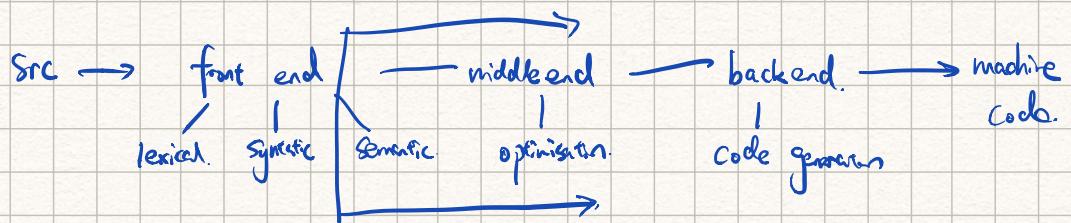
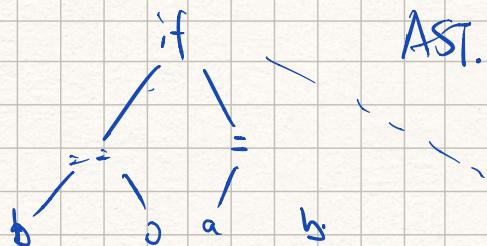


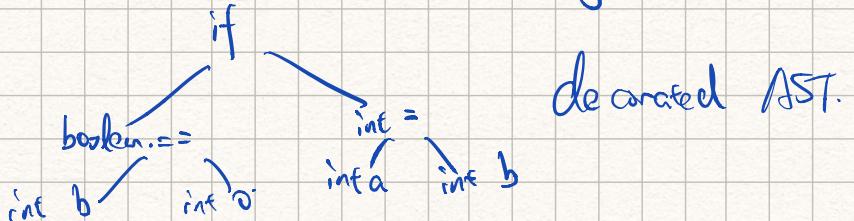
Compilers Back-End Notes



From front-end, we are here.

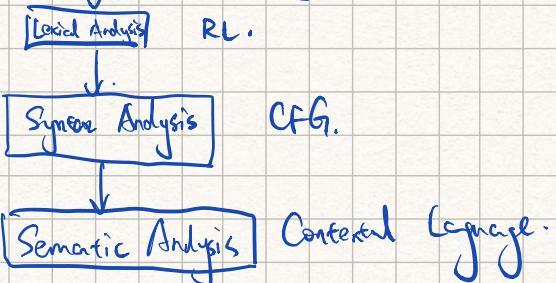


1. Semantic analysis



decorated AST.

The Language Hierarchy



All syntactically correct programs are CFG.
Compilable programs are not.

CFG

→ regardless of the context of a non-terminal,
scope/ types are irrelevant to the parser.

Context-Sensitive Grammer.

$$G = (N, \Sigma, P, S)$$

? nonterm terminals
↑ production rules

is context sensitive if all $\beta \in P$ are

$$\alpha A \beta \rightarrow \alpha \gamma \beta$$

$$AGN: \sim \rightarrow F(N \cup \Sigma)^*$$

$$\sim F(N \cup \Sigma)^+$$

Limitations of CFGs.

We can add. context in CFG.

but very complicated!

If we need grammar we need. Type 0 grammar.
Unrestricted — TM.

Contextual (Semantic) Analysis

Outside of grammar check.

- * undefined symbols
- * type check errors
- * scope errors.
- * classes and method ambiguities.
- * inheritance.

- Look at Scoping Rules / Typing Rules.

Scoping Rules.

Declaration.

- explicit var i : integer;
- implicit any variable start with I is a integer.

Maybe independent from Name declarations.

Scope rules determine which declaration applies in the context.

Binding

an assoc name \leftrightarrow things

broken for (int n, m) {

 tmp;
 if (n < m) tmp = true
 else --- false
 Return tmp
}.

foo() → int × int. → bool.

n → int.

m → int.

