## Review for the Midterm

COMS W4115

Prof. Stephen A. Edwards

Spring 2007

Columbia University

Department of Computer Science

## The Midterm

- 70 minutes
- 4–5 problems
- Closed book
- One sheet of notes of your own devising
- Comprehensive: Anything discussed in class is fair game
- Little, if any, programming.
- Details of ANTLR/C/Java/Prolog/ML syntax not required
- Broad knowledge of languages discussed

## Topics

- Structure of a Compiler
- Scripting Languages
- Scanning and Parsing
- Regular Expressions
- Context-Free Grammars
- Top-down Parsing
- Bottom-up Parsing
- ASTs
- Name, Scope, and Bindings
- Control-flow

## Compiling a Simple Program

```
int gcd(int a, int b)
{
  while (a != b) {
    if (a > b) a -= b;
    else b -= a;
  }
  return a;
}
```

## What the Compiler Sees

```
int gcd(int a, int b)
{
  while (a != b) {
    if (a > b) a -= b;
    else b -= a;
  }
  return a;
}
```

```
i n t sp g c d ( i n t sp a , sp i
n t sp b ) nl { nl sp sp w h i l e sp
( a sp ! = sp b ) sp { nl sp sp sp i
f sp ( a sp > sp b ) sp a sp - = sp b
; nl sp sp sp sp e l s e sp b sp - = sp
a ; nl sp sp } nl sp sp r e t u r n sp
a ; nl } nl
```
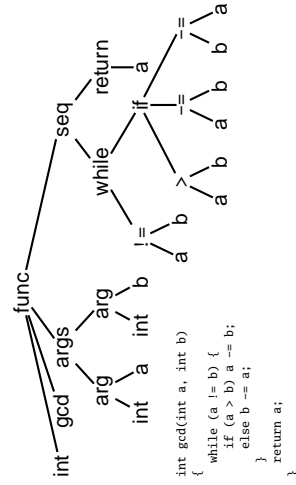
Text file is a sequence of characters

## Lexical Analysis Gives Tokens

```
int gcd(int a, int b)
{
  while (a != b) {
    if (a > b) a -= b;
    else b -= a;
  }
  return a;
}
```

Tokens: `int` `gcd` `(` `int` `a` `,` `int` `b` `)` `-` `while` `(` `a` `!=` `b` `)` `-` `if` `(` `a` `>` `b` `)` `a` `-=` `b` `;` `else` `b` `-=` `a` `;` `return` `a` `;`
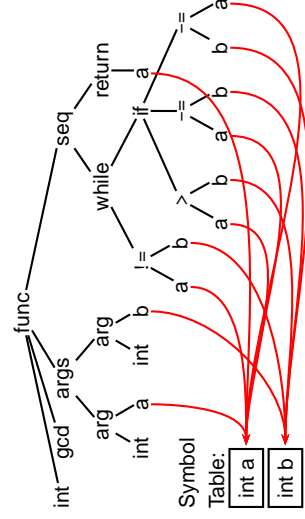
A stream of tokens. Whitespace, comments removed.

## Parsing Gives an AST

Abstract syntax tree built from parsing rules.

## Semantic Analysis Resolves Symbols

Symbol Table: `int a`  `int b`

Types checked; references to symbols resolved

## Translation into 3-Address Code

```
L0: sne   $1,  a, b
    seq   $0, $1, 0
    btrue $0, L1      % while (a != b)
    sl    $3,  b, a
    seq   $2, $3, 0
    btrue $2, L4      % if (a < b)
    sub   a,   a, b   % a -= b
    jmp   L5
L4: sub   b,   b, a   % b -= a
L5: jmp   L0
L1: ret   a
```

Idealized assembly language w/ infinite registers

# Generation of 80386 Assembly

```
gcd:    pushl  %ebp          % Save frame pointer
        movl   %esp,%ebp
        movl   8(%ebp),%eax   % Load a from stack
        movl   12(%ebp),%edx  % Load b from stack
.L8:    cmpl   %edx,%eax
        je     .L3            % while (a != b)
        jle    .L5            % if (a < b)
        subl   %edx,%eax      % a -= b
        jmp    .L8
.L5:    subl   %eax,%edx      % b -= a
        jmp    .L8
.L3:    leave                 % Restore SP, BP
        ret
```

