

Review for the Final

COMS W4115

Prof. Stephen A. Edwards

Fall 2003

Columbia University

Department of Computer Science

The Final

Like 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 (1)

Structure of a Compiler

Scripting Languages

Scanning and Parsing

Regular Expressions

Context-Free Grammars

Top-down Parsing

Bottom-up Parsing

ASTs

Topics (2)

Name, Scope, and Bindings

Types

Control-flow constructs

Code Generation

Logic Programming: Prolog

Functional Programming: ML and the Lambda Calculus

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 s p g c d ( i n t s p a , s p i
n t s p b ) n l { n l s p s p w h i l e s p
( a s p ! = s p b ) s p { n l s p s p s p i
f s p ( a s p > s p b ) s p a s p - = s p b
; n l s p s p s p e l s e s p b s p - = s p
a ; n l s p s p } n l s p s p r e t u r n s p
a ; n l } n l
```

Text file is a sequence of characters

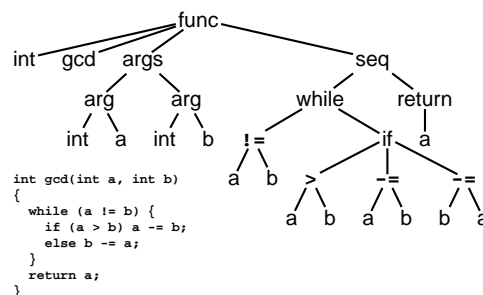
Lexical Analysis Gives Tokens

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

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

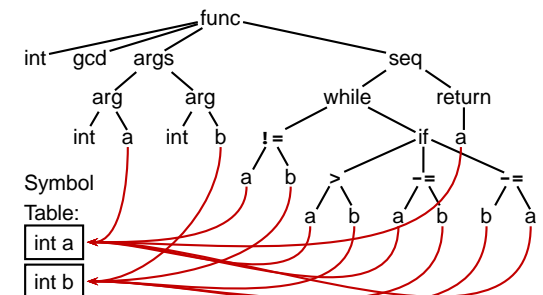
A stream of tokens. Whitespace, comments removed.

Parsing Gives an AST



Abstract syntax tree built from parsing rules.

Semantic Analysis Resolves Symbols



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)
    s1 $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
    
```

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

Scanning and Automata

Deterministic Finite Automata

A state machine with an initial state

Arcs indicate "consumed" input symbols.

States with double lines are accepting.

If the next token has an arc, follow the arc.

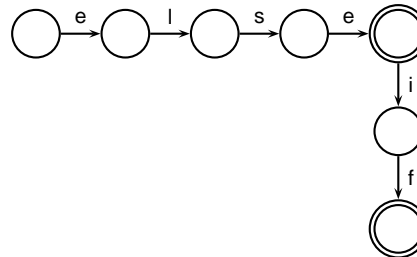
If the next token has no arc and the state is accepting, return the token.

If the next token has no arc and the state is not accepting, syntax error.

Deterministic Finite Automata

```

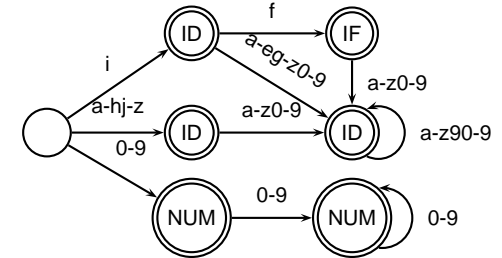
ELSE: "else" ;
ELSEIF: "elseif" ;
    
```