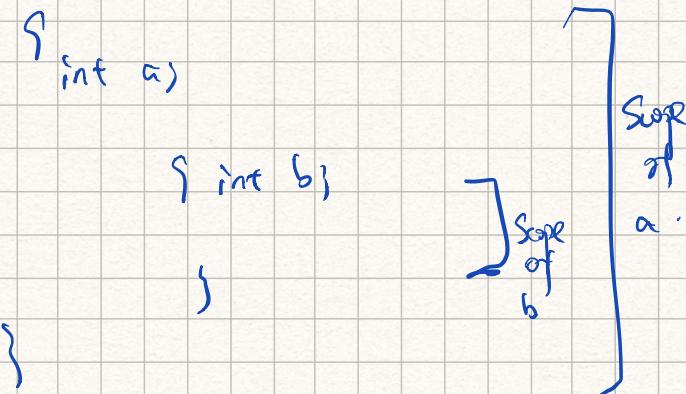
tmp → bool.

Scope: textual region where a binding is active.

Static Scope.

Dynamic Scope.

Scope of vars.



Scope of arguments of func.

public int factorial (int n) { (Scope of n)

}.

OOP makes things + complex.

Class defines a scope that covers those of all methods.

```

class A {
    private x;
    public void g();
}

```

Class B extends A.

```

public void h() { g(); }
}

```

Static Scope.

Lexical Scope. → Storing info in the src.
 Common choice. Scope determined at compile time.
 typically, most recent encounter. when falling S.E.

Simplest design. - one global scope.

two scope. - global / local to a function

Nested function. / more complications

Hole and Qualifier.

* a binding hidden by a same-name nested declaration.

Hole.

Var x = 2.

if x > 0

 int x=5] hidden (Hole).
 { ~~x=~~ xf=1. }

}

* Some language still allows access to outer / more global var by applying qualifier or Scope operator

int cout > 42;
std::cout << cout << endl;
 ↑ ↓
 Qualifier function / not dt.

Dynamic Scope.

binding depends on execution flow.

- * whatever declaration recently take place.
- * binding copies when execution leaves lexical scope.

In general, flow of execution is not predictable at compile time.

Dynamic Scoping more implemented in interpreted languages.

Dynamic Scoping

- * can be dangerous
- * easier to implement.

Referencing environment.

* Scoping rules → where to bind.

↔ Referencing Environment → when to bind.

↓
use it (seq. of scopes to execute to find
the current binding of a given name)

Shallow deep binding.

(dynamic). Shallow binding happens when the function is actually called.

deep binding

happens when the function is first referred to

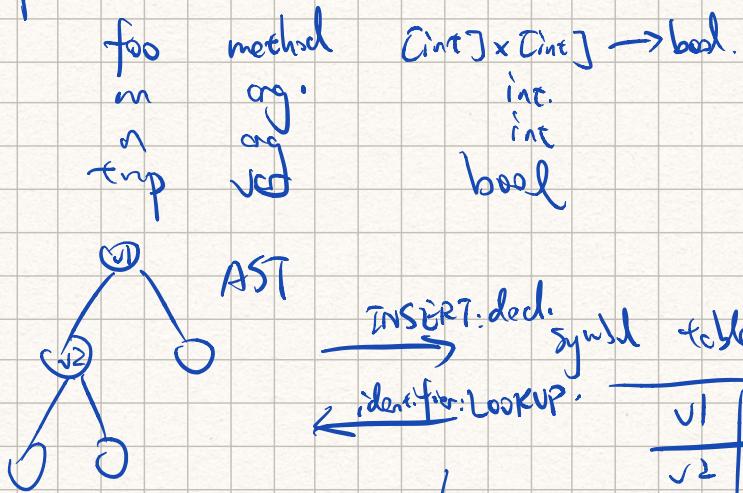
most use deep binding.

Implementation.

Symbol Tables

focus on static Scoping (deep binding rules)

Example.



