_i
- The path have no assignments to x

A statement $x := \dots$ is dead code if x is dead after the assignment (not live)

- That statement can be deleted

=

$$\text{Rule 1} \quad L(p, x, \text{out}) = \vee \{ L(s, x, \text{in}) \mid s \text{ successor of } p \}$$

$$\text{Rule 2} \quad \downarrow \leftarrow x \text{ is live} \quad L(s, x, \text{in}) = \text{True if } s \text{ refers to } x \text{ or truths}$$

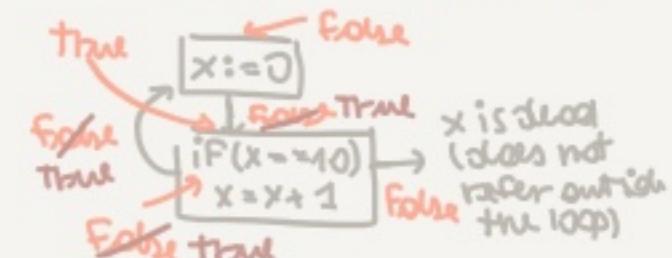
$$\text{Rule 3} \quad \downarrow \leftarrow x \text{ not live}$$

$$x := \alpha \quad L(x := \alpha, x, \text{in}) = \text{False if } \alpha \text{ does not refer to } x$$

$$\text{Rule 4} \quad \downarrow \quad L(s, x, \text{in}) = L(s, x, \text{out}) \text{ IF } s \text{ does not refer to } x$$

Algorithm

1. Let all $L(\dots) = \text{False}$ initially
2. Repeat until all statements s satisfies rules 1-4
- Pick s where one of 1-4 does not hold and update using appropriate rule



Each value can only increment (ordering False < True)
Since it can only change once, termination is guaranteed
=

Constant propagation is forward analysis

- Information pushed from inputs to outputs

Liveness is backward analysis

- From outputs back toward inputs

Register allocation

Intermediate code uses unlimited temporaries

- Simplifies code gen and optimisation
- But translation to assembly is harder

Problem:

Rewrite intermediate code to use no more temporaries than there are registers

Solution: Assign multiple temporaries to one register (without changing program's behaviour)

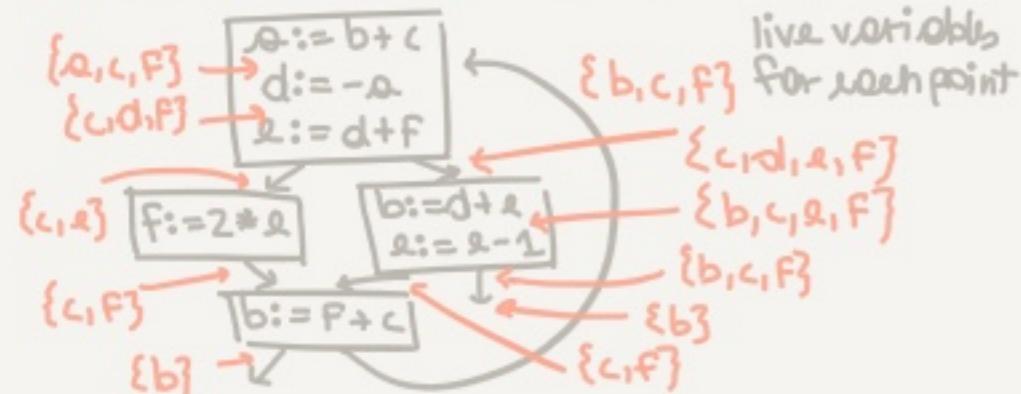
$$\begin{array}{ll} A := c + d & t_1 := r_2 + r_3 \\ L := A + b & t_2 := r_1 + r_4 \\ F := L - 1 & r_1 := r_2 - 1 \end{array}$$

- Announce A and L are dead after their use
- Allocate A , L and F all to one register

Register allocation in the original FORTRAN project was a bottleneck in the quality of the compiler

The Register allocation scheme based on graph coloring solved the issue years later

=
Temporaries t_1 and t_2 can share a register if they are not live at the same time



Construct an undirected graph with a node for each temporary and an edge between if the temporaries are live at the same time at some point in the program

This is called register interference graph (RIG)

- Two temporaries can be allocated to the same register, if there is no edge between



RIG for our example

the RIG gives a global picture of the register requirements

After RIG, the register allocation algorithm is architecture independent

Graph coloring

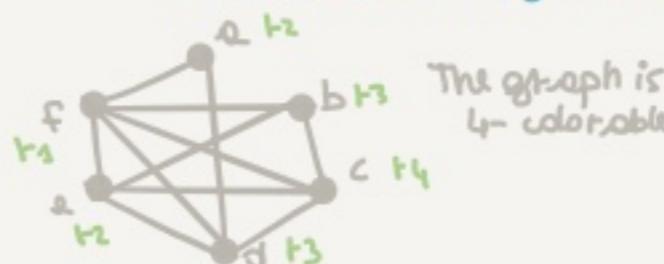
Assign colors to nodes such that nodes such that nodes connected by an edge have different color