# Describing Tokens

**Alphabet**: A finite set of symbols

Examples: $\{0, 1\}$, $\{A, B, C, \ldots, Z\}$, ASCII, Unicode

**String**: A finite sequence of symbols from an alphabet

Examples: $\epsilon$ (the empty string), Stephen, $\alpha\beta\gamma$

**Language**: A set of strings over an alphabet

Examples: $\emptyset$ (the empty language), $\{1, 11, 111, 1111\}$, all English words, strings that start with a letter followed by any sequence of letters and digits

# Scanning and Automata

# Operations on Languages

Let $L = \{\epsilon, \text{wo}\}$, $M = \{\text{man, men}\}$

**Concatenation**: Strings from one followed by the other

$LM = \{\text{man, men, woman, women}\}$
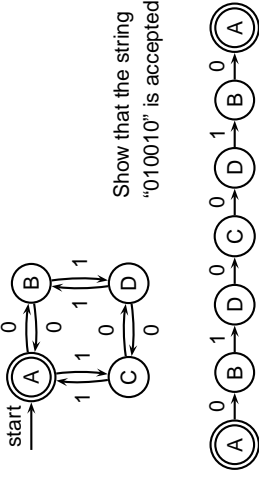
**Union**: All strings from each language

$L \cup M = \{\epsilon, \text{wo, man, men}\}$

**Kleene Closure**: Zero or more concatenations

$M^* = \{\epsilon, M, MM, MMM, \ldots\} =$
$\{\epsilon, \text{man, men, manman, manmen, menman, menmen,}$
$\text{manmanman, manmanmen, manmenman, \ldots}\}$

# Nondeterministic Finite Automata

"All strings containing an even number of 0's and 1's"



1. Set of states $S$: $\left\{ \text{A}, \text{B}, \text{C}, \text{D} \right\}$
2. Set of input symbols $\Sigma$: $\{0, 1\}$
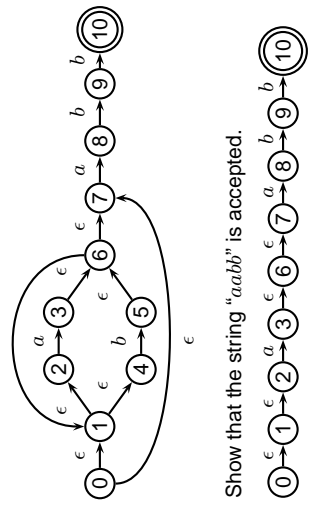3. Transition function $\sigma : S \times \Sigma_\epsilon \longrightarrow 2^S$

| state | $\epsilon$ | 0 | 1 |
|-------|------------|-----|-----|
| A | – | {B} | {C} |
| B | – | {A} | {D} |
| C | – | {D} | {A} |
| D | – | {C} | {B} |

4. Start state $s_0$ : A
5. Set of accepting states $F$: $\left\{ \text{A} \right\}$

# Regular Expressions over an Alphabet $\Sigma$

A standard way to express languages for tokens.

1. $\epsilon$ is a regular expression that denotes $\{\epsilon\}$
2. If $a \in \Sigma$, $a$ is an RE that denotes $\{a\}$
3. If $r$ and $s$ denote languages $L(r)$ and $L(s)$,

   • $(r)|(s)$ denotes $L(r) \cup L(s)$
   • $(r)(s)$ denotes $\{tu : t \in L(r), u \in L(s)\}$
   • $(r)^*$ denotes $\cup_{i=0}^{\infty} L^i$ $(L^0 = \emptyset$ and $L^i = LL^{i-1})$

# The Language induced by an NFA

An NFA accepts an input string $x$ iff there is a path from the start state to an accepting state that "spells out" $x$.



Show that the string "010010" is accepted.



# Translating REs into NFAs



# Translating REs into NFAs

Example: translate $(a|b)^*abb$ into an NFA



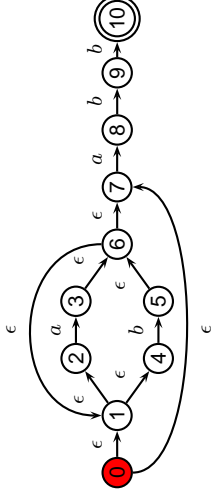Show that the string "$aabb$" is accepted.

# Simulating NFAs

Problem: you must follow the "right" arcs to show that a string is accepted. How do you know which arc is right?

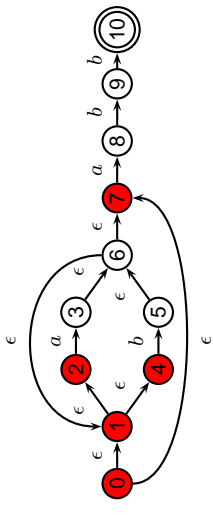Solution: follow them all and sort it out later.

"Two-stack" NFA simulation algorithm:

1. Initial states: the $\epsilon$-closure of the start state

2. For each character $c$,

   • New states: follow all transitions labeled $c$

   • Form the $\epsilon$-closure of the current states

3. Accept if any final state is accepting

# Simulating an NFA: $\cdot aabb$, Start

# Simulating an NFA: $\cdot aabb$, $\epsilon$-closure

# Simulating an NFA: $a\cdot abb$

# Simulating an NFA: $a\cdot abb$, $\epsilon$-closure

# Simulating an NFA: $aa\cdot bb$
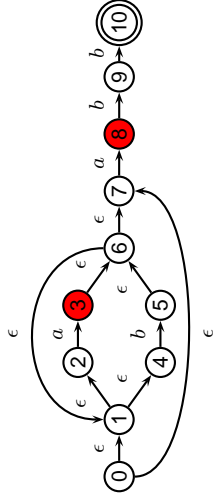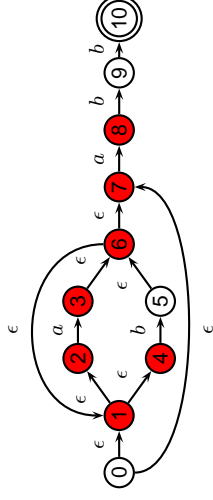
# Simulating an NFA: $aa\cdot bb$, $\epsilon$-closure
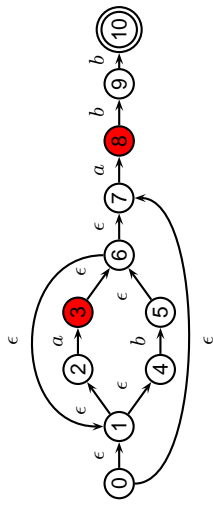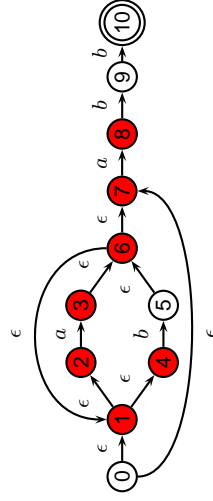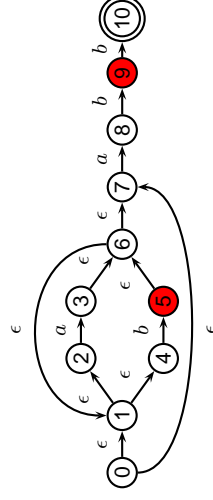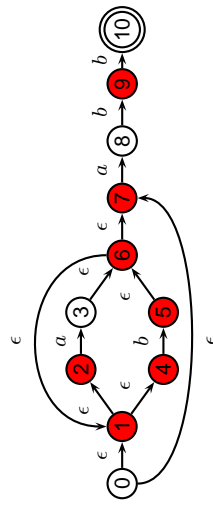
# Simulating an NFA: $aab\cdot b$

# Simulating an NFA: $aab\cdot b$, $\epsilon$-closure

## Simulating an NFA: $aabb\cdot$

## Simulating an NFA: $aabb$, Done

## Deterministic Finite Automata

Restricted form of NFAs:

- No state has a transition on $\epsilon$
- For each state $s$ and symbol $a$, there is at most one edge labeled $a$ leaving $s$.

Differs subtly from the definition used in COMS W3261 (Sipser, *Introduction to the Theory of Computation*)

Very easy to check acceptance: simulate by maintaining current state. Accept if you end up on an accepting state. Reject if you end on a non-accepting state or if there is no transition from the current state for the next symbol.

## Deterministic Finite Automata

```
ELSE:    "else" ;
ELSEIF:  "elseif" ;
```

## Deterministic Finite Automata

```
IF:  "if" ;
ID:  'a'..'z' ('a'..'z' | '0'..'9')* ;
NUM: ('0'..'9')+ ;
```

## Building a DFA from an NFA

Subset construction algorithm

Simulate the NFA for all possible inputs and track the states that appear.

Each unique state during simulation becomes a state in the DFA.

## Subset construction for $(a|b)^*abb$ (1)

## Subset construction for $(a|b)^*abb$ (2)

## Subset construction for $(a|b)^*abb$ (3)

# Grammars and Parsing

## Subset construction for $(a|b)^*abb$ (4)



## Ambiguous Grammars

A grammar can easily be ambiguous. Consider parsing

$$3 - 4 * 2 + 5$$

with the grammar

$$e \rightarrow e + e \,|\, e - e \,|\, e * e \,|\, e / e$$



## Fixing Ambiguous Grammars

Original ANTLR grammar specification

```
expr
  : expr '+' expr
  | expr '-' expr
  | expr '*' expr
  | expr '/' expr
  | NUMBER
  ;
```

Ambiguous: no precedence or associativity.

## Assigning Precedence Levels

Split into multiple rules, one per level

```
expr : expr '+' expr
     | expr '-' expr
     | term ;

term : term '*' term
     | term '/' term
     | atom ;

atom : NUMBER ;
```

Still ambiguous: associativity not defined

## Assigning Associativity

Make one side or the other the next level of precedence

```
expr : expr '+' term
     | expr '-' term
     | term ;

term : term '*' atom
     | term '/' atom
     | atom ;

atom : NUMBER ;
```

## A Top-Down Parser

```
stmt : 'if' expr 'then' expr
     | 'while' expr 'do' expr
     | expr ':=' expr ;

expr : NUMBER | '(' expr ')' ;

AST stmt() {
  switch (next-token) {
  case "if": match("if"); expr(); match("then"); expr();
  case "while" : match("while"); expr(); match("do"); expr();
  case NUMBER or "(" : expr(); match(":="); expr();
  }
}
```

## Writing LL(k) Grammars

Cannot have left-recursion

```
expr : expr '+' term  |  term ;
```
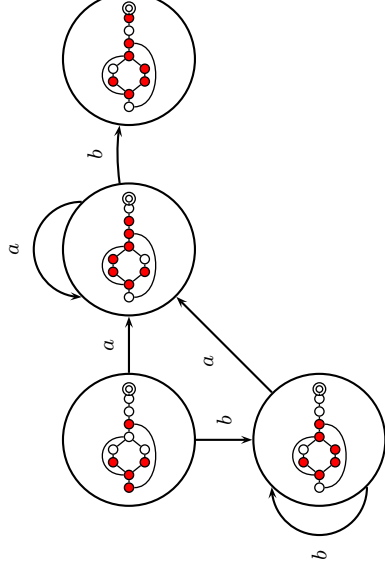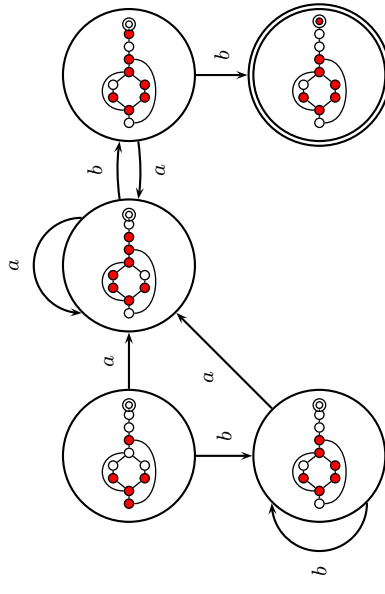
becomes

```
AST expr() {
  switch (next-token) {
  case NUMBER : expr(); /* Infinite Recursion */
}
```

## Writing LL(1) Grammars

Cannot have common prefixes

```
expr : ID '(' expr ')'
     | ID '=' expr
```

becomes

```
AST expr() {
  switch (next-token) {
  case ID : match(ID); match('('); expr(); match(')');
  case ID : match(ID); match('='); expr();
}
```

# Eliminating Common Prefixes

Consolidate common prefixes:

```
expr
    : expr '+' term
    | expr '-' term
    | term
    ;
```

becomes

```
expr
    : expr ('+' term | '-' term )
    | term
    ;
```

# Eliminating Left Recursion

Understand the recursion and add tail rules

```
expr
    : expr ('+' term | '-' term )
    | term
    ;
```

becomes

```
expr : term exprt ;
exprt : '+' term exprt
      | '-' term exprt
      | /* nothing */
      ;
```

# Bottom-up Parsing

# Rightmost Derivation

1 : $e \rightarrow t + e$
2 : $e \rightarrow t$
3 : $t \rightarrow \text{Id} * t$
4 : $t \rightarrow \text{Id}$

A rightmost derivation for **Id** * **Id** + **Id**:

$e$
$t + e$
$t + \text{Id}$
$\text{Id} * \text{Id} + \text{Id}$
$\text{Id} * \text{Id} + \text{Id}$

Basic idea of bottom-up parsing: construct this rightmost derivation backward.
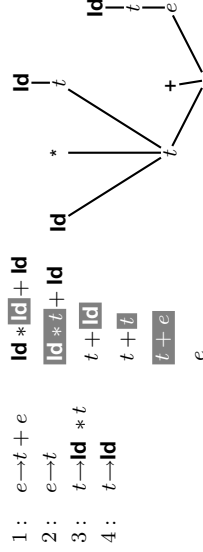
Here, I've drawn a box around each symbol to expand.

# Handles

1 : $e \rightarrow t + e$
2 : $e \rightarrow t$
3 : $t \rightarrow \text{Id} * t$
4 : $t \rightarrow \text{Id}$

$\text{Id} * \text{Id} + \text{Id}$
$\text{Id} * \text{Id} + \text{Id}$
$t + \text{Id}$
$t + t$
$t + e$
$e$

This is a reverse rightmost derivation for **Id** * **Id** + **Id**.

Each highlighted section is a handle.

Taken in order, the handles build the tree from the leaves to the root.

# Shift-reduce Parsing

1 : $e \rightarrow t + e$
2 : $e \rightarrow t$
3 : $t \rightarrow \text{Id} * t$
4 : $t \rightarrow \text{Id}$

| stack | input | action |
|---|---|---|
| | Id * Id + Id | shift |
| Id | * Id + Id | shift |
| Id* | Id + Id | shift |
| Id * Id | + Id | reduce (4) |
| Id * t | + Id | reduce (3) |
| t | + Id | shift |
| t+ | Id | shift |
| t + Id | | reduce (4) |
| t + t | | reduce (2) |
| t + e | | reduce (1) |
| e | | accept |

Scan input left-to-right, looking for handles.

An oracle tells what to do

# LR Parsing

1 : $e \rightarrow t + e$
2 : $e \rightarrow t$
3 : $t \rightarrow \text{Id} * t$
4 : $t \rightarrow \text{Id}$

| | action | | | | goto | |
|---|---|---|---|---|---|---|
| | **Id** | + | * | $ | e | t |
| 0 | s1 | | | | 7 | 2 |
| 1 | r4 | r4 | s3 | r4 | | |
| 2 | r2 | s4 | r2 | r2 | | |
| 3 | s1 | | | | | 5 |
| 4 | s1 | | | | 6 | 2 |
| 5 | r3 | r3 | r3 | r3 | | |
| 6 | r1 | r1 | r1 | r1 | | |
| 7 | | | | acc | | |

| stack | input | action |
|---|---|---|
| 0 | **Id** * **Id** + **Id** $ | shift, goto 1 |

1. Look at state on top of stack
2. and the next input token
3. to find the next action
4. In this case, shift the token onto the stack and go to state 1.