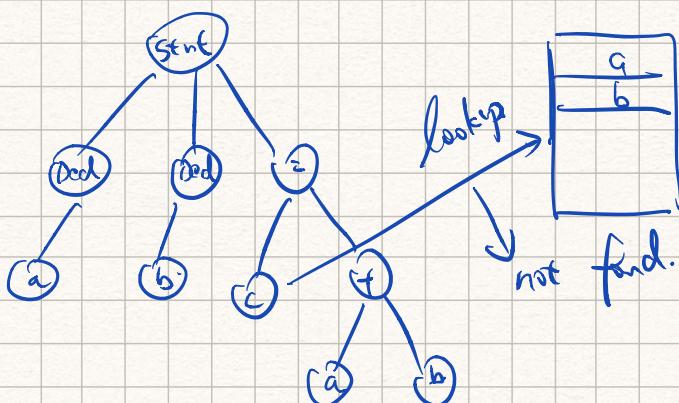
INSERT: decl. symbol table.

lookup

v1	1
v2	-

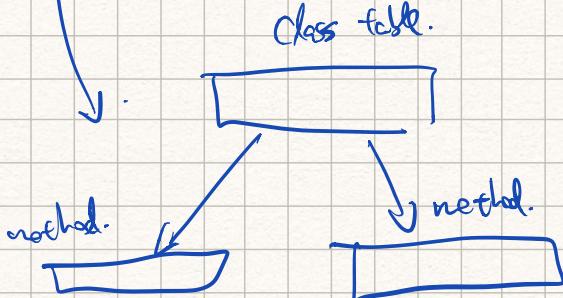
If cannot find id. in symbol table, "unidentified id."

Example: check for undecl. identifiers



Single table is not sufficient.

- Nested scopes - inner binding hide outer bindings
- Forward references. - use when not declared. (e.g. functions)



find which declaration is active.

- * Start from current.
- * go up until find it.
- * don't go downwards.

The hierarchy solves the name collision.

found in deeper tables.

Forward References

- * construct tree hierarchy of tables first.
- * then, traverse the symbol table tree.

Implementation.

Two ways:

- * one table, add a column with scope information.

- * multiple tables: stack tree of tables

Symbol Table Op.

- Create a new empty symbol table with a given parent table.

- Insert a new id to a table.

- Lookup

Implementation: Individual Tables

* Linked list.

* BST.

* Hashtable.

- good balance.

- good hashing

Type checking.

e.g. has x been used according to its type
are parameters passed in of expected types

Type System.: Define type expressions and rules

all possible types

constraints on types

Type expressions.

Basic types (primitive types)

int, double, boolean.

Build types from basic types using-

- type constructs.

- function types.

* class types

OOP.

```
class A {  
    int i;  
    public float f;  
};
```

Type: $T_1 \times T_2 \times T_3 \dots \times T_n \rightarrow T_r$.

foo: int \times double \rightarrow double.

{ array
Struct.
points. } \rightarrow
int[i]
int[10]

struct q
int a;
float b;
};

int * z.

public A foo (int i, B[] b, double [] d.)

foo: i:int, b:array(B), d:array(double) → A.

public class A { }.

int [0] a;

boolean b;

double [] c;

}.

{ a: array(int, 10), b: boolean,

c: array(C) }.

Type Info in symbol table.

Name ↪ type.

parameter ↪ type.

method ↪ para, result.

class ↪ instance var / method decl:

int pow (int n, int m) { }.

int i=0

int result=1

{

|

}.

pow int × int → int	
n	int
m	int
i	int
result	int

class A.

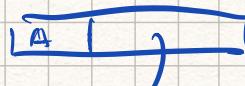
private void m1 - -

public double m2 - -

j int int

m1 void → void FALSE

m2 int × double → double TRUE

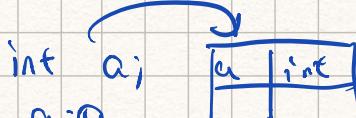


public?

FALSE

TRUE

Type rules



int a; \rightarrow int
 double b; \rightarrow double.
 int c;
 $c = a + b \rightarrow (a + b) \rightarrow$. int + double
 $=$ double + double.
 $=$ double.

for each assign stmt. LHS. must be Sgr type.

Comparison. should have same type.

7 22 C

result is boolean

int n = pow(i+j, j)

Type of $(\text{pow}(\text{it}^j, j))$ = int \times int
→ int

matdos typeof (n) = int.