Deterministic Finite Automata

```

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



Nondeterministic Finite Automata

DFAs with ϵ arcs.

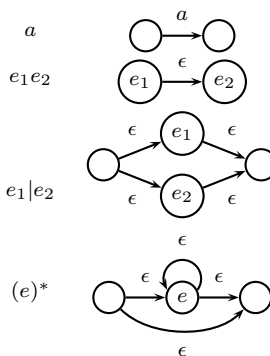
Conceptually, ϵ arcs denote state equivalence.

ϵ arcs add the ability to make nondeterministic (schizophrenic) choices.

When an NFA reaches a state with an ϵ arc, it moves to every destination.

NFAs can be in multiple states at once.

Translating REs into NFAs

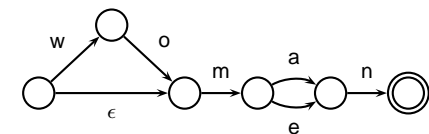


RE to NFAs

Building an NFA for the regular expression

$(wo|e)m(a|e)n$

produces

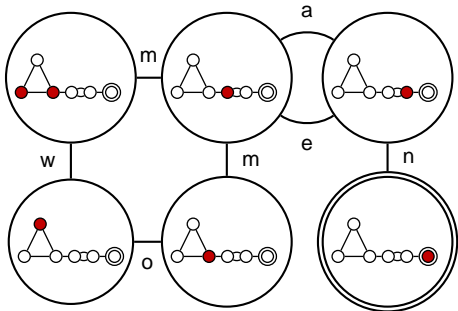


after simplification. Most ϵ arcs disappear.

Subset Construction

How to compute a DFA from an NFA.

Basic idea: each state of the DFA is a *marking* of the NFA



Subset Construction

An DFA can be exponentially larger than the corresponding NFA.

n states versus 2^n

Tools often try to strike a balance between the two representations.

ANTLR uses a different technique.

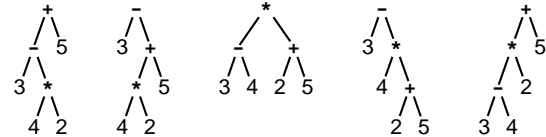
Ambiguous Grammars

A grammar can easily be ambiguous. Consider parsing

$3 - 4 * 2 + 5$

with the grammar

$e \rightarrow e + e \mid e - e \mid e * e \mid e / e$



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 ;
    
```

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 */
    
```


LR Parsing

1:	$e \rightarrow t + e$	stack	input	action
2:	$e \rightarrow t$		ld * ld + ld \$	shift, goto 1
3:	$t \rightarrow \mathbf{ld} * t$		* ld + ld \$	shift, goto 3
4:	$t \rightarrow \mathbf{ld}$		ld + ld \$	shift, goto 1
		action	goto	
		ld + * \$	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.

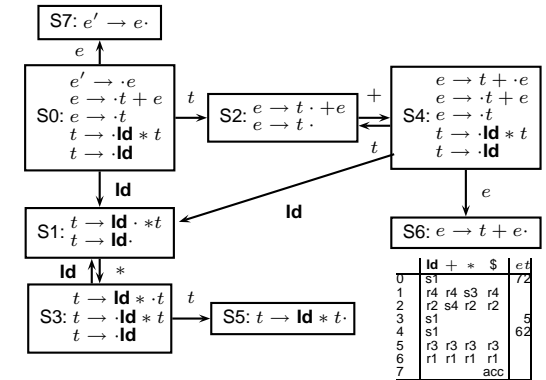
- $e \rightarrow t + e$
- $e \rightarrow t$
- $t \rightarrow \mathbf{ld} * t$
- $t \rightarrow \mathbf{ld}$

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

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

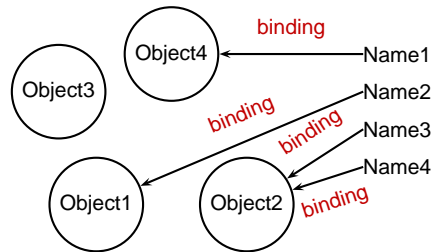
The first is a placeholder. The second are the two possibilities when we're just before e . The last two are the two possibilities when we're just before t .

Constructing the SLR Parsing Table

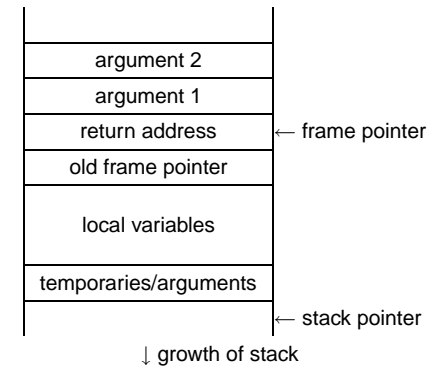


Names, Objects, and Bindings

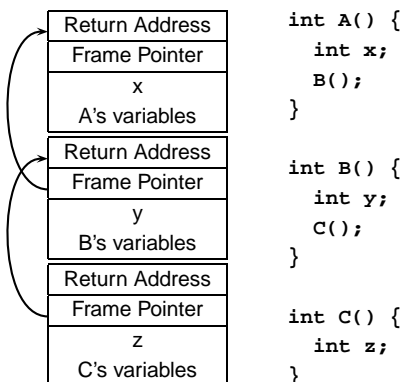
Names, Objects, and Bindings



Activation Records



Activation Records



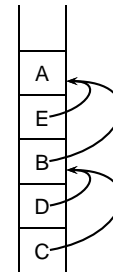
Nested Subroutines in Pascal

```

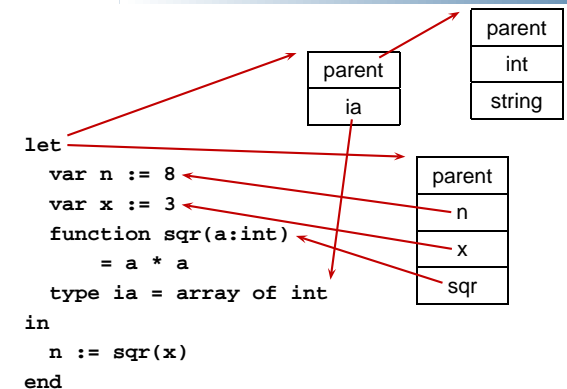
procedure A;
  procedure B;
    procedure C;
      begin .. end
  end B;

  procedure D;
    begin C end
  end D;

  procedure E;
    begin B end
  end E;
begin E end
    
```



Symbol Tables in Tiger



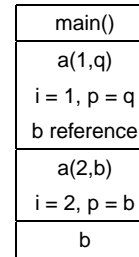
Shallow vs. Deep binding

```
typedef int (*ifunc)();
ifunc foo() {
    int a = 1;
    int bar() { return a; }
    return bar;
}
int main() {
    ifunc f = foo();
    int a = 2;
    return (*f)();
}
```

	static	dynamic
foo	1	2
main	1	1

Shallow vs. Deep binding

```
void a(int i, void (*p)()) {
    void b() { printf("%d", i); }
    if (i=1) a(2,b) else (*p)();
}
void q() {}
int main() {
    a(1,q);
}
```



	static
a	2
main	1

Layout of Records and Unions

Modern memory systems read data in 32-, 64-, or 128-bit chunks:

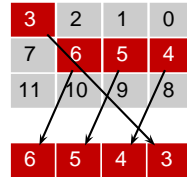


Reading an aligned 32-bit value is fast: a single operation.



Layout of Records and Unions

Slower to read an unaligned value: two reads plus shift.



SPARC prohibits unaligned accesses.

MIPS has special unaligned load/store instructions.

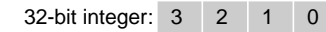
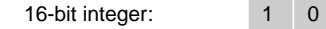
x86, 68k run more slowly with unaligned accesses.

Layout of Records and Unions

Modern processors have byte-addressable memory.



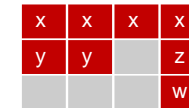
Many data types (integers, addresses, floating-point numbers) are wider than a byte.



Layout of Records and Unions

Most languages "pad" the layout of records to ensure alignment restrictions.

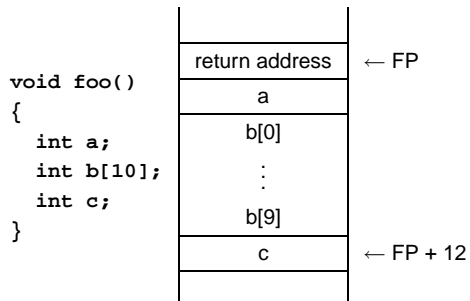
```
struct padded {
    int x; /* 4 bytes */
    char z; /* 1 byte */
    short y; /* 2 bytes */
    char w; /* 1 byte */
};
```



: Added padding

Allocating Fixed-Size Arrays

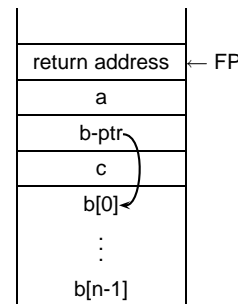
Local arrays with fixed size are easy to stack.



Allocating Variable-Sized Arrays

As always:
add a level of indirection

```
void foo(int n)
{
    int a;
    int b[n];
    int c;
}
```



Variables remain constant offset from frame pointer.

Static Semantic Analysis

Static Semantic Analysis

Lexical analysis: Make sure tokens are valid

```
if i 3 "This"           /* valid */
#a1123                  /* invalid */
```

Syntactic analysis: Makes sure tokens appear in correct order

```
for i := 1 to 5 do 1 + break /* valid */
if i 3                      /* invalid */
```

Semantic analysis: Makes sure program is consistent

```
let v := 3 in v + 8 end    /* valid */
let v := "f" in v(3) + v end /* invalid */
```

Implementing multi-way branches

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

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.

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

Static Semantic Analysis

Basic paradigm: recursively check AST nodes.

