



## 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
            jmp  .L8
      .L3:  leave
            ret
    
```

## Operations on Languages

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

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

$LM = \{\text{man, men, woman, women}\}$

**Union:** All strings from each language

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

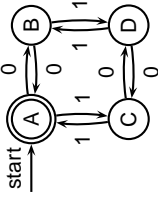
**Kleene Closure:** Zero or more concatenations

$M^* = \{\epsilon, M, MM, MMM, \dots\} =$

$\{\epsilon, \text{man, men, manman, manmen, manmen, manmen, manmanman, manmanmen, manmenman, manmanman, manmanman, manmanman, \dots}\}$

## 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.



## Describing Tokens

**Alphabet:** A finite set of symbols

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

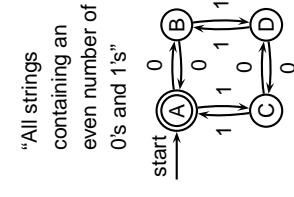
**String:** A finite sequence of symbols from an alphabet

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

## Nondeterministic Finite Automata



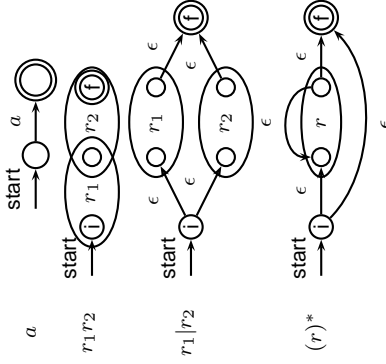
1. Set of states  $S: \{\textcircled{A}, \textcircled{B}, \textcircled{C}, \textcircled{D}\}$  containing an even number of 0's and 1's
2. Set of input symbols  $\Sigma: \{0, 1\}$
3. Transition function  $\sigma: S \times \Sigma_\epsilon \rightarrow 2^S$
4. Start state  $s_0: \textcircled{A}$
5. Set of accepting states  $F: \{\textcircled{C}\}$

## Regular Expressions over an Alphabet $\Sigma$

A standard way to express languages for tokens.

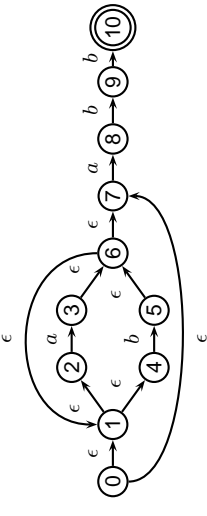
1.  $\epsilon$  is a regular expression that denotes  $\{\epsilon\}$
2. If  $a \in \Sigma$ ,  $a$  is an RE that denotes  $\{a\}$
3. If  $r$  and  $s$  denote languages  $L(r)$  and  $L(s)$ ,
  - $(r)|(s)$  denotes  $L(r) \cup L(s)$
  - $(r)(s)$  denotes  $\{tu : t \in L(r), u \in L(s)\}$
  - $(r)^*$  denotes  $\cup_{i=0}^{\infty} L^i$  ( $L^0 = \emptyset$  and  $L^i = LL^{i-1}$ )

## 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

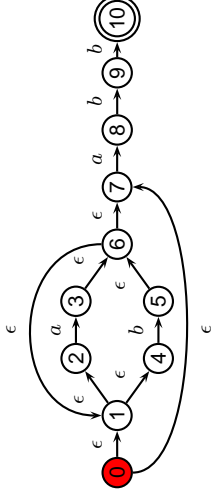
Problem: you must follow the "right" arcs to show that a string is accepted. How do you know which arc is right?

Solution: follow them all and sort it out later.

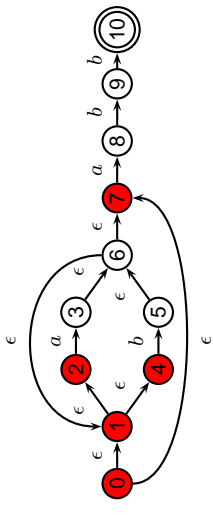
"Two-stack" NFA simulation algorithm:

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

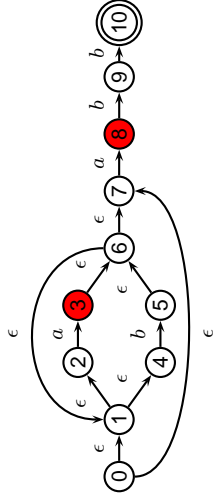
### Simulating an NFA: $aabb$ , Start



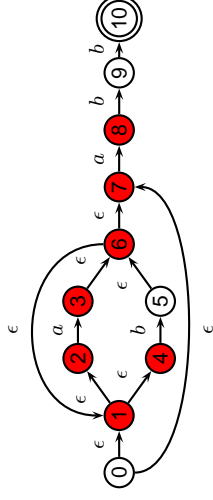
### Simulating an NFA: $aabb$ , $\epsilon$ -closure



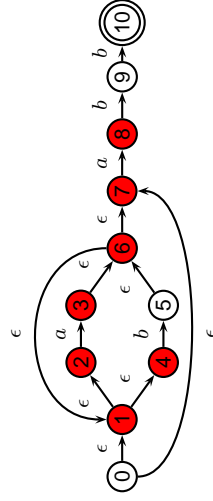
### Simulating an NFA: $aabb$



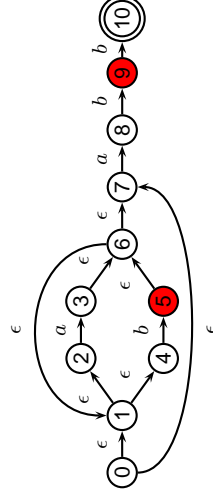
### Simulating an NFA: $aabb$ , $\epsilon$ -closure



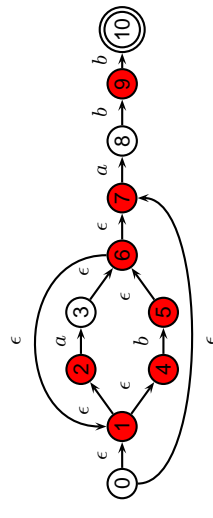