Implementation

Construct Symbol Table.

* traverse (AST.)

* use type info in the table

Slope and type decks

* After STable built., visit to AST.

and access to
Symbol table.

* Ensure Scope / type. are respected.

If no formal referencing, 1 single visit is enough

Error Handling.

When error found. → report error and
continue.

also may need to create right type.
if not went a cascading errors

Semantic Analyzer.

1. process decl.

add. to STable.

report multiple decl.

2. process stmt.

undeclared var.

update [↑] guesses of ID

update ID nodes of AST to STable.

3. Go through all stmt. again.

check for type errors.

Intermediate Representation.

LLVM.

High Level IR.

extract original src intent easier.

Low Lvl. IR.

generate machine code easier

IR should.

independent of source or target lang.

Convenient to produce during semantic analysis

- translate to multiple targets

AST \rightarrow IR \rightarrow Machine Code

instruction for a VM

abstract

GIMPLE / LLVM / Java Bytecode

Ideal for Optimisation

, AST about 40 different node types

IR a dozen instructions
IR for 2 types of machines
Register Machines Stack Machines

$2+3$,

Load R1 2.
Load R2 3

ADD R1, R2

↓

5 goes to R1

Identifier based temporary vars.

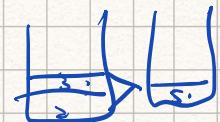
TAC (Three address Code)

$2+3$

PUSH 2
PUSH 3

ADD

pop 2
pop 3 → result
push result



JVM.

Performed on Stack.

data recovery :

Data : store vars. / referenced by

constant pool.

addr
or
ID

Instruction : code.

Register Machine.

→ TAC. has at most three addr (or less).

$$a = (b \times c)^* (-d)$$

$$+1 = b \times c \quad +2 = -d \quad a = f1 \times f2.$$

Type of TAC.

Binary $a = b \text{ OP } c$.

Arithmetic: ADD, SUB, DIV, MUL, MOD

Logical: AND, OR, XOR

Comparisons: EQ, NEQ, LE, LEQ - - -

Unary $a \text{ OP } c$.

MINUS, NEG

COPY: $a = b$.

Array access, $a = b[i]$ $a[i] = b$

Field access: $a = b.f$ $b.f = a$.

Control flow:

Labelling label L.

Unconditional jump: jump L.

Conditional jump: $\text{if } a \text{ L}$

$\text{if } a \text{ L}$

Function calls

param pi

call f

$a = \text{call } f$

return

return y.

Temporary vars

operands can be.

vars.

consts,

temp vars

use temp vars to store intermediate values

TAC Example.

$n = 0$
while ($n < 10$) {
 $n = n + 1$
}

$n = 0$.
label L.
 $t2 = n \text{ LEQ } 10$
fjump $t2$ END.
 $n = n \text{ ADD } 1$.
jmp L.
label END

Translation to TAC.

Need an algo to do so.

- * Start from AST
- * define rules for each AST node type
- * Recursively apply template

Given e.

$T[e]$ is a translated low level IR.

$t = T[e]$ means t stores the result of low level IR e.

$$t = T[e] \Leftrightarrow t = e.$$

Unary operators.

$$t = T[e] \quad t = \text{OP} \quad t1.$$

Binary Ops

$$t1 = T[e1].$$

$$t2 = T[e2]$$

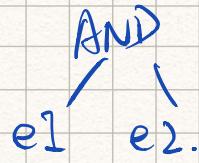
$$t = t1 \text{ OP } t2.$$

Short-circuit OR.

$$\begin{array}{c} \text{OR} \\ e1 \quad e2 \end{array}$$

$t = T[e1]$
 $t \text{ jump } t \text{ END} \leftarrow \text{if } e1 \text{ true}$
 $t = T[e2]$
must be

short-circuit AND.



label END.