A graph is k -colorable if it has a coloring with k colors

=

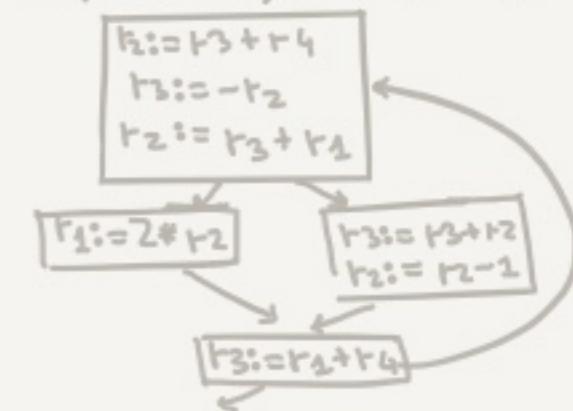
Let $k = \text{number of registers}$

If the RIG is k colorable, there is a register assignment that uses no more than k registers



Control flow graph for RIG.

- Replace temporaries with registers



=
Graph coloring is NP-hard

- Compilers usually use heuristics

Observation: DIVIDE AND CONQUER

- Pick a node with fewer than k neighbours in RIG
- Eliminate it and its edges in RIG
- If result is k -colorable, so is the original graph

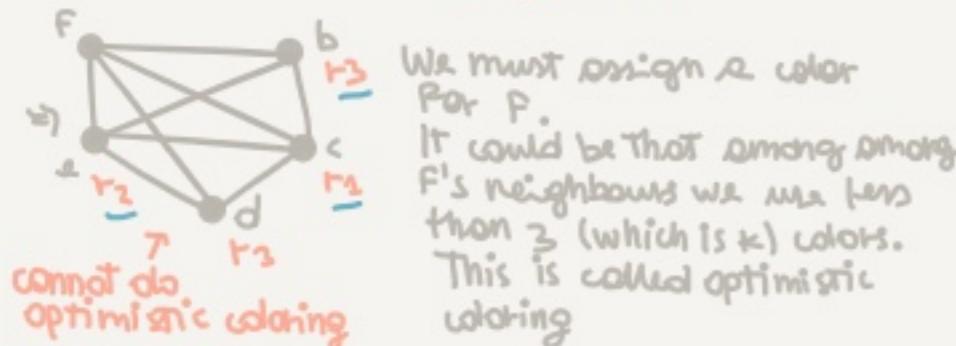
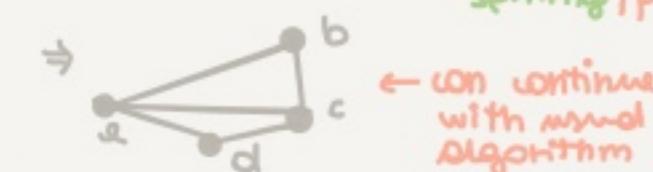
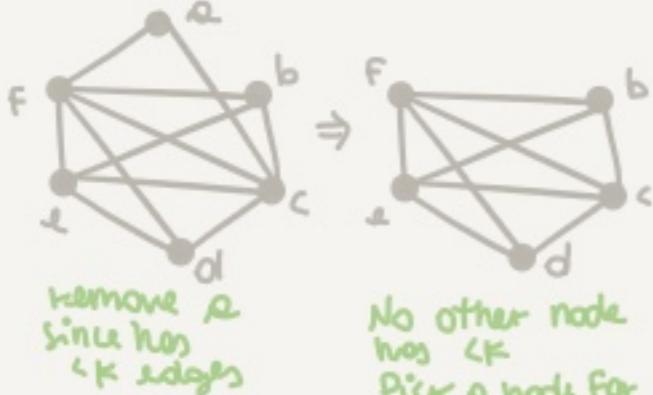
Let $c_1 \dots c_k$ be the colors assigned to to t neighbours, since $n < k$ we can pick some color for t that is different from those of its neighbours

Alg.: =

- Pick node t with fewer than k neighbours
- Put t on the stack and remove it from RIG
- Repeat until graph is empty
- Start with last node colored
- At each step, pick a color from those assigned to already coloured neighbours

Spilling

If all nodes or the RIG have k or more neighbours. Then there are not enough registers and some values are spilled in memory