### Simulating an NFA: $aabb$ , $\epsilon$ -closure



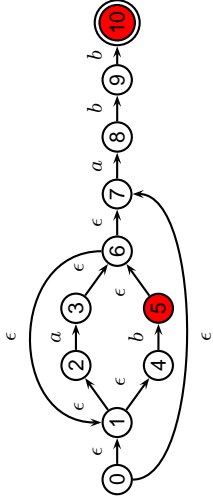
### Simulating an NFA: $aabb$



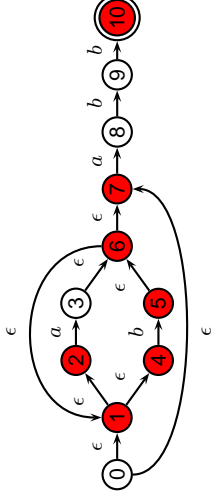
### Simulating an NFA: $aabb$ , $\epsilon$ -closure



## Simulating an NFA: $aabb$ .



## Simulating an NFA: $aabb$ ., Done



## Deterministic Finite Automata

Restricted form of NFAs:

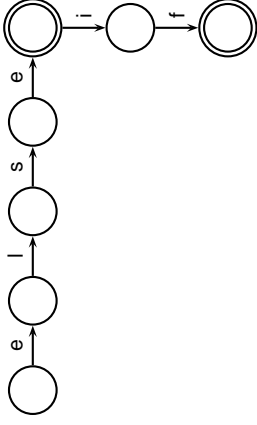
- No state has a transition on  $\epsilon$
- For each state  $s$  and symbol  $a$ , there is at most one edge labeled  $a$  leaving  $s$ .

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.

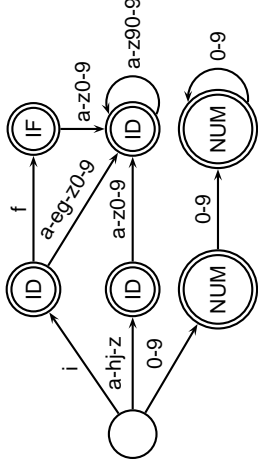
## Deterministic Finite Automata

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

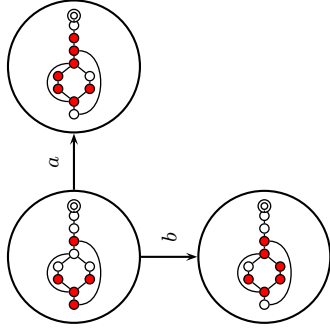


## Deterministic Finite Automata

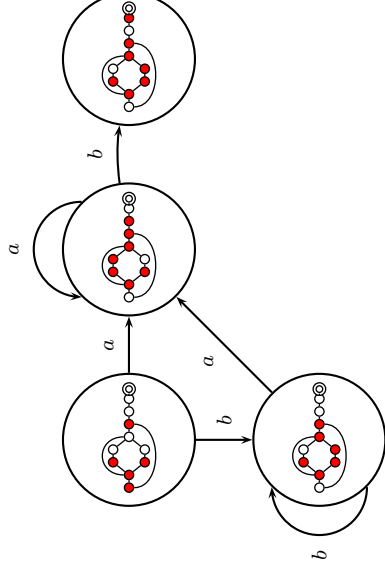
```
IF: "if" ;
ID: 'a'..'z' ('a'..'z' | '0'..'9')* ;
NUM: ('0'..'9')+ ;
```



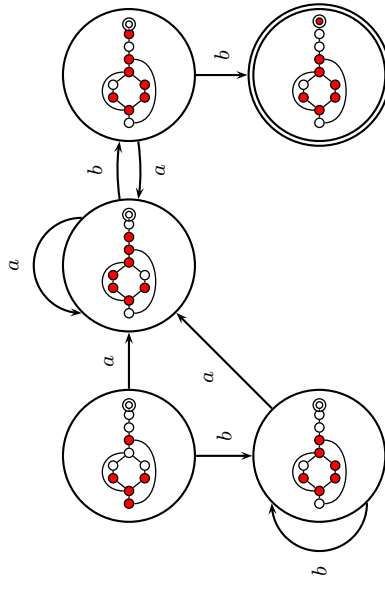
## Subset construction for $(a|b)^*abb$ (1)



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



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

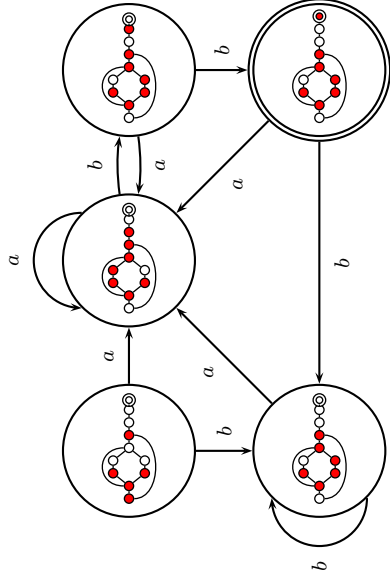


Subset construction algorithm

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.

## 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();
  }
}