$$t = T[e_1]$$

$t \text{ jump } t \text{ NEXT} \leftarrow$ if e_1 true
jump END

label NEXT

$$t = T[e_2]$$

label END

then evaluate
 e_2

otherwise
must false

Array access

$$t = T[v[e]]$$

$$t_1 = T[e]$$

$$t = v[t_1]$$

Field access

$$t = T[e_1.e_2]$$

$$t_1 = T[e_1]$$

$$t_2 = T[e_2]$$

$$t = t_1.t_2$$

List of statements

$$T[s_1; s_2; \dots].$$

$$T[s_1]$$

$$T[s_2]$$

Concat of IR instructions

Assignments

$$T[v = e].$$

$$v = T[e].$$

$$T[v[e_1] = e_2]$$

$$t_1 = T[e_1]$$

$$t_2 = T[e_2]$$

$$v[t_1] = t_2.$$

Control Flow.

If-then.

$$T[\text{if}(e) \text{ then } s]$$



$$t_1 = T[e]$$

fjump t_1 END

$$T[s].$$

END

If then else.
e / \ s1 s2.

$T[\text{if}(e) \text{ then } s1 \text{ else } s2]$

$\Leftrightarrow t1 = T(e)$

fjmp t1 ELSE

$T[s1]$

jump END

ELSE

$T[s2]$

END

Function calls

$T[\text{call } f(e1, e2 \dots en)]$

param $T[e1]$

param $T[e2]$

:

param $T[en]$

call f.

return

$t = T[e]$
return t.

Example.

if (c) {
 if (d) {
 a = b
 }
}
{print(x)}

$T[- - -]$

U.

$T[\text{if}(c) \dots]$
 $T[\text{print}(x)],$

$t1 = T(c)$

fjmp t1 END

$T[\text{if}(d) \dots]$

label END
param x.

call print

U

(

)

TAC Exercise.

```

t1 = n >= 0
if n >= 0 return 0
else S n2 = 1. if n = 1
    return 1.
else S.

```

$t3 = n - 1$
~~BAR(n-1)~~
 $t4 = n - 2$.
~~BAR(n-2)~~
return ~~BAR(n-1) + BAR(n-2)~~

calc Fibonacci
Seq.

A is de input

Implementing TAC.

- * Many ops don't use reg 2.

- * jugs put labels in result.

Triples : results are priority used for jumps

Problem:

Optimisation often makes code unreadable.

broke triples (indexes no longer preserved)

* Indirect Triples

addr	inst
45	[20]
46	[21]
47	[22]
⋮	⋮
⋮	⋮

+ b [20].

* quadruples

* easier to optimise.

* redundant temporaries

* triples

* moving Stmt. change references.

* Indirect triples.

* moving stmt change the referencing.
no change to triples

* save space if temp. is reused.

Stack Machines

* Operands results in the stack

* Register allocation not an issue

* ASM IR Java ByteCode

Java ByteCode. (Const. pooled out. begin of class)

* load and store.

* arithmetic and logic.

* type conversion.

* object creation and manipulation.

* operand stack management.

* control flow transfer.

* method invocation and return.

* exception/ sync and etc.

Load / store.

a: reference

i: integer

l: long

Pros: Stack-based IR

* compact code. 1 dram.

e.g. illed #S

* simpler compiler / code gen is simpler.
independent
from other
code.

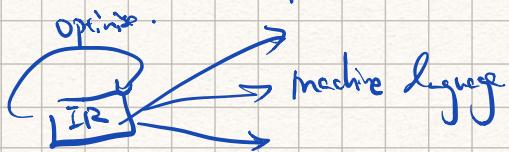
* simpler interpreter. / mem access centralized
less variants of
instructions

Cons:

- * more mem references. (no tmp registers)
- * common subexpression are hard to factor out.
- * less flexible (hard to reuse code)
when optimizing.

Optimisation outline.

- * High level lvl. optimisation
- * Time-Space Trade off.
- * Specific Optimisation Techniques
- * Prephase - vs - Super optimisation



- + precise semantics
- + faster.
- + small footprint.
- + little Space

Semantic preserving changes improves
time / space characteristics

optimization → improvement.
only for "avg. expected cases".

Lvl of Optimisation.

