

Code Generation

© Copyright 2005, Matthew B. Dwyer and Robby. The syllabus and lectures for this course are copyrighted materials and may not be used in other course settings outside University of Nebraska-Lincoln and Kansas State University in their current form or modified form without express written permission of one of the copyright holders. During this course, students are prohibited from selling notes to or being paid for taking notes by any person or commercial firm without the express written permission of one of the copyright holders.

Compiler Architecture



Code Generation

Issues

□ generating an internal representation of machine codes for source code abstractions:

- statements and expressions
- scopes and types
- procedures, parameters, return values
- control flow constructs etc.

□ runtime-system

optimizing the code (ignored for now); andemitting the code to files in binary format.

Code Generation

Not addressed now ..

□ low level issues

- register allocation
- multicore architectures
- • • •

functional / logic / domain-specific languages
 ...

Runtime system

- Compiler builds
 - an environment and assumes the program be executed in this environment
 - a set of procedures that are needed for program execution, e.g., for
 - interaction with the OS
 - data access and memory administration
 - error handling
 - uses virtual memory addresses

Runtime system

- Managing data accesses
 - Stack: Procedure (method) invocations
 - Heap: Garbage collection
 - o File system / IO
- Scheduling
 - o Concurrency
 - Non-determinism
- Monitoring
 - o Exceptions
 - o Debugging

Memory architecture



Control flows between procedure activations in LIFO (last-in-first-out) discipline (\Rightarrow stack principle)

procedure P(...)

Q(...); ... end;

procedure Q (...);

R(...); ... S(...); ... end;

Compiler

P Q Q R S Activation tree processed by pre-order traversal

Stack-like storage allocation for activations:















Activation records (AR)

Data needed within a procedure invocation is collected in an *activation record* or *stack frame*.

Control/Dynamic link

Static link

Saved state

Formal parameters

Function result

Local variables

Temporary storage

Reference to AR of caller

Ref. to AR of callee's static parent

In particular return-address









Determining static link for new activation record:

- → If m-n+1 = 0 use dynamic link; otherwise follow static link chain of caller for m-n+1 levels (performed at run time)
- m caller's static level (known at compile time)
- n callee's static level (known at compile time)



Activation records (AR)

Data needed within a procedure invocation is collected in an *activation record* or *stack frame*.

Control/Dynamic link

Static link

Saved state

Formal parameters

Function result

Local variables

Temporary storage

Reference to AR of caller

Ref. to AR of callee's static parent

In particular return-address

An important issue: Memory space needed for data?!

Memory layout: elementary data types

- Boolean (sometimes absent or identified with integers, e.g. C, LISP)
- Integer (sometimes divided into several types like in Java's integer-types byte, short, integer, long, (BigInteger))
- Real (floating standard IEEE 745, sometimes divided into several types like Java's float, double)
- Characters
- Enumerations (e.g.: Pascal declaration

type month = (jan, feb,mar,apr,may,jun,jul,aug,sep,oct,nov,dec); with implicit order jan < feb <...<dec ;

booleans and characters can be seen as predefined enumerations)

Memory layout: composite data types

- Arrays
- Records
- Sets
- Variant records
- Pointers

Memory layout: arrays

One-dimensional array: a: array[low..high] of T;

Layout in store



Memory layout: arrays

Memory layout of arrays:

Two-dimensional array: a: array $[l_1..h_1, l_2..h_2]$ of T;





Row-major array layout

Memory layout

Addr(a[i, j]) = base + ((i-l₁) * (h_2 -l₂+1) + (j-l₂)) * size(T) = base - (l₁ * d₂ + l₂) * size(T) + (i * d₂ + j) * size(T)

= $base_{red}$ + (i * d₂ + j) * size(T)

Memory layout

n-dimensional array: a: array[$I_1..h_1,...,I_n..h_n$] of T;

Let $d_i = h_i - l_i + 1$

$$Addr(a[i_{1},..,i_{n}]) = base + \left(\sum_{j=1}^{n} (i_{j}-l_{j}) *\prod_{k=j+1}^{n} d_{k}\right) * size(T)$$

= base - $\left(\sum_{j=1}^{n} l_{j} * \prod_{k=j+1}^{n} d_{k}\right) * size(T) + \left(\sum_{j=1}^{n} i_{j} * \prod_{k=j+1}^{n} d_{k}\right) * size(T)$
= base_{red} + $\left(\sum_{j=1}^{n} i_{j} * d^{(j)}\right) * size(T)$

Memory layout: arrays

Taxonomy of arrays (wrt. to allocation time, place):

- o Static arrays:
 - Allocated at compile time
 - (Example in C: static int a[10];)
 - Not included in PASCAL
 - Rapid address calculation (\rightarrow previous slides)

o Arrays with static bounds:

- Bounds known at compile-time
 - (Example in C: int a[10];
 - Example in PASCAL: a = array[0..10] of integer;)
- Stack-allocated at every procedure
- base unknown at compile time
- Fast address calculation

Memory layout: arrays

- o Dynamic arrays:
 - Bounds known only at run-time
 - Stack-allocated:
 - Size determined once for every procedure activation (Example in ALGOL 60: ARRAY A[0:N,0:N];)
 - Technical issue: Storage has to be allocated at the end of a stack frame
 - More flexible than arrays with static bounds, but address calculation gets slower
 - Heap-allocated:
 - Size determined at any allocation point (Example JAVA: int[] a;.....; a = new int[10];)
 - Very flexible but additional run-time overhead due to dereferencing

Memory layout: records







- Pointers are addresses in data memory (heap, stack)
- Only heap pointers in PASCAL, JAVA.
 Non-null pointers generated only via new-instructions.
- In C, C++ stack as well as heap pointers Explicit address-off operator &.
 Pointer arithmetics

Problems with manual pointer deallocation:

- ① **Dangling pointers**: Pointers to deallocated heap locations
- ② Memory leaks: Unaccessible heap locations

Problems with manual pointer deallocation:

- ① **Dangling pointers**: Pointers to deallocated heap locations
- ② Memory leaks: Unaccessible heap locations



Problems with manual pointer deallocation:

- ① Dangling pointers: Pointers to deallocated heap locations
- ② Memory leaks: Unaccessible heap locations



p := q;

Problems with manual pointer deallocation:

- ① Dangling pointers: Pointers to deallocated heap locations
- ② Memory leaks: Unaccessible heap locations

- The dangling pointer problem is avoided by using automatic memory deallocation only (*garbage collection*) like in JAVA.
- Memory leaks can also be reliably avoided by advanced garbage collection methods (not reference counting).

However, garbage collected programming languages are not well-suited for real-time sensitive applications.

C.A.R. Hoare (1973): "Their introduction into highlevel languages may be a step backward from which we may never recover."



Runtime system

- Managing data accesses
 - Stack: Procedure (method) invocations
 - Heap: Garbage collection
 - o File system / IO
- Scheduling
 - Concurrency
 - Non-determinism
- Monitoring
 - Exceptions
 - Debugging

Garbage collection

- Nature of object oriented programming: Almost everything is an object (Smalltalk, Java) and thus heap-located.
- Heap management is a central issue
- Heap-allocated objects that are not reachable from stack-located program variables (possibly through chains of pointers) are called garbage
Reference counting

- Keep track on the number of references to a heap allocated record. Reference count stored within each record.
- Garbage collector emits extra instructions.
 If pointer p \rightarrow r₁ is changed towards p \rightarrow r₂ perform:

1) Increment the reference count of r_2

2) Decrement the reference count of r_1 . If reference count of r_1 reaches zero, then r_1 is put on the freelist.

New heap allocated objects use locations from freelist.

Reference counting

Pros:

- Simple
- Smooth. Suited for real-time sensitive applications

Cons

- Cycles of garbage cannot be reclaimed → memory leakage (see example on the next slide)
- Expensive!

Even increment/decrement-operations take several machine instructions.

Naive implementation has to perform updates on every assignment

(some can be eliminated using static program analysis)

Reference counting



Reference counting



Mark and sweep

- Program variables and heap-allocated records form a directed graph.
- A graph-traversal algorithm like depth-first search (DFS) is used in order to mark all reachable nodes.
 (→ mark phase)
- Nodes that are not marked must be garbage and should be reclaimed (put on freelist). (→ sweep phase)
- New heap allocated objects use locations from freelist.

Mark and sweep 22 Stack Heap 37 a 2

Mark and sweep





Mark and sweep

Pros:

- No memory leakage
- Simple to implement

Cons

- Expensive!

Large portions of the heap have to be explored (Some improvements exist)

- Not smooth
- Overhead to manage freelist

- Essentially: Mark & sweep + compaction
 No freelist
 - Two separate heaps (*from-area*, *to-area*) switching roles
- Reachable records of from-area are copied into consecutive front segment of to-area
- Triggered by heap-limits

Copying collection



Compiler

Code Generation

Copying collection







Compiler

Code Generation

Copying collection



Code Generation









Copying collection

... and so on

Generational garbage collection

- Basic observation (empirically justified):

 Newly created objects are likely to die soon.
 Old objects will most likely survive many more collections.
- Divide the heap into generations G₀, G₁, G₂,...,G_k
 Often only 2 or 3 generations (nursery,...)