```
1 + break           1 - 5
```



```
check(+)           check(-)
check(1) = int     check(1) = int
check(break) = void check(5) = int
FAIL: int ≠ void   Types match, return int
```

Ask yourself: at a particular node type, what must be true?

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:
```

Applicative- vs. and Normal-Order

Most languages use applicative order.

Macro-like languages often use normal order.

```
#define p(x) (printf("%d ",x), x)
#define q(a,b,c) total = (a), \
    printf("%d ", (b)), \
    total += (c)
```

```
q( p(1), 2, p(3) );
```

Prints 1 2 3.

Some functional languages also use normal order evaluation to avoid doing work. "Lazy Evaluation"

Mid-test Loops

```
while true do begin
    readln(line);
    if all_blanks(line) then goto 100;
    consume_line(line);
end;
100:

LOOP
    line := ReadLine;
WHEN AllBlanks(line) EXIT;
    ConsumeLine(line)
END;
```

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;
}
```

What is printed by

```
q( p(1), 2, p(3) );
```

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

Currying

Functions are first-class objects that can be manipulated with abandon and treated just like numbers.

```
- fun max a b = if a > b then a else b;
val max = fn : int -> int -> int
- val max5 = max 5;
val max5 = fn : int -> int
- max5 4;
val it = 5 : int
- max5 6;
val it = 6 : int
-
```

Reduce

Another popular functional language construct:

```
fun reduce (f, z, nil) = z
  | reduce (f, z, h::t) = f(h, reduce(f, z, t));
```

If f is “-”, $\text{reduce}(f, z, a::b::c)$ is $a - (b - (c - z))$

```
- reduce( fn (x,y) => x - y, 0, [1,5]);
val it = ~4 : int
- reduce( fn (x,y) => x - y, 2, [10,2,1]);
val it = 7 : int
```

Pattern Matching

More fancy binding

```
fun map (_,[]) = []
  | map (f,h :: t) = f h :: map(f,t);
```

“_” matches anything

$h :: t$ matches a list, binding h to the head and t to the tail.

Recursion

ML doesn't have variables in the traditional sense, so you can't write programs with loops.

So use recursion:

```
- fun sum n =
=   if n = 0 then 0 else sum(n-1) + n;
val sum = fn : int -> int
- sum 2;
val it = 3 : int
- sum 3;
val it = 6 : int
- sum 4;
val it = 10 : int
```

But why always name functions?

```
- map( fn x => x + 5, [10,11,12]);
val it = [15,16,17] : int list
```

This is called a *lambda* expression: it's simply an unnamed function.

The `fun` operator is similar to a lambda expression:

```
- val add5 = fn x => x + 5;
val add5 = fn : int -> int
- add5 10;
val it = 15 : int
```

The Lambda Calculus

More recursive fun

```
- fun map (f, l) =
=   if null l then nil
=   else f (hd l) :: map(f, tl l);
val map = fn : ('a -> 'b) * 'a list -> 'b list

- fun add5 x = x + 5;
val add5 = fn : int -> int

- map(add5, [10,11,12]);
val it = [15,16,17] : int list
```

Pattern Matching

Functions are often defined over ranges

$$f(x) = \begin{cases} x & \text{if } x \geq 0 \\ -x & \text{otherwise.} \end{cases}$$

Functions in ML are no different. How to cleverly avoid writing if-then:

```
fun map (f,[]) = []
  | map (f,l) = f (hd l) :: map(f,tl l);
```

Pattern matching is order-sensitive. This gives an error.