High lvl. src \rightarrow AST.

Intermediate lvl. AST \rightarrow IR.

Low lvl. IR \rightarrow machine code.

Different lvl gives different possibilities.

Lvl of Optimisation

High lvl.

* type info.

Intermediate lvl.

* control flow / data flow analysis

Low lvl.

* Target code selection

* Register allocation

Time vs. Space.

Small is fast as well.

current data small L1 cache.

Common Techniques

Common Subexpression

$$a1 := b1 + c1 \cdot d.$$

$$a2 := b2 + c2 \cdot d.$$

Algebraic identities ! caveat.

Computer arithmetic does not follow mathematics
(floating point)

Dead code elimination.

* code that will not be executed.

* will never have any effect

* Analyse on control and data flow.

Constant Folding

evaluate sub expression at compile time.
reduce total number of instructions.

Propagation:

can be propagated through multiple variables.

$$\begin{array}{l} x = 10 \\ y = x + 1 \\ z = x + 2 \end{array}$$

↓ 10
↓ 11
↓ 12.

Constant folding: Caveat.

Be careful about semantics / precision.

Compile time precision should match. For the precision.

Strength reduction. use cheaper operators than expensive ones

Comparisons

Add / subtract / bit ops

mult.

divide / mod.

faster
↓ slower.
Int float

e.g. $\sum x \Rightarrow x < 2 + x$.

$x^2 \Rightarrow x + x$.

corner cases

$i = 10$.

while ($i > 0$) {

$i = i - 1$,

$t = 4 * i$,

$a[i] = b[t]$

}

$i = 10$

$t = 4 * 10$

while ($i > 0$) {

$i = i - 1$

$t = t - 4$

$a[i] = b[t]$

}

Inlining

avoid function call by copy and pasting function in when called.

caveats! : should be small funcs
careful with recursion.

Loop optimization.

* unrolling

- repeat the loop if ops are known
less jumps but grows code size.

* Fusion / Fission

fusion \Rightarrow put loops together fewer jumps

```
for (int i=1; i<100)  
    a[i] = i  
    b[i] = i+1  
    i++
```

fission \Rightarrow split loops up. multicore take a loop each

```
for (int i = 1; i < 100)  
    a[i] = i
```

```
for (int i = 1; i < 100)  
    b[i] = i+1
```

Code Motion. (Move loop independent code out of it.)

Pro: less instructions cont.

Cons: results in register spill.

Tiling (Blocking).

Create smaller loop around data to prevent register spill.
(Consine for tile).

```
for (i=0; i<N; i++)  
    j++
```

$\text{for } (j \geq 0 \text{ } j < N \text{ } j += B)$
 $\text{for } (i = j, i < \min(N, j + B) \text{ } i += 1)$
 - - -
 }.

Principle of Locality

- * Temporal locality reuse of data/resources within small time interval.
- * Spatial locality use. — — — within same storage location.
- * seq --- use — — — in a linear order.

Inversion.

convert a while loop into a do-while loop

reduces goto jumps

if --- initial check

```

do
|   ← change. / no need to jump
|   back and check
while () ← check here.
  
```

Interchanging.

change inner/outer loop for better locality.

$\text{for (int } i = \dots)$ $\text{for (j} \rightarrow)$
 $\text{for (int } j = \dots)$. \Rightarrow $\text{for (i} \rightarrow)$

Unswitching.

Take if-branching out by taking if outside of loops

Peephole Optimisation.

* only works on small portions of code.

* use better machine - - .

- * keep scanning until no more optimisations.
- * one peephole may bring another.

Find Result dependent on. order applied.

best order is hard to find.

Super-optimisation

exhaustive search over the space of all possible programs
to find the shortest program

→ used to find performance-critical inner loop optimisation

→ used to find peephole optimisation

Routine Organisation.

Machine Code. is a set of instructions understandable by CPUs

ASM: Readable low level.

Corresponds 1-1 machine code

Easily transferrable by assembler.

