

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

Review for the Midterm

COMS W4115
 Prof. Stephen A. Edwards
 Fall 2006
 Columbia University
 Department of Computer Science

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

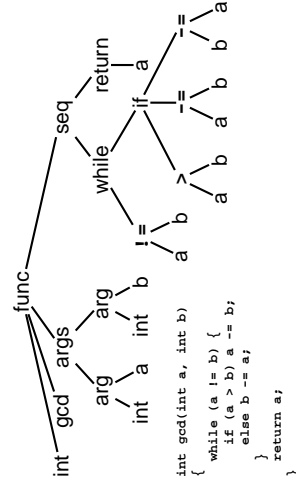
Lexical Analysis Gives Tokens

```
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 ; ;

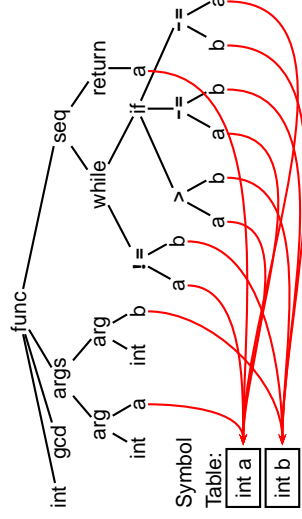
A stream of tokens. Whitespace, comments removed.

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
```

Idealized assembly language w/ infinite registers

Simulating NFAs

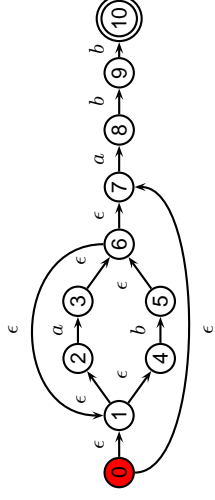
Problem: you must follow the "right" arcs to show that a string is accepted. How do you know which arc is right?

Solution: follow them all and sort it out later.

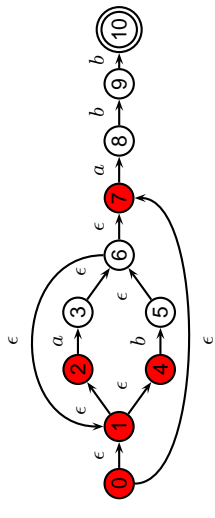
"Two-stack" NFA simulation algorithm:

1. Initial states: the ϵ -closure of the start state
2. For each character c :
 - New states: follow all transitions labeled c
 - Form the ϵ -closure of the current states
3. Accept if any final state is accepting

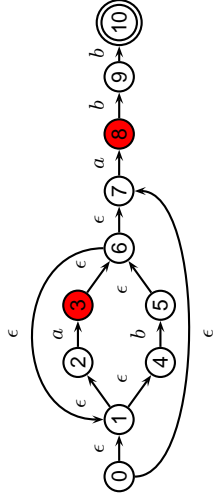
Simulating an NFA: $aabb$, Start



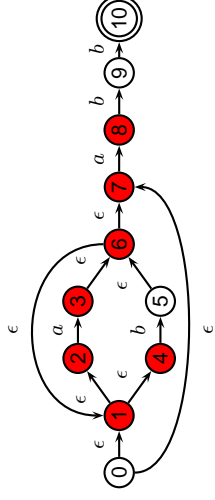
Simulating an NFA: $aabb$, ϵ -closure



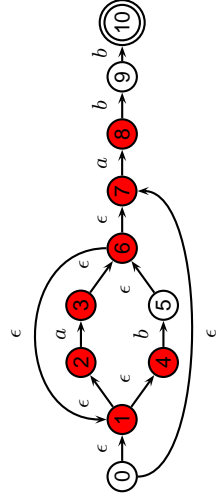
Simulating an NFA: $aabb$



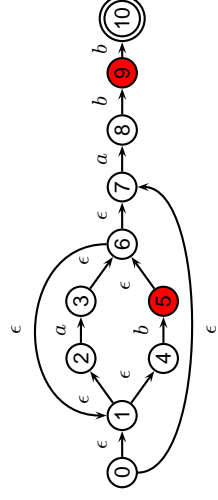
Simulating an NFA: $aabb$, ϵ -closure



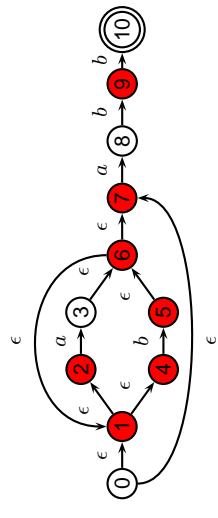
Simulating an NFA: $aabb$, ϵ -closure



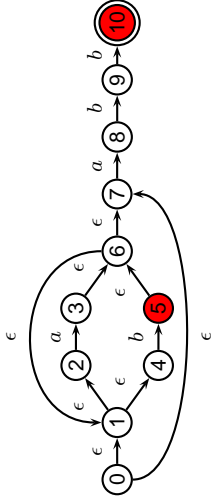
Simulating an NFA: $aabb$



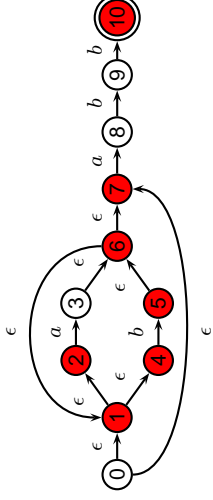
Simulating an NFA: $aabb$, ϵ -closure



Simulating an NFA: $aabb$.



Simulating an NFA: $aabb$, Done



Deterministic Finite Automata

Restricted form of NFAs:

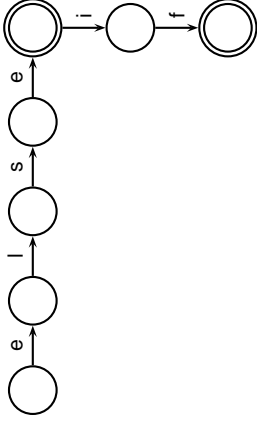
- No state has a transition on ϵ
- For each state s and symbol a , there is at most one edge labeled a leaving s .

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.

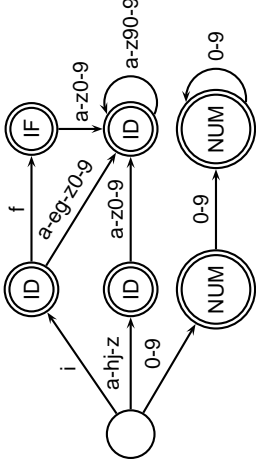
Deterministic Finite Automata

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

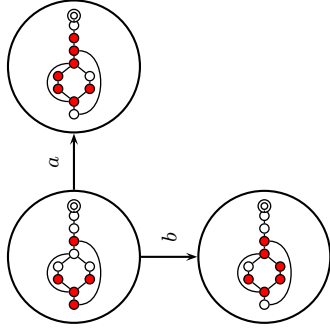


Deterministic Finite Automata

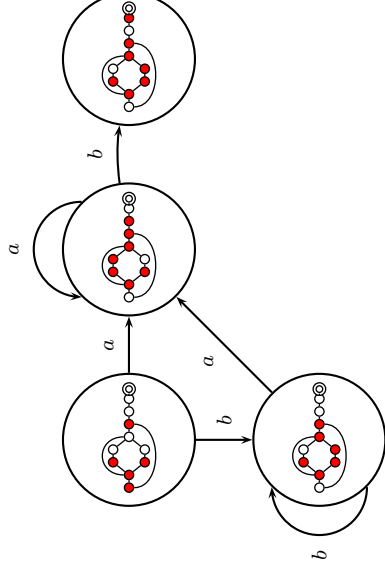
```
IF: "if" ;
ID: 'a'..'z' ('a'..'z' | '0'..'9')* ;
NUM: ('0'..'9')+ ;
```



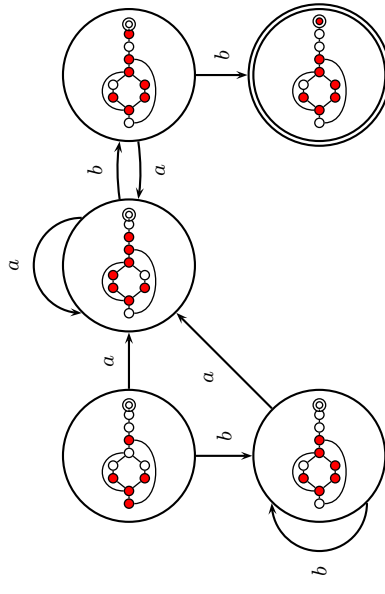
Subset construction for $(a|b)^*abb$ (1)