| stack | input | action |
|---|---|---|
| 0 | **Id** * **Id** + **Id** $ | shift, goto 1 |
| 0 Id[1] | * **Id** + **Id** $ | shift, goto 3 |
| 0 Id[1] *[3] | **Id** + **Id** $ | shift, goto 1 |
| 0 Id[1] *[3] Id[1] | + **Id** $ | reduce w/ 4 |

Action is reduce with rule 4 ($t \rightarrow$ **Id**). The right side is removed from the stack to reveal state 3. The goto table in state 3 tells us to go to state 5 when we reduce a $t$:

| stack | input | action |
|---|---|---|
| 0 Id[1] *[3] t[5] | + **Id** $ | |

# LR Parsing

1 : $e \rightarrow t + e$
2 : $e \rightarrow t$
3 : $t \rightarrow \text{Id} * t$
4 : $t \rightarrow \text{Id}$

| stack | input | action |
|---|---|---|
| 0 | **Id** * **Id** + **Id** $ | shift, goto 1 |
| 0 Id[1] | * **Id** + **Id** $ | shift, goto 3 |
| 0 Id[1] *[3] | **Id** + **Id** $ | shift, goto 1 |
| 0 Id[1] *[3] Id[1] | + **Id** $ | reduce w/ 4 |
| 0 Id[1] *[3] t[5] | + **Id** $ | reduce w/ 3 |
| 0 t[2] | + **Id** $ | shift, goto 4 |
| 0 t[2] +[4] | **Id** $ | shift, goto 1 |
| 0 t[2] +[4] Id[1] | $ | reduce w/ 4 |
| 0 t[2] +[4] t[2] | $ | reduce w/ 2 |
| 0 t[2] +[4] e[6] | $ | reduce w/ 1 |
| 0 e[7] | $ | accept |

| | action | | | | goto | |
|---|---|---|---|---|---|---|
| | **Id** | + | * | $ | e | t |
| 0 | s1 | | | | 7 | 2 |
| 1 | r4 | r4 | s3 | r4 | | |
| 2 | r2 | s4 | r2 | r2 | | |
| 3 | s1 | | | | | 5 |
| 4 | s1 | | | | 6 | 2 |
| 5 | r3 | r3 | r3 | r3 | | |
| 6 | r1 | r1 | r1 | r1 | | |
| 7 | | | | acc | | |

# Constructing the SLR Parse Table

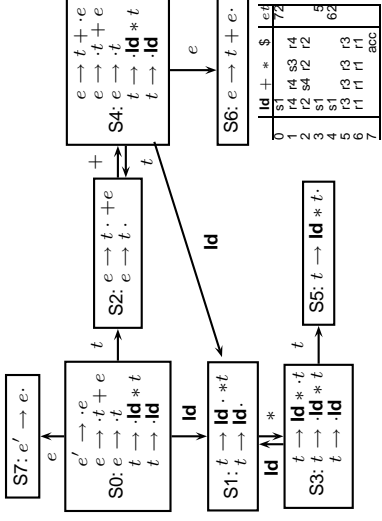The states are places we could be in a reverse-rightmost derivation. Let's represent such a place with a dot.

1 : $e \rightarrow t + e$
2 : $e \rightarrow t$
3 : $t \rightarrow \mathbf{Id} * t$
4 : $t \rightarrow \mathbf{Id}$

Say we were at the beginning ($\cdot e$). This corresponds to

$e' \rightarrow \cdot e$
$e \rightarrow \cdot t + e$
$e \rightarrow \cdot t$
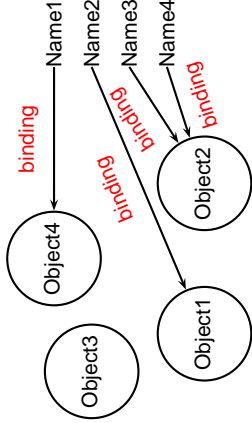$t \rightarrow \cdot \mathbf{Id} * t$
$t \rightarrow \cdot \mathbf{Id}$

The first is a placeholder. The second are the two possibilities when we're just before $e$. The last two are the two possibilities when we're just before $t$.
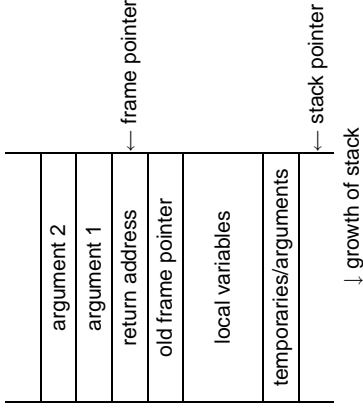
# Constructing the SLR Parsing Table



| | Id | + | * | \$ | e | t |
|---|---|---|---|---|---|---|
| 0 | s1 | | | | 7 | 2 |
| 1 | r4 | s3 | r4 | | | |
| 2 | r2 | s4 | r2 | | | |
| 3 | s1 | | | | | 5 |
| 4 | s1 | | | | 6 | 2 |
| 5 | r3 | r3 | r3 | r3 | | |
| 6 | r1 | r1 | r1 | r1 | | |
| 7 | | | | acc | | |

# Names, Objects, and Bindings



# Activation Records



argument 2
argument 1
$\leftarrow$ frame pointer
return address
old frame pointer
local variables
temporaries/arguments
$\leftarrow$ stack pointer

$\downarrow$ growth of stack

# Nested Subroutines in Pascal

```
procedure A;
  var a : integer;
  procedure B;
    var b : integer;
    procedure C;
      var c : integer;
      begin .. end
    procedure D;
      var d : integer;
      begin
        C;
      end;
  begin { Body of B }
    D;
  end;
  procedure E;
    var e : integer;
    begin
      B;
    end;
begin { Body of A }
  E;
end;
```



# Symbol Tables in a Functional Lang.



```
let
  var n := 8
  var x := 3
  function sqr(a:int)
    = a * a
  type ia = array of int
in
  n := sqr(x)
end
```

# Names, Objects, and Bindings

# Activation Records

```
int A() {
  int x;
  B();
}

int B() {
  int y;
  C();
}

int C() {
  int z;
}
```



# Control-Flow

# Side-effects

```
int x = 0;

int foo() { x += 5; return x; }

int a = foo() + x + foo();
```

GCC sets a=25.

Sun's C compiler gave a=20.

C says expression evaluation order is implementation-dependent.

# Multi-way Branching

```
switch (s) {
case 1: one(); break;
case 2: two(); break;
case 3: three(); break;
case 4: four(); break;
}
```

Switch sends control to one of the case labels. Break terminates the statement.

# Applicative- and Normal-Order Evaluation

```
int p(int i) { printf("%d ", i); return i; }
void q(int a, int b, int c)
{
    int total = a;
    printf("%d ", b);
    total += c;
}
q( p(1), 2, p(3) );
```

Applicative: arguments evaluated before function is called.

Result: 1 3 2

Normal: arguments evaluated when used.

Result: 1 2 3

# Misbehaving Floating-Point Numbers

1e20 + 1e-20 = 1e20

1e-20 $\ll$ 1e20

$(1 + 9e\text{-}7) + 9e\text{-}7 + 1 + (9e\text{-}7 + 9e\text{-}7)$

$9e\text{-}7 \ll 1$, so it is discarded, however, 1.8e-6 is large enough

$1.00001(1.000001 - 1) \neq 1.00001 \cdot 1.000001 - 1.00001 \cdot 1$

$1.00001 \cdot 1.000001 = 1.0000110001$ requires too much intermediate precision.

# Implementing multi-way branches

Obvious way:

```
if (s == 1) { one(); }
else if (s == 2) { two(); }
else if (s == 3) { three(); }
else if (s == 4) { four(); }
```

Reasonable, but we can sometimes do better.

# Nondeterminism

Nondeterminism is not the same as random:

Compiler usually chooses an order when generating code.

Optimization, exact expressions, or run-time values may affect behavior.