```

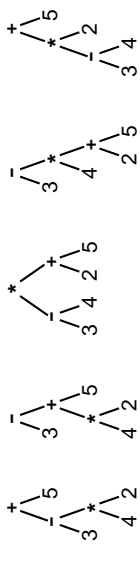
## Ambiguous Grammars

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 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

## Writing LL(k) Grammars

Cannot have left-recursion

```

expr : expr '+' term | term ;

```

becomes

```

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

```

## 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 ;

```

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

A rightmost derivation for  $\text{ld} * \text{ld} + \text{ld}$ :

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

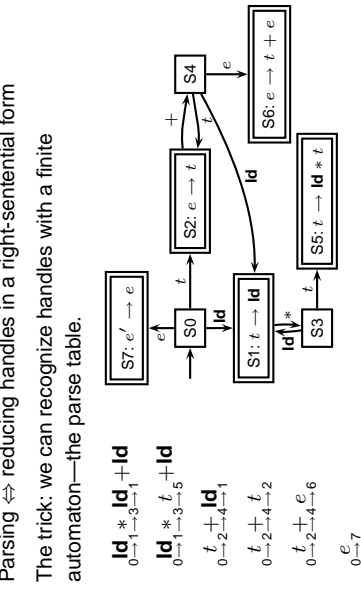
Basic idea of bottom-up parsing: construct this rightmost derivation **backward**.  
 The outlined terms are what we are expanding, *not handles*.

## Handles

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

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

This is a reverse rightmost derivation for  $\text{ld} * \text{ld} + \text{ld}$ .  
 Each highlighted section is a **handle**.  
 Taken in order, the handles build the tree from the leaves to the root.



## Handle Hunting

## Shift-reduce Parsing

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

Stack:  $\text{ld} * \text{ld} + \text{ld}$   
 Input:  $\text{ld} * \text{ld} + \text{ld}$   
 Action: shift, reduce (4), reduce (3), shift, reduce (2), reduce (1), 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:  $t \rightarrow \text{ld} * t$   
 4:  $t \rightarrow \text{ld}$

Stack:  $\text{ld} * \text{ld} + \text{ld} \$$   
 Input:  $\text{ld} * \text{ld} + \text{ld} \$$   
 Action: 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.

goto	$e$	$t$
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	

## LR Parsing

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

Stack:  $\text{ld} * \text{ld} + \text{ld} \$$   
 Input:  $\text{ld} * \text{ld} + \text{ld} \$$   
 Action: reduce w/ 4

Action is reduce with rule 4  
 ( $t \rightarrow \text{ld}$ ). 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
$\text{ld}$	$\text{ld} * \text{ld} + \text{ld} \$$	shift, goto 1
$\text{ld} *$	$\text{ld} + \text{ld} \$$	shift, goto 3
$\text{ld} * \text{ld}$	$+ \text{ld} \$$	shift, goto 1
$\text{ld} * \text{ld} +$	$\text{ld} \$$	reduce w/ 4

## 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	$\text{id} * \text{id} + \text{id} \$$	shift, goto 1
0 1	$* \text{id} + \text{id} \$$	shift, goto 3
0 1 3	$\text{id} + \text{id} \$$	shift, goto 1
0 1 3 1	$+ \text{id} \$$	reduce w/4
0 1 3 1 4	$+ \text{id} \$$	reduce w/3
0 1 3 1 4 3	$+ \text{id} \$$	shift, goto 4
0 1 3 1 4 3 4	$\text{id} \$$	shift, goto 1
0 1 3 1 4 3 4 1	$\$$	reduce w/4
0 1 3 1 4 3 4 1 4	$\$$	reduce w/2
0 1 3 1 4 3 4 1 4 2	$\$$	reduce w/1
0 1 3 1 4 3 4 1 4 2 6	$\$$	accept

goto	$e$	$t$
0	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	

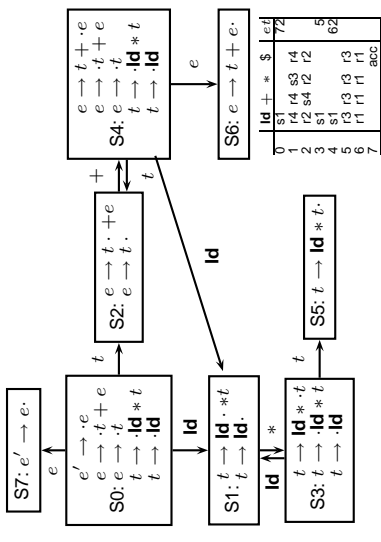
## Constructing the SLR Parse Table

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

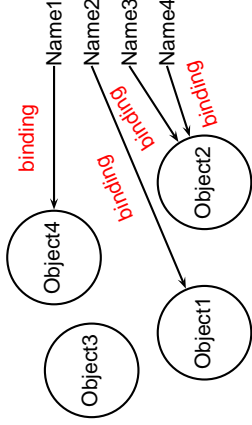
- $e' \rightarrow t + e$
- $e \rightarrow t$
- $t \rightarrow \text{id} * t$
- $t \rightarrow \text{id}$

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

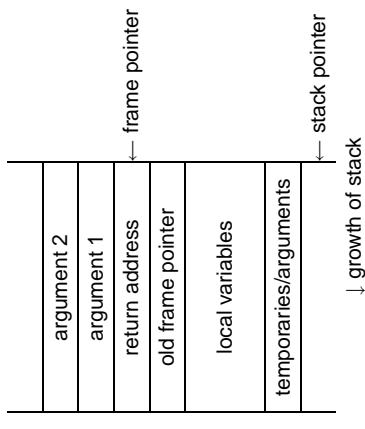
## Constructing the SLR Parsing Table



## Names, Objects, and Bindings



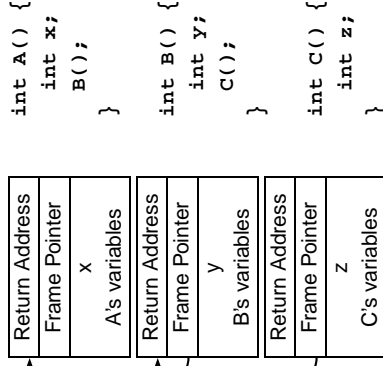
## Activation Records



## Names, Objects, and Bindings

### Bindings

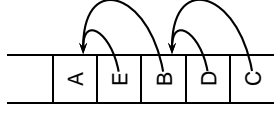
## Activation Records



## Nested Subroutines in Pascal

```

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



## Symbol Tables in Tiger