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



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

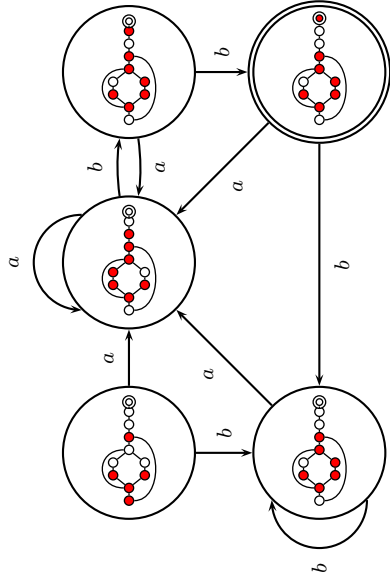


Subset construction algorithm

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 $(a|b)^*abb$ (4)



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

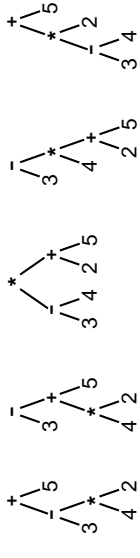
Ambiguous Grammars

A grammar can easily be ambiguous. Consider parsing

$3 - 4 * 2 + 5$

with the grammar

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



Grammars and Parsing

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();

```

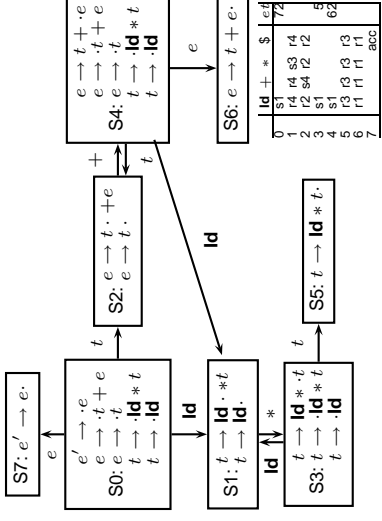

Constructing the SLR Parse Table

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: $t \rightarrow \text{ld} * t$
- 4: $t \rightarrow \text{ld}$

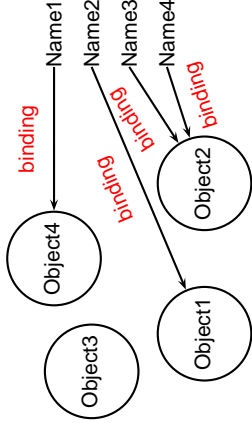
Say we were at the beginning ($t \cdot e$). This corresponds to $e' \rightarrow t \cdot e$. The first is a placeholder. The second are the two possibilities when we're just before e . The last two are the two possibilities when we're just before t .

Constructing the SLR Parsing Table

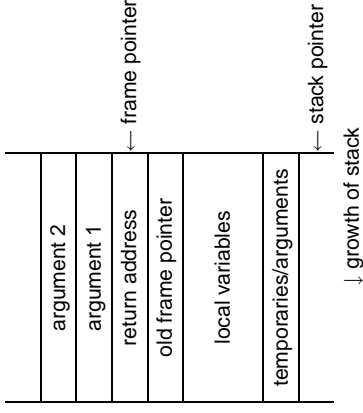


Names, Objects, and Bindings

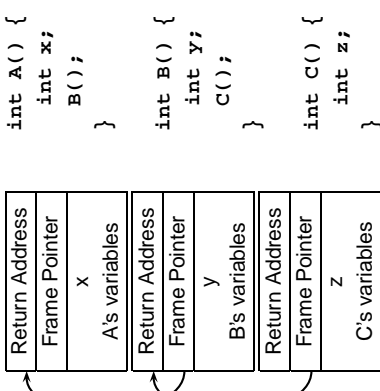
Names, Objects, and Bindings



Activation Records



Activation Records



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

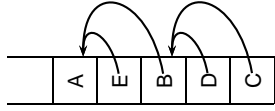
int B() {
  int y;
  C();
}

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

Nested Subroutines in Pascal

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

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



Symbol Tables in Tiger