Bottom line: don't know what code will do, but often know set of possibilities.

```
int p(int i) { printf("%d ", i); return i; }
int q(int a, int b, int c) {}
q( p(1), p(2), p(3) );
```

Will *not* print 5 6 7. It will print one of

1 2 3, 1 3 2, 2 1 3, 2 3 1, 3 1 2, 3 2 1

# Gotos vs. Structured Programming

Break and continue leave loops prematurely:

```
for ( i = 0 ; i < 10 ; i++ ) {
    if ( i == 5 ) continue;
    if ( i == 8 ) break;
    printf("%d\n", i);
}

Again: if ( !(i < 10)) goto Break;
    if ( i == 5 ) goto Continue;
    if ( i == 8 ) goto Break;
    printf("%d\n", i);
Continue: i++; goto Again;
Break:
```

# Implementing multi-way branches

If the cases are *dense*, a branch table is more efficient:

```
switch (s) {
case 1: one(); break;
case 2: two(); break;
case 3: three(); break;
case 4: four(); break;
}

labels l[] = { L1, L2, L3, L4 }; /* Array of labels */
if (s>=1 && s<=4) goto l[s-1];   /* not legal C */
L1: one(); goto Break;
L2: two(); goto Break;
L3: three(); goto Break;
L4: four(); goto Break;
Break:
```

# Implementing Inheritance

Simple: Add new fields to end of the object

Fields in base class always at same offset in derived class

Consequence: Derived classes can never remove fields

| C++ | Equivalent C |
|---|---|
| `class Shape {` | `struct Shape {` |
| `    double x, y;` | `    double x, y;` |
| `};` | `};` |
| `class Box : Shape {` | `struct Box {` |
| `    double h, w;` | `    double x, y;` |
| `};` | `    double h, w;` |
| | `};` |

## Virtual Functions

```
class Shape {
    virtual void draw();   // Invoked by object's class
};                          // not its compile-time type.
class Line : public Shape {
    void draw();
};
class Arc : public Shape {
    void draw();
};

Shape *s[10];
s[0] = new Line;
s[1] = new Arc;
s[0]->draw();   // Invoke Line::draw()
s[1]->draw();   // Invoke Arc::draw()
```

## Virtual Functions

The Trick: Add a "virtual table" pointer to each object.

```
struct A {                    A's Vtbl       B's Vtbl
    int x;
    virtual void Foo();       A::Foo         B::Foo
    virtual void Bar();       A::Bar         A::Bar
};                            a1             B::Baz
struct B : A {
    int y;                    vptr           b1
    virtual void Foo();       x
    virtual void Baz();       a2             vptr
};                                           x
                              vptr
A a1, a2; B b1;               x              y
```

## Virtual Functions

```
struct A {                     B's Vtbl
    int x;
    virtual void Foo();        B::Foo
    virtual void Bar()         A::Bar
    { do_something(); }        B::Baz
};
struct B : A {
    int y;                     *a
    virtual void Foo();        vptr
    virtual void Baz();        x
};                             y

A *a = new B;
a->Bar();
```
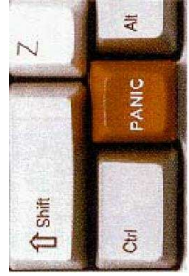
## Exceptions

A high-level replacement
for C's setjmp/longjmp.

```
struct Except { };

void baz()←──throw Except; }
void bar()←──baz(); }

void foo() {
    try {
        bar();
    } catch (Except e) {
        printf("oops");
    }
}
```

## One Way to Implement Exceptions

```
try {                push(Ex, Handler);

    throw Ex;        throw(Ex);
                     pop();
                     goto Exit;
} catch (Ex e) {     Handler:
    foo();               foo();
}                    Exit:
```

push() adds a handler to a stack

pop() removes a handler

throw() finds first matching handler

Problem: imposes overhead even with no exceptions

## Implementing Exceptions Cleverly

Real question is the nearest handler for a given PC.

| Lines | Action |
|-------|--------|
| 1–2   | Reraise |
| 3–5   | H1 |
| 6–9   | Reraise |
| 10–12 | H2 |
| 13–14 | Reraise |

```
1  void foo() {
2
3      try {                    3. look in table
4          bar();
5      } catch (Ex1 e) { H1: a(); }
6
7  }    2. H2 doesn't handle Ex1, reraise
8  void bar() {                 1. look in table
9
10     try {
11         throw Ex1();
12     } catch (Ex2 e) { H2: b(); }
13
14 }
```