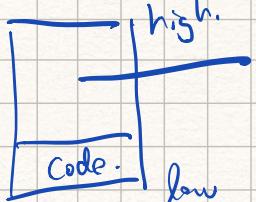
Routine organisation.

→ manage routine resources

→ correspondence between static (compile) - dynamic (run-time) structure

→ storage organisation

OS runs the program.
allocates the memory



1. Procedure Activation record.
2. Local data.
3. Global data.
4. Dynamic allocated data

An invocation of P is an activation of P

Similarly lifetime of var. of is the portion of execution in which x is defined (dynamic)

lifetime of an activation. * all the steps to execute P steps in procedure P calls

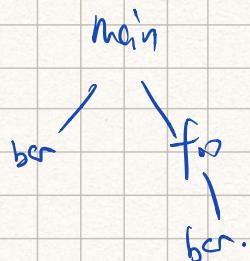
Activation Tree.

P → Q. Q should return before P.

lifetime of activation should be nested.

↓
tree.

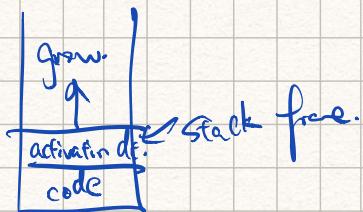
foo() { bar();
bar(); }
main() {
 bar();
 foo();
}



Activation tree dependent on runtime behavior.

may be different for diff. inputs.

can use a stack to track the active procedures



We need activation record/call frame to execute correctly

F()
G() ← call G.

F paused.

G's AR contains param -- info to execute G.

goes back to F() with return.

General AR.

- * Return / Results
- * Params
- * Control link.
- * Return addr. (where to go back to)
- * Temporaries (Temp values)
- * Local vars.
- * Access link. (address data needed by procedure)

Example of AR design.

- * Store enough info to execute correctly
- * compilers try to hold as much of frames as possible

Global data

Heap struct.

No global var in an AR

fixed addr assigned by compiler.

statically allocated.



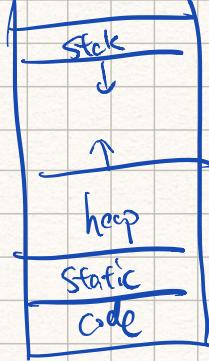
Global data/objects with fixed addr

a object. that survives the procedure

(cannot use AR)

S.t. Stored in heap struct. (malloc free).

Stack and heap.



both grow.
must make sure not grow into.

if grow into (out of memory)

Memory management.

garbage = value that will not be used anymore.
what is garbage? hard

Solutions

C/C++ style.

programmer deal with it.

memory leaks

dangling references

ptr1 ptr2.

double free bugs, freed (freed again)

GC.

GC: avoidance ~~uses~~.

but. const additional resource / performance.

Code Gen.

Conversion of IR \rightarrow machine instructions

Input : IR.

Output : absolute machine language.

relocatable machine language. o file.

assembly language. - S.

Code gen = Instruction Selection, + scheduling
+ register allocation.
+ debug data generation.

Three Main Architecture.

- * Register.
- * Stack \leftarrow operation between top of stack and other stacks
- * Accumulator \leftarrow where calc occur.

Instruction set.

ISA $\xrightarrow{\text{Software}}$ defines valid executions of hardware.
ISA $\xrightarrow{\text{hardware}}$

Architectures.

native DT.
Instructions
Registers
Addr modes
mem. architecture
interrupt & exception.
external I/O

$\xrightarrow{\text{RISC}}$ Reduced Instruction Set Computer

$x = y + z$
CISC MOV y, R0
 ADD z, R0
 MOV R0, x

$\xrightarrow{\text{RISC}}$

MOV y, R0
MOV z, R1
ADD R0, R1
MOV R1, x

IR \rightarrow Machine Need to know type of the architecture.

Goal: gen code that stores into registers as much as possible
much faster than memory.

Registers.

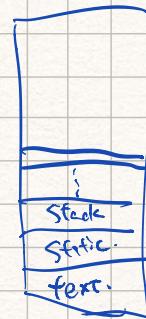
general registers. — data
segment registers // pointer
control registers index.