```
fun map (f,l) = f (hd l) :: map(f,tl l)
  | map (f,[]) = [];
```

The Lambda Calculus

Fancy name for rules about how to represent and evaluate expressions with unnamed functions.

Theoretical underpinning of functional languages.
Side-effect free.

Very different from the Turing model of a store with evolving state.

ML:

```
fn x => 2 * x;
```

English:

The Lambda Calculus:

```
 $\lambda x. * 2 x$ 
```

“the function of x that returns the product of two and x ”

Evaluating Lambda Expressions

Pure lambda calculus has no built-in functions; we'll be impure.

To evaluate $(+ (* 5 6) (* 8 3))$, we can't start with $+$ because it only operates on numbers.

There are two *reducible expressions*: $(* 5 6)$ and $(* 8 3)$. We can reduce either one first. For example:

$(+ (* 5 6) (* 8 3))$	
$(+ 30 (* 8 3))$	
$(+ 30 24)$	Looks like deriving a
54	sentence from a grammar.

Evaluating Lambda Expressions

We need a reduction rule to handle λ s:

$$(\lambda x. * 2 x) 4$$
$$(* 2 4)$$
$$8$$

This is called β -reduction.

The formal parameter may be used several times:

$$(\lambda x. + x x) 4$$
$$(+ 4 4)$$
$$8$$

Reduction Order

The order in which you reduce things can matter.

$$(\lambda x. \lambda y. y) ((\lambda z. z z) (\lambda z. z z))$$

We could choose to reduce one of two things, either

$$(\lambda z. z z) (\lambda z. z z)$$

or the whole thing

$$(\lambda x. \lambda y. y) ((\lambda z. z z) (\lambda z. z z))$$

Reduction Order

Reducing $(\lambda z. z z) (\lambda z. z z)$ effectively does nothing because $(\lambda z. z z)$ is the function that calls its first argument on its first argument. The expression reduces to itself:

$$(\lambda z. z z) (\lambda z. z z)$$

So always reducing it does not terminate.

However, reducing the outermost function does terminate because it ignores its (nasty) argument:

$$(\lambda x. \lambda y. y) ((\lambda z. z z) (\lambda z. z z))$$
$$\lambda y. y$$

Reduction Order

The *redex* is a sub-expression that can be reduced.

The *leftmost* redex is the one whose λ is to the left of all other redexes. You can guess which is the *rightmost*.

The *outermost* redex is not contained in any other.

The *innermost* redex does not contain any other.

For $(\lambda x. \lambda y. y) ((\lambda z. z z) (\lambda z. z z))$,

$(\lambda z. z z) (\lambda z. z z)$ is the leftmost innermost and

$(\lambda x. \lambda y. y) ((\lambda z. z z) (\lambda z. z z))$ is the leftmost outermost.

Applicative vs. Normal Order

Applicative order reduction: Always reduce the leftmost **innermost** redex.

Normative order reduction: Always reduce the leftmost **outermost** redex.

For $(\lambda x. \lambda y. y) ((\lambda z. z z) (\lambda z. z z))$, applicative order reduction never terminated but normative order did.

Applicative vs. Normal Order

Applicative: reduce leftmost innermost

"evaluate arguments before the function itself"

eager evaluation, call-by-value, usually more efficient

Normative: reduce leftmost outermost

"evaluate the function before its arguments"

lazy evaluation, call-by-name, more costly to implement, accepts a larger class of programs

Normal Form

A lambda expression that cannot be reduced further is in *normal form*.

Thus,

$$\lambda y. y$$

is the normal form of

$$(\lambda x. \lambda y. y) ((\lambda z. z z) (\lambda z. z z))$$

Normal Form

Not everything has a normal form

$$(\lambda z. z z) (\lambda z. z z)$$

can only be reduced to itself, so it never produces a non-reducible expression.

"Infinite loop."

The Church-Rosser Theorems

If $E_1 \leftrightarrow E_2$ (are interconvertible), then there exists an E such that $E_1 \rightarrow E$ and $E_2 \rightarrow E$.

“Reduction in any way can eventually produce the same result.”

If $E_1 \rightarrow E_2$, and E_2 is in normal form, then there is a normal-order reduction of E_1 to E_2 .

“Normal-order reduction will always produce a normal form, if one exists.”

Church-Rosser

Amazing result:

Any way you choose to evaluate a lambda expression will produce the same result.

Each program means exactly one thing: its normal form.

The lambda calculus is deterministic w.r.t. the final result.

Normal order reduction is the most general.