If optimistic coloring fails we spill F

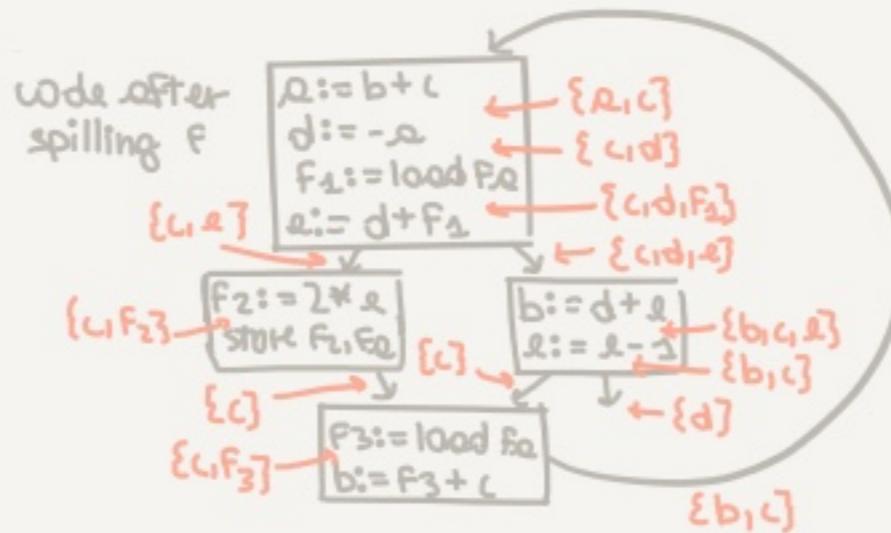
- Allocate a memory location for F
- Typically in the current stack frame
- Call this f_a

Before each operation that reads F insert

$f := \text{load } f_a$

After each operation that writes F

Store f, f_a

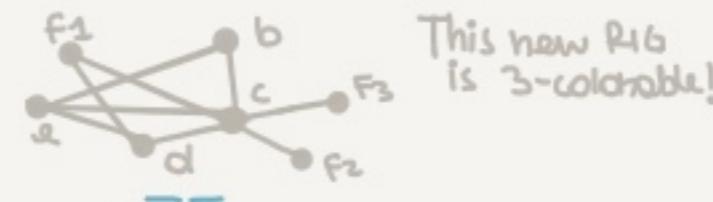


f_i is live only

- Between $f_i = \text{load } f_a$ and the next instruction
- Between a store f_i, f_a and the preceding instr.

Spilling reduces live range of F

- Fewer RIG neighbours



Additional spills may be required before a coloring is found

Some spilling choices lead to better code but any choice is correct

- Spill temporary with most conflicts
- Spill temporaries with few users
- Avoid spilling inner loops

Managing caches

Power usage, limits size and speed of registers and caches

The cost of a cache miss is really high and typically requires more than one level of cache to bridge a fast processor with large main memory

Compilers are really good at managing registers but not so good at managing caches

- Up to the programmer to make good use of it

Consider the loop

```
for (j:=1; j<10; j++)
    for (i:=1; i<1000000; i++)
        a[i] *= b[i]
```

Every iteration will be a cache miss

```
a[1] b[1] a[2] b[2] ...
```

```
for (i:=1; i<1000000; i++)
    for (j:=1; j<10; j++)
        a[i] *= b[i]
```

- Computes the same thing
- But with much better cache behaviour
- More than 10x faster

A compiler can perform this **loop interchange** optimisation

Not all compilers implement this because it might be hard to decide whether you can do it. Usually it's up to the programmer

Garbage collection

C and C++ have manual memory management

Many storage bugs

- Forget to free unused memory
- Dereferencing dangling pointers
- Overwriting parts of a data structure
- Others...

These bugs are hard to find!

=
This problem has been studied since
1950s for LISP

Become mainstream with Java

=
A program can only use the objects
that it can find

An object x is reachable if

- a register contains a pointer to x
- another reachable object y contains
a pointer to x

Reachable objects can be found by
starting from registers and
following all the pointers

Unreachable objects can never be used
and are called garbage

```
x ← new A;  
y ← new B;  
x ← y;  
if alwaysTrue()  
then x ← new A  
else x.fool() fi
```

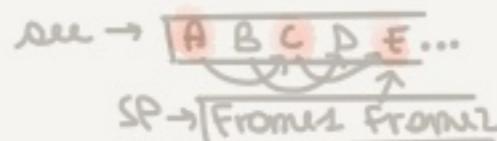
$x \not\rightarrow [A]$ ← A is garbage
 $y \rightarrow [B]$ ← collected
garbage collect here

B is reachable
(but it is going
to be overwritten,
meaning it's never used)

Reachability is an approximation,
just because it is reachable it
doesn't guarantee that it will be used

Java uses an accumulator

- points to an object which may point to others
- ... and a stack pointer
- both stack frame contains pointers
 - e.g. method parameters
- both stack frames also contains non-pointers
 - e.g. return address



B and D are
unreachable.
(D has a pointer
but from an
unreachable object)

=
Each garbage collection scheme has these steps

- 1 Allocate space for new objects
- 2 When space runs out
 - compute what objects might be reused
 - Free space used by objects not found in (a)

Some strategies perform garbage collection
before the space runs out

Java Arrays

Assume $B \subset A$

```
B[] b = new B[10]; b → B's
A[] a = b;
a[0] = new A();
b[0].methodNotDeclaredInA();
```

Runtime error!

Having multiple classes to updateable locations with different types is unusual

This problem is common among programming languages

Standard solution: $B[] \subset A[]$ if $B=A$
 Subtyping rule for arrays

Java checks each array assignment at runtime

Adopts overloads on array computations
primitive type arrays are not affected

Java Exceptions

In a section of your code, you find an error

- out of memory
- a list is supposed to be sorted but is not
- etc.

Add a new type (class) of exceptions

Add new forms

```
try {...} catch(x) {...}
      throw exception
```

$T(v) = \text{exception thrown with value } v$

$E \vdash e_1 : v_1$

$\frac{}{E \vdash \text{try}\{e_1\} \text{ catch}(x)\{e_2\}:v_1}$

$E \vdash e_2 : T(v_1)$

$E[x \leftarrow v_1] \vdash e_2 : v_2$

$\frac{}{E \vdash \text{try}\{e_1\} \text{ catch}(x)\{e_2\}:v_2}$

When we encounter try

- Mark current stack position

When we throw an exception

- Unwind the stack to the first try
- Execute the corresponding catch

More complex techniques reduce the cost of try and catch

—

In Java, exceptions are part of the method interface and checked by the compiler

public void x() throws MyException

In the original Java project they observed that there are many exceptions that could be caused by Java programs and people usually limit themselves of exceptions to handle

Some exceptions do not need to be part of the method signature

- dereferencing null
- others ...

throw must be applied to an object of type Exception
Interfaces

Specify relationships between classes without using inheritance

```
interface PointInterface{void move(int dx,
int dy);}
```

```
class Point implements PointInter.
{ void move(int dx,int dy){...}}
```

From the Java language manual

"Java programs can use interfaces to make it unnecessary for related classes to share a common abstract superclass or to add methods to Object".
i.e. they play the same role as multiple inheritance in C++, because a class can implement multiple interfaces

Inheritance is more efficient and should be used where possible

- Classes implementing interfaces need not be at fixed offsets

```
class Point implements PointInterface {
    void move (int dx,int dy){...}}
```

```
class Point2 implements PointInterface {
    void dummy(){...}
    void move (int dx,int dy){...}}
```

—

Dispatch L.F() where L has an interface is more complex than usual

One approach:

- Each class implementing an interface has a lookup table method names \rightarrow methods
- Hash method name for faster lookup

Java coercions

Java allows primitive types to be coerced in certain contexts

$1 + 2.0$ int 1 is widened to float 1.0

A coercion is a primitive function that the compiler inserts for you

Java distinguishes coercions and casts in

- Widening, always succeeds ($\text{int} \rightarrow \text{float}$)
- Narrowing, may fail if data can't be converted to desired type ($\text{float} \rightarrow \text{int}$)

Narrowing casts must be explicit
widening coercions/casts can be implicit

Casts in can lead to surprising behaviour

Java threads

Each thread in Java has its own program counter and stack

(local variables and
definition records)

Threads are objects of type Thread

- start and stop methods

Synchronized (x) { ... }

- locks x, executes x, and unlocks x

Or can be put on methods (more common)

Synchronized f () { ... }

- object implicitly locked

class Simple {

 int a=1, b=2;

 void synchronized to () { ... }

 void fto () { ... }

 void synchronized to () { ... }

WRONG!

Both methods
must have the lock

That does nothing!

This is synchronised
instead

A variable should only hold values written by some thread

- Writes of values are atomic

Cannot switch thread while writing and only write some bits

- Java guarantees that, except for doubles

It consumes two words and translates into two machine instructions
Values created from two different threads are called "out of thin air values"

By declaring double volatile, then the writes will be atomic

Other Topics

Java allows classes to be loaded at runtime

- Type checking of the source takes place at compile time
- Bytecode verification takes place at runtime

Loading policies are handled by the ClassLoader

— =

A class is initialised when a symbol in the class is first used

- Not when the class is loaded
- Delays initialisation errors to a predictable point

- 1 Lock the object for the class
- 2 Check if it is already initialised

3 If initialised, return

4 Else make initialisation as in progress by this thread and unlock the class

5 Initialise superclasses (lexical order)

- static and final fields first

- Give every field default value before init.

6. Any errors result in an incorrectly initialised class, mark as erroneous

7. If no errors, lock class, label class as initialised, notify threads waiting on class object, unlock class