X86 Register Single Model.

Data Reg. EAX Accumulator
 ECX Counter
 Der. Reg. Stack Ptr. ESP
 Base Ptr. EBP
 Instruction Ptr. EIP

EBX Base.
 EDX Data.

Index Reg. SegIndxBSI
 Desc. Index EDI



X86 ASM Instructions

two operands.
 AT&T syntax (Gcc/GAS)
 Intel syntax.

(op src dst)
 op dst src.
 mov %eax \$1

X86 ASM Instructions

Arithmetic add sub inc mod idiv imul ...

at3.

data accumulator
 movl %eax %eax
 addl %eax %eax
 addl %eax %eax

Logic and or not zero.

Data mov : mov. movl \$5 (%rax)

Stack pushq. %rax popq %rax.

Comparison : cmp / test

Control : jump ! jg, jl, jf, jle, je.

DTs.

Suffix

b.	byte	byte 8 bit
s	short	16 bit
w	word	16 bit
l	long	32 bit 4 byte
ll	quad	64 bit 8 byte

4 byte offset.

movl 4(%rbp) %rax makes a quad.
 pushq %rbp

ref.

Access Static Global Data.

Static variables not in stack.

have fixed addresses throughout.

directly refer. using symbolic names

int g=5;

foo { int a=1 < g; }

: global -g.
: data.

-g: . long \$

movl -g %eax.

addl \$1 %eax.

Simple Code Gen.

Basic Block.

seq. of consecutive statements where flow ↓ ↑ without branching

TAC of $x = y + z \rightarrow \text{defn } x \text{ reference } y \text{ add } z$

1. determine leaders

R1 : first statement is a leader.

R2 : target of a jump is a leader.

R3 : immediately follows is a leader.

2. for a leader a basic block is all commands to the next

Example.

prod0
P1=1.

t1 = i L00 20

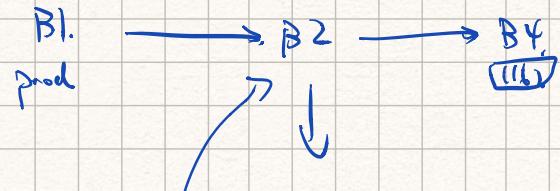
t2 = NOT t1.

cjmp t2 C16)

t3 = i

t4 = a[t3]

Flow Graphs. Rep. of basic blocks



$t5 := i$
 $t6 = b[t5]$
 $t7 = t4 * t6$
 $t8 = \text{prod} + t7$
 $\text{prod} = t8$
 $t9 = i + 1$
 $i = t9$
 $\text{jmp } B3$

B3
 t3
 jmp B2.

Next use analysis.

Generate code for a basic block.

- * For each var in a block, want to know if the value will be used again.
 - * Yes. Keep the register (var is live)
 - * No. deallocate. — dead

We don't know how we arrived.

cannot assume var in registers (relabel them!?)

don't know -- go.

store them in memory!

Next Use Algorithm : Computing next User

Init Assgn. :

all not temp or live
all temp or dead.

Scan over from beginning → end.
end → beginning

Suppose we find a TAC (i) $x = y$ op z during back scan.

attach(i) the info about nextuse of x y z .

x = not live.

y z = live nextuse = (i).

Single Code Gen alg.

from TAC.

Inside each block. one Stmt \rightarrow machine code.

register descriptor (what's in a register)

new register needed.

address descriptor (where is a var)

decide access method.

$$d = (a - b) + (a - c) + (a - d)$$

$$t1 = a - b$$

$$t2 = a - c$$

$$t3 = t1 + t2$$

$$d = t3 + t2.$$

MOV	a	R0	SUB	b	R0
MOV	a	R1	SUB	c	R1.
ADD	R0	R1.			
ADD	R1	R0			
MOV	R0	d.			

$$R0 \rightarrow \$ear$$

$$R1 \rightarrow \$ebr$$

$$a = -4\$ (rbp)$$

$$b = -8\$ (rbp).$$

$$c = -16\$ (rbp)$$

$$d = -20\$ (rbp)$$