Generational garbage collection

Idea:

- $^{_{\rm D}}$ Each generation is exponentially larger than the one before, e.g.: $G_0\equiv \frac{1}{2}~$ MB, $G_1\equiv 2~$ MB , $G_2\equiv 8~$ MB,
- Objects are promoted to the next generation, if they survived two or three collections
- ^D G_0 is collected most often, G_0 together with G_1 more rarely, G_0 together with G_1 and G_2 even more rarely, etc.

Code Generation

Idea:

- Use intermediate code
 - Parse tree, AST, Symbol Table
 - Three-Address-Code
- Three-Address-Code
 - can be generated from AST
 - turning the tree representation into a linear sequence of instructions
 - use (at most) three *symbolic* addresses
 - Symbolic addresses are mapped to memory addresses or registers later

Code Generation

Three-Address-Code (TAC): Expressions a + a * (b-c) + b * d



TAC:

$$t_0 = b - c$$

 $t_1 = a * t_0$
 $t_2 = a + t_1$
 $t_3 = b * d$
 $t_4 = t_2 + t_3$

Code Generation

Three-Address-Code (TAC): Expressions a + a * (b-c) + b * d

TAC:

$$t_0 = b - c$$

 $t_1 = a * t_0$
 $t_2 = a + t_1$
 $t_3 = b * d$
 $t_4 = t_2 + t_3$

generated code, e.g.:

ор	dest	src	[src ₂]
LD	RO	b	
SUB	R0	R0	С
LD	R1	а	
MUL	R1	R1	R0
LD	R2	а	
ADD	R2	R2	R1
LD	R3	b	
MUL	R3	RЗ	d
ADD	R4	R2	R3

Input stream:

cost = (price + tax) * 6

Token stream:

<ID,1> <=> <(> <ID,2> <+> <ID,3> <)> <*> <NUMBER,4>

Symbol Table:

1	cost
2	price
3	tax
4	6

<ID,1> <=> <(> <ID,2> <+> <ID,3> <)> <*> <NUMBER,4>

=
<ID,1>
+
<NUMBER,4>
<ID,2>
<ID,3>

• Symbol Table:

1	cost	local	double
2	price	local	double
3	tax	local	double
4	6	constant	int

■ AST:

<ID,1> <=> <(> <ID,2> <+> <ID,3> <)> <*> <NUMBER,4>

modified AST:

Symbol Table:



1	cost	local	double
2	price	local	double
3	tax	local	double
4	6	constant	int

after type checking

<ID,1> <=> <(> <ID,2> <+> <ID,3> <)> <*> <NUMBER,4>



TAC:

- $t_1 = inttodouble(6)$ $t_2 = id_2 + id_3$ $t_3 = t_2 * t_1$ $id_1 = t_3$
- optimized TAC:

 $t_1 = id_2 + id_3$ $id_1 = t_1 * 6.0$

<ID,1> <=> <(> <ID,2> <+> <ID,3> <)> <*> <NUMBER,4>

optimized TAC
 generated code, e.g.:

$$t_1 = id_2 + id_3$$

 $id_1 = t_1 * 6.0$

LD R0 ID₂ ADD R0 R0 ID₃ MUL R0 R0 #6.0 ST ID₁ R0

Code generation: conditionals

if (cond) S_{then} else S_{else} // result in r_0 code (cond) if r₀=0 goto l_{else} $code(S_{then})$ goto l_{exit} l_{else}: code (S_{else}) l_{exit}:

Code generation: loops

while (cond) S_{body}

 $l_{loop}:$ code(cond) // result in r₀
if r₀=0 goto l_{exit}
code(S_{body})
goto l_{loop}
l_{exit}:

Code generation: switch-case



Code generation: switch-case

code(e)	// result in r0
goto l _{base} +r ₀	// computed jump
l _{base} :	
:	
goto l ₁	
:	// Jumptable
goto l _n	
goto l _{default}	
$l_1: code(S_1)$	
:	
l _n : code(S _n)	
l _{default} : code(S _d)	

Stack-based code generation: assignments

$$x = 2 * (x - y);$$

$$\frac{\text{Register-based}}{(\text{TAC})}$$

$$t1 = 2$$

(TA

t1

X

t2 = x-y

t1 = t1 * t2

= +1

Bootstrapping

- Modern compilers translate to target languages for which compilers or interpreters exist.
- Generation of machine code is unnecessary.
- Appropriate composition of existing compilers and/or interpreters with "simple-to-write translators" (if necessary)



T-diagrams

For compilers:



To be read as: Compiler from S to T written in I
T-diagrams

For interpreters:



To be read as: Interpreter for S written in I

Provided: P VM P VM VM

ETH Zürich Pascal portable compiler

(P = Pascal, VM = Pascal P-code (virtual machine))







Compiler

2. Compile interpreter:



3. Yields (interpretive) compiler:



4. Hand-written backend:



5. Bootstrap backend:



6. Bootstrap compiler:







