

Lexical Analysis & Parsing (1)

15-411/15-611 Compiler Design

Seth Copen Goldstein

September, 27 2021

Today

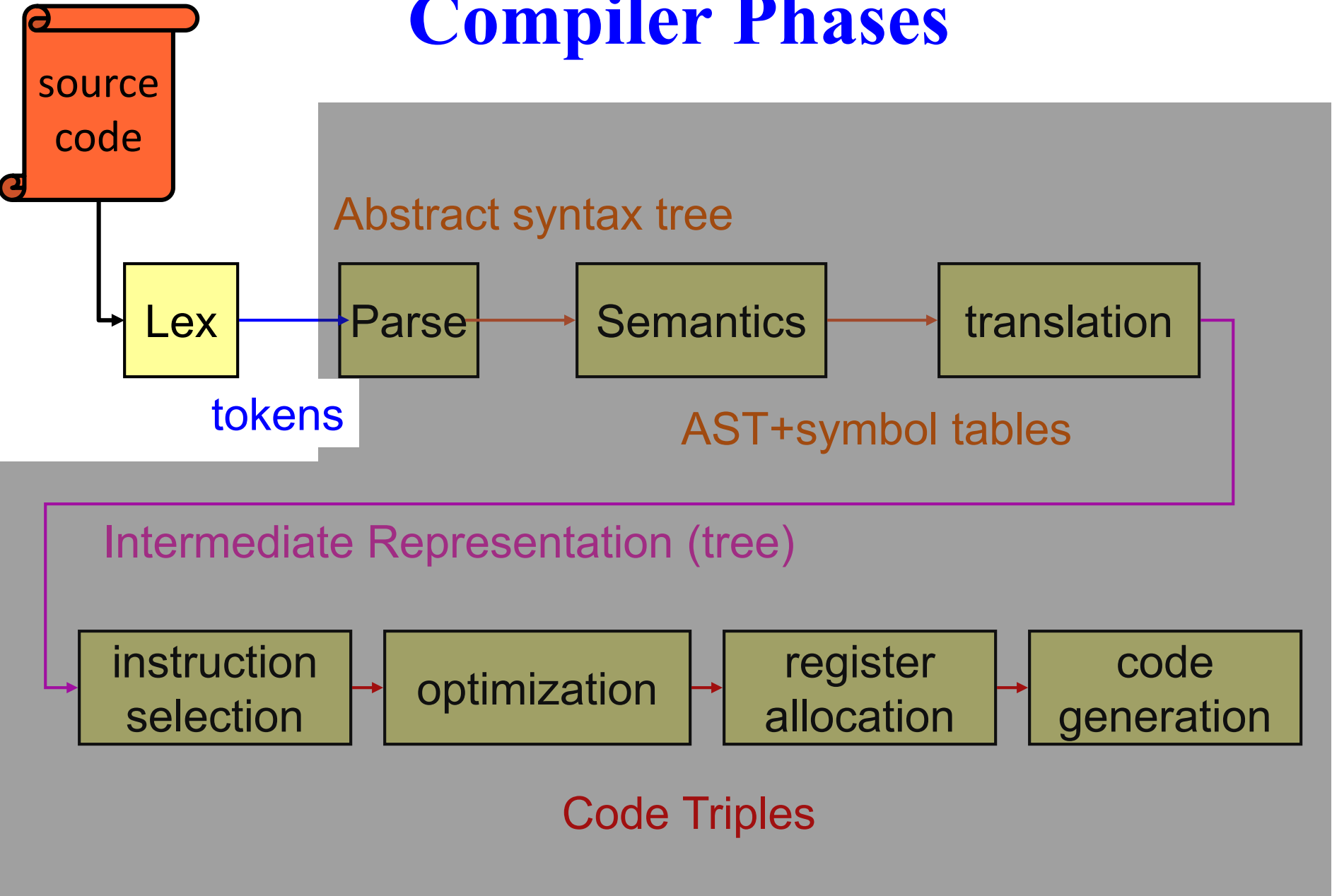
- Lexing
- Parsing

Today – part 1

Lexing

- Flex & other scanner generators
- Regular Expressions
- Finite Automata
- RE \rightarrow NFA
- NFA \rightarrow DFA
- DFA \rightarrow Minimized DFA
- Limits of Regular Languages

Compiler Phases



The Lexer

- Turn stream of characters into a stream of tokens

```
// create a user friendly descriptor for this arg.  
// if key is absent, then use it.  Otherwise use longkey  
  
char*  
ArgDesc::helpkey(WhichKey keytype, bool includebraks)  
{  
    static char buffer[128]; /* format buffer */  
    char* p = buffer;  
    ...  
}
```

```
CHAR STAR ID DOUBLE_COLON ID LPARIN ID ID COMMA BOOL ID  
RPARIN LBRACE STATIC CHAR ID LBRAK INTCONST RBRAK SEMI  
CHAR STAR ID EQ ID SEMI ...
```

The Lexer

- Turn stream of characters into a stream of tokens
 - Strips out “unnecessary characters”
 - comments
 - whitespace
 - Classify tokens by type
 - keywords
 - numbers
 - punctuation
 - identifiers
 - Track location
 - Associate with syntactic information

The Lexer

- Turn stream of characters into a stream of tokens

```
// create a user friendly descriptor for this arg.  
// if key is absent, then use it.  Otherwise use longkey  
  
char*  
ArgDesc::helpkey(WhichKey keytype, bool includebraks)  
{  
    static char buffer[128]; /* format buffer */  
    char* p = buffer;  
    ...  
}
```

```
CHAR STAR ID DOUBLE COLON ID LPARIN ID ID COMMA BOOL ID  
RPARIN LBRACE STATIC CHAR ID LBRAK INTCONST RBRAK SEMI  
CHAR STAR ID EQ ID SEMI ...
```

The Lexer

- Turn stream of characters into a stream of tokens

```
// create a user friendly descriptor for this arg.  
// if key is absent, then use it.  Otherwise use longkey  
  
char*  
ArgDesc::helpkey(WhichKey keytype, bool includebraks)  
{  
    static char buffer[128]; /* format buffer */  
    char* p = buffer;
```

Position: 4,0

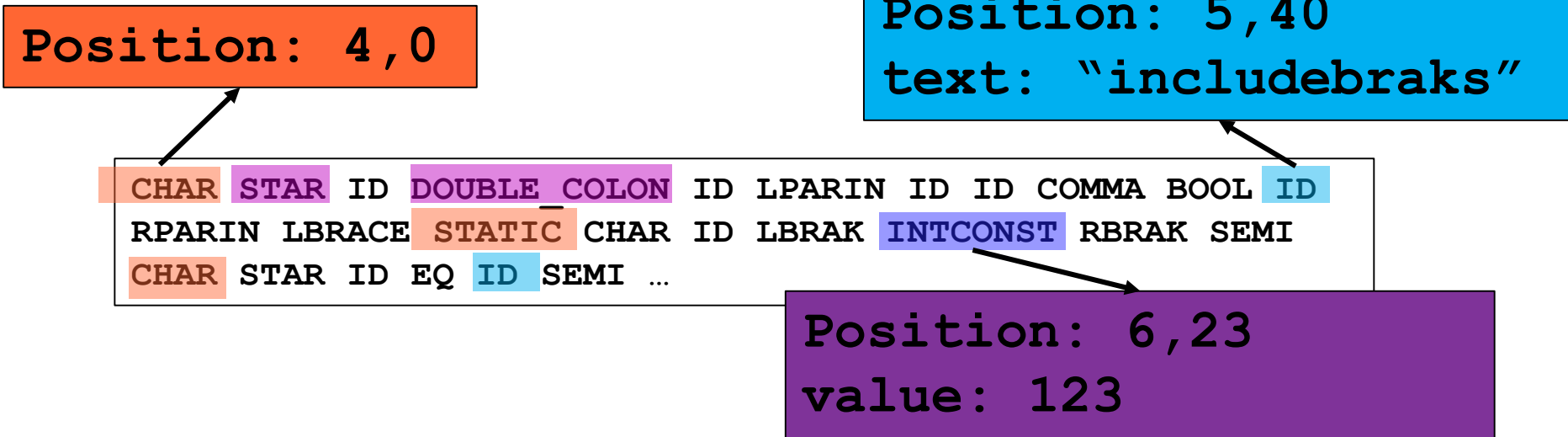
Position: 5,40
text: "includebraks"

```
CHAR STAR ID DOUBLE COLON ID LPARIN ID ID COMMA BOOL ID  
RPARIN LBRACE STATIC CHAR ID LBRAK INTCONST RBRAK SEMI  
CHAR STAR ID EQ ID SEMI ...
```

Position: 6,23
value: 123

The Lexer

- Turn stream of characters into a stream of tokens
 - More concise
 - Easier to parse

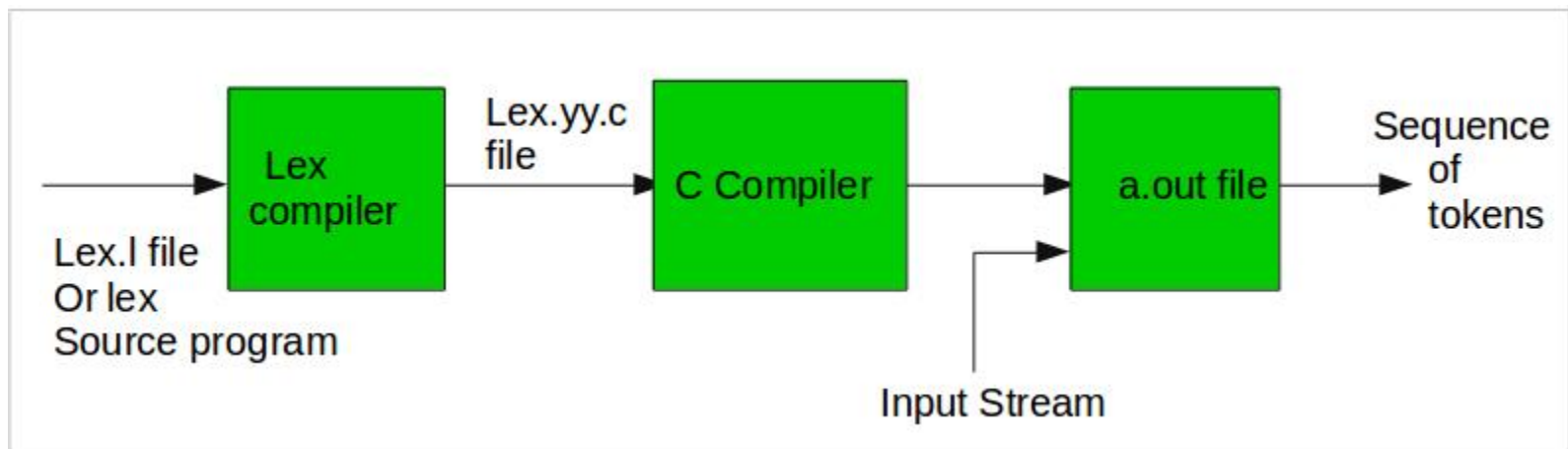


Lexical Analyzers

- Input: stream of characters
- Output: stream of tokens (with information)
- How to build?
 - By hand is tedious
 - Use Lexical Analyzer Generator, e.g., flex
- Define tokens with regular expressions
- Flex turns REs into Deterministic Finite Automata (DFA) which recognizes and returns tokens.

FLEX

- Define tokens
- Generate scanner code
- Main interface: `yylex()` which reads from `yyin` and returns tokens til EOF



2. Flex Program Format

- A flex program has three sections:

Definitions

%%

RE rules & actions

%%

User code

wc As a Flex Program

```
%{
    int charCount=0, wordCount=0, lineCount=0;
}%
word    [^ \t\n]+
%%
{word} {wordCount++; charCount += yyleng; }
[\\n]  {charCount++; lineCount++;}
.      {charCount++;}
%%
int main(void) {
    yylex();
    printf("Chars %d, Words: %d, Lines: %d\\n",
        charCount, wordCount, lineCount);
    return 0;
}
```

A Flex Program

```
%{  
    int charCount=0, wordCount=0, lineCount=0;  
}%  
word    [^ \t\n]+  
%%
```

```
{word} {wordCount++; charCount += yyleng; }  
[\n]   {charCount++; lineCount++;}  
.      {charCount++;}  
%%
```

```
int main(void) {  
    yylex();  
    printf("Chars %d, Words: %d, Lines: %d\n",  
          charCount, wordCount, lineCount);  
    return 0;  
}
```

1) Definitions

2) Rules & Actions

3) User Code

skip

Section 1: RE Definitions

- Format:

name RE

- Examples:

`digit` [0-9]

`letter` [A-Za-z]

`id` {letter} ({letter}|{digit})*

`word` [^\t\n]+

Regular Expressions in Flex

x	match the char x
\.	match the char .
"string"	match contents of string of chars
.	match any char except \n
^	match beginning of a line
\$	match the end of a line
[xyz]	match one char x , y , or z
[^xyz]	match any char except x , y , and z
[a-z]	match one of a to z

Regular Expressions in Flex (cont)

<code>r*</code>	closure (match 0 or more r's)
<code>r+</code>	positive closure (match 1 or more r's)
<code>r?</code>	optional (match 0 or 1 r)
<code>r1 r2</code>	match <i>r1</i> then <i>r2</i> (concatenation)
<code>r1 r2</code>	match <i>r1</i> or <i>r2</i> (union)
<code>(r)</code>	grouping
<code>r1 \ r2</code>	match <i>r1</i> when followed by <i>r2</i>
<code>{ <i>name</i> }</code>	match the RE defined by name

Some number REs

$[0-9]$

A single digit.

$[0-9]^+$

An integer.

$[0-9]^+ (\.[0-9]^+)?$ An integer or fp number.

$[+-]? [0-9]^+ (\.[0-9]^+)? ([eE][+-]?[0-9]^+)?$
Integer, fp, or scientific notation.

Section 2: RE/Action Rule

- A rule has the form:

```
name      { action }  
re        { action }
```

- the name must be defined in section 1
 - the action is any C code
-
- If the named RE matches* an input character sequence, then the C code is executed.

* Some caveats here

Rule Matching

- Longest match rule.

```
"int"      { return INT; }  
"integer"  { return INTEGER; }
```

- If rules can match same length input, first rule takes priority.

```
"int"      { return INT; }  
[a-z]+     { return ID; }  
[0-9]+     { return NUM; }
```

Section 3: C Functions

- Added to end of the lexical analyzer

Removing Whitespace

name → `whitespace`
`%%`
`{whitespace}`

RE → `▪`
`%%`

`[\t\n]`

empty action → `;`

`{ ECHO; }`

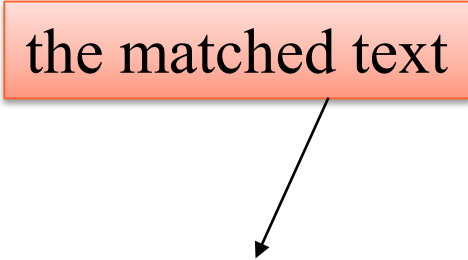
ECHO macro → `{ ECHO; }`

```
int main(void)
{
    yylex();
    return 0;
}
```

Printing Line Numbers

```
%{
    int lineno = 1;
}%
%%
^(.*)\n    { printf("%4d\t%s", lineno, yytext);
              lineno++;}
%%
int main(int argc, char *argv[])
{
    // proper arg processing & error handling, ...
    yyin = fopen(argv[1], "r");
    yylex();
    return 0;
}
```

the matched text

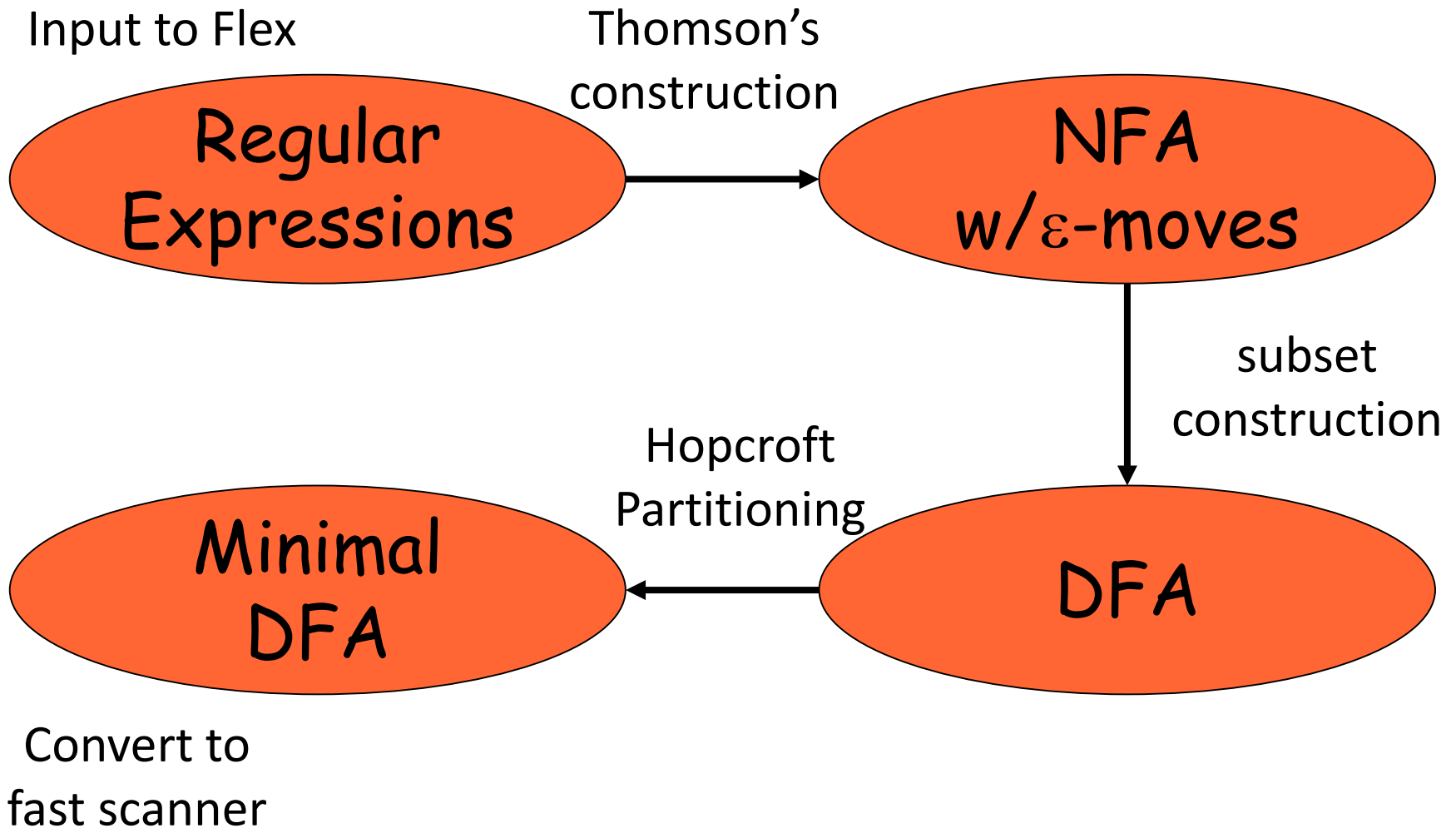


Today – part 1

- Lexing
- Flex & other scanner generators
- **Regular Expressions**
- Finite Automata
- RE \rightarrow NFA
- NFA \rightarrow DFA
- DFA \rightarrow Minimized DFA
- Limits of Regular Languages

Under The Covers

- How to go from REs to a working scanner?



Regular Languages

- Finite Alphabet, Σ , of symbols.
- word (or string), a finite sequence of symbols from Σ .
- Language over Σ is a set of words from Σ .
- Regular Expressions describe Regular Languages.
 - easy to write down, but hard to use directly
- The languages accepted by Finite Automata are also Regular.

Regular Expressions defined

- Base Cases:

- A single character a

- The empty string ε

- Recursive Rules:

If R_1 and R_2 are regular expressions

- Concatenation R_1R_2

- Union $R_1 \mid R_2$

- Closure R_1^*

- Grouping (R_1)

- REs describe Regular Languages.

RE Examples

- even a's
- odd b's
- even a's or odd b's
- even a's followed by odd b's

RE Examples

- even a's

$$b^* (a b^* a b^*)^*$$

- odd b's

$$a^* b a^* (b a^* b a^*)^*$$

- even a's or odd b's
- even a's followed by odd b's

RE Examples

- even a's

$$R^A = b^* (a b^* a b^*)^*$$

- odd b's

$$R^B = a^* b a^* (b a^* b a^*)^*$$

- even a's or odd b's

$$R^A \mid R^B$$

- even a's followed by odd b's

$$R^A R^B$$

Today – part 1

- Lexing
- Flex & other scanner generators
- Regular Expressions
- **Finite Automata**
- **RE \rightarrow NFA**
- NFA \rightarrow DFA
- DFA \rightarrow Minimized DFA
- Limits of Regular Languages

Finite Automata

- finite set of states
- set of edges from states to states labeled by letter from Σ
- initial state
- set of accepting states
- How it works:
 - Start in initial state, on each character transition goto state using edge labeled for that character.
 - If at end of word we are in accepting state, the word is in language
 - Language accepted are strings that cause FA to end in an accepting state

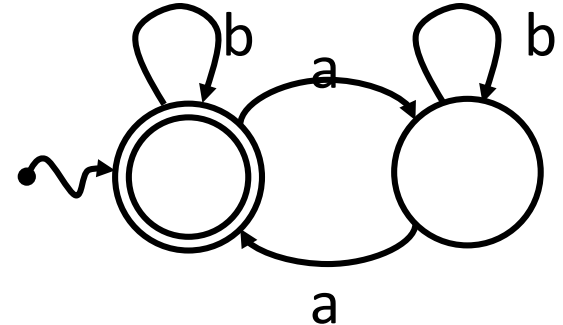
Example REs \rightarrow FA

- even a's $b^* (a b^* a b^*)^*$

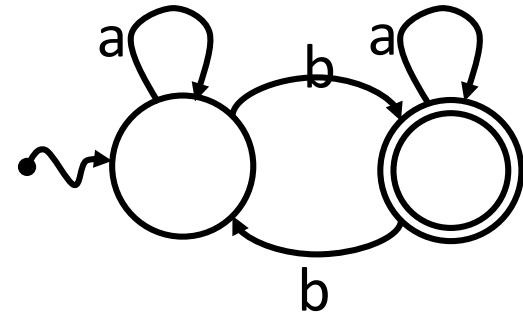
- odd b's $a^* b a^* (b a^* b a^*)^*$

Example REs \rightarrow FA

- even a's $b^* (a b^* a b^*)^*$



- odd b's $a^* b a^* (b a^* b a^*)^*$

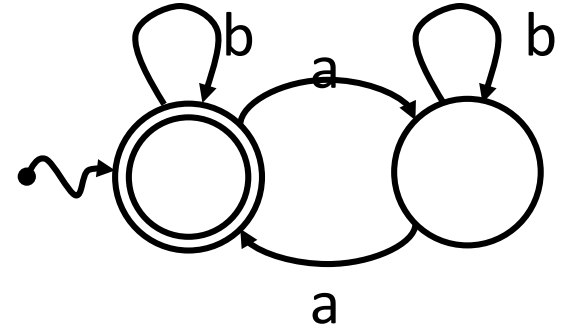


Deterministic Finite Automata
DFA

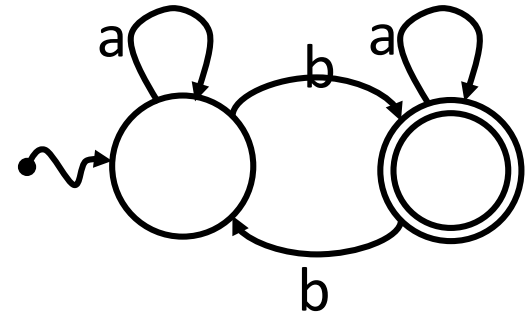
Ad Hoc

Example REs \rightarrow FA

- even a's $b^* (a b^* a b^*)^*$



- odd b's $a^* b a^* (b a^* b a^*)^*$



- even a's or odd b's

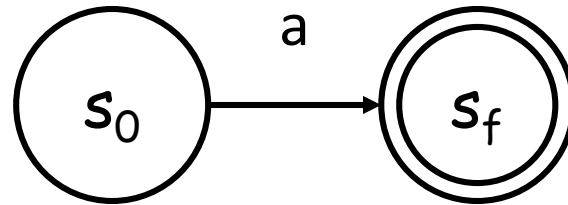
$$R^A \mid R^B$$

- even a's followed by odd b's

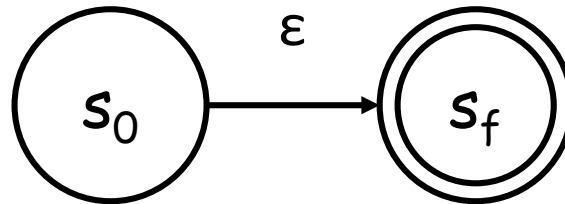
$$R^A R^B$$

Converting RE to NFA: Base Case

- for $a \in \Sigma$ the NFA $M_a = \{\Sigma, \{s_0, s_f\}, \delta, s_0, \{s_f\}\}$

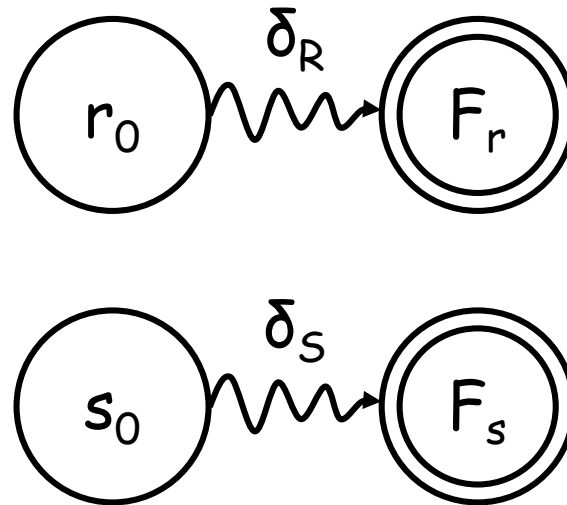


- for ε the NFA $M_\varepsilon = \{\Sigma, \{s_0, s_f\}, \delta, s_0, \{s_f\}\}$



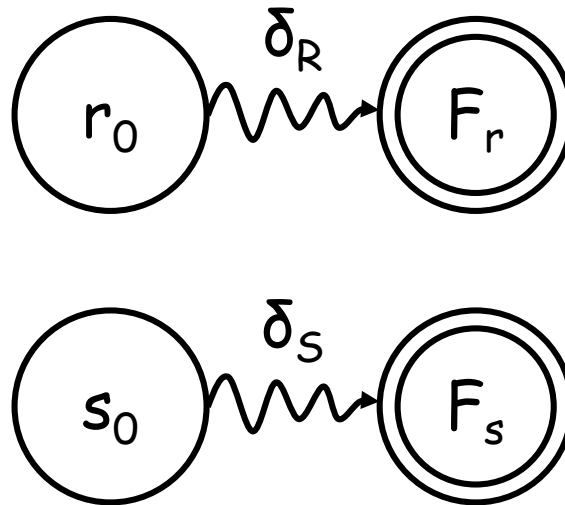
Recursive Case

- for RE R with $M_R = \{\Sigma, s_R, \delta_R, r_0, F_r\}$ and
RE S with $M_S = \{\Sigma, s_S, \delta_S, s_0, F_s\}$



R|S

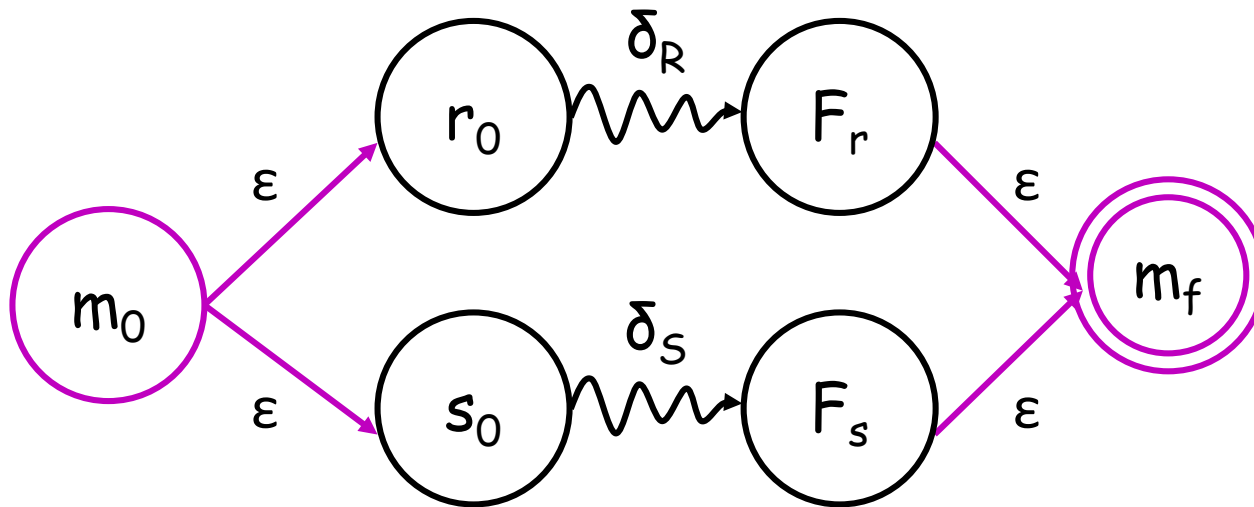
- for RE R with $M_R = \{\Sigma, s_R, \delta_R, r_0, F_r\}$ and RE S with $M_S = \{\Sigma, s_S, \delta_S, s_0, F_s\}$



- $M_{R|S} = \{\Sigma, s_R \cup s_S \cup \{m_0, m_f\}, \delta_{R|S}, m_0, m_f\}$

R|S

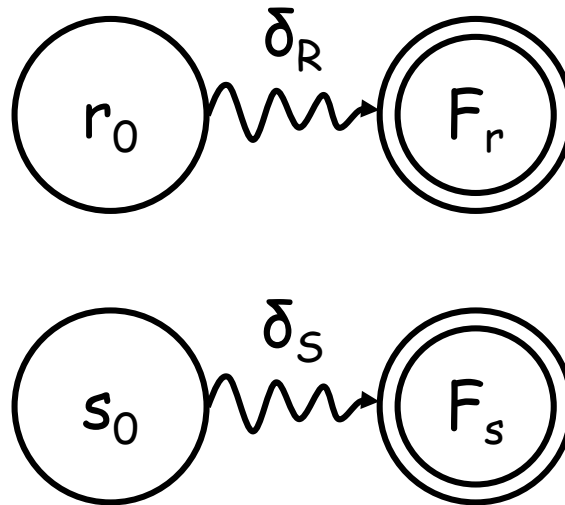
- for RE R with $M_R = \{\Sigma, s_R, \delta_R, r_0, F_r\}$ and RE S with $M_S = \{\Sigma, s_S, \delta_S, s_0, F_s\}$



- $M_{R|S} = \{\Sigma, s_R \cup s_S \cup \{m_0, m_f\}, \delta_{R|S}, m_0, m_f\}$

RS

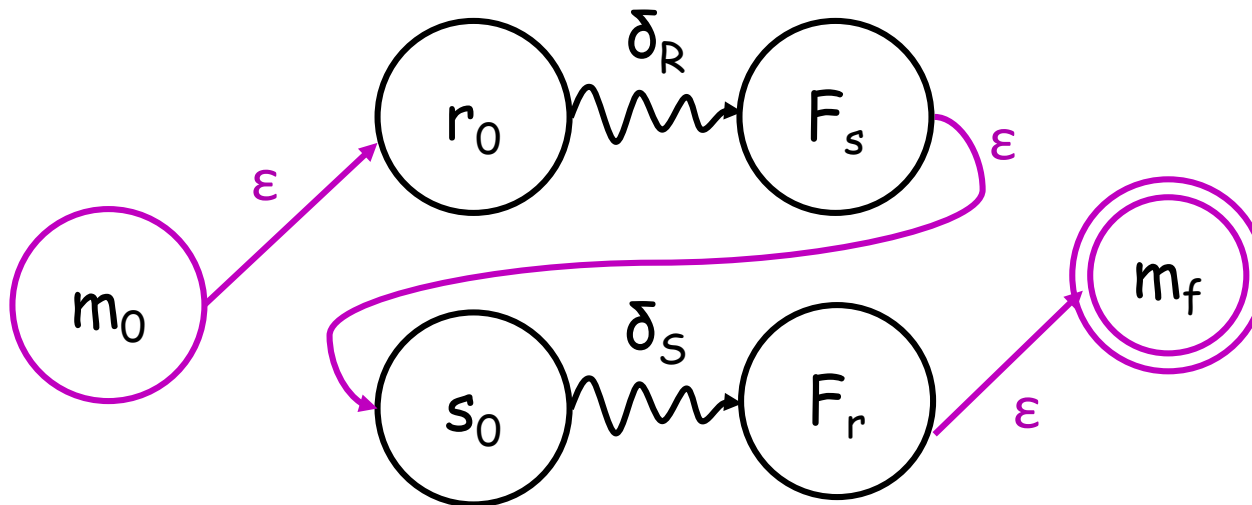
- for RE R with $M_R = \{\Sigma, s_R, \delta_R, r_0, F_r\}$ and RE S with $M_S = \{\Sigma, s_S, \delta_S, s_0, F_s\}$



- $M_{RS} = \{\Sigma, s_R \cup s_S \cup \{m_0, m_f\}, \delta_{RS}, m_0, m_f\}$

RS

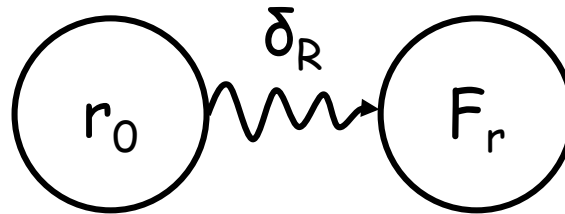
- for RE R with $M_R = \{\Sigma, s_R, \delta_R, r_0, F_r\}$ and RE S with $M_S = \{\Sigma, s_S, \delta_S, s_0, F_s\}$



- $M_{RS} = \{\Sigma, s_R \cup s_S \cup \{m_0, m_f\}, \delta_{RS}, m_0, m_f\}$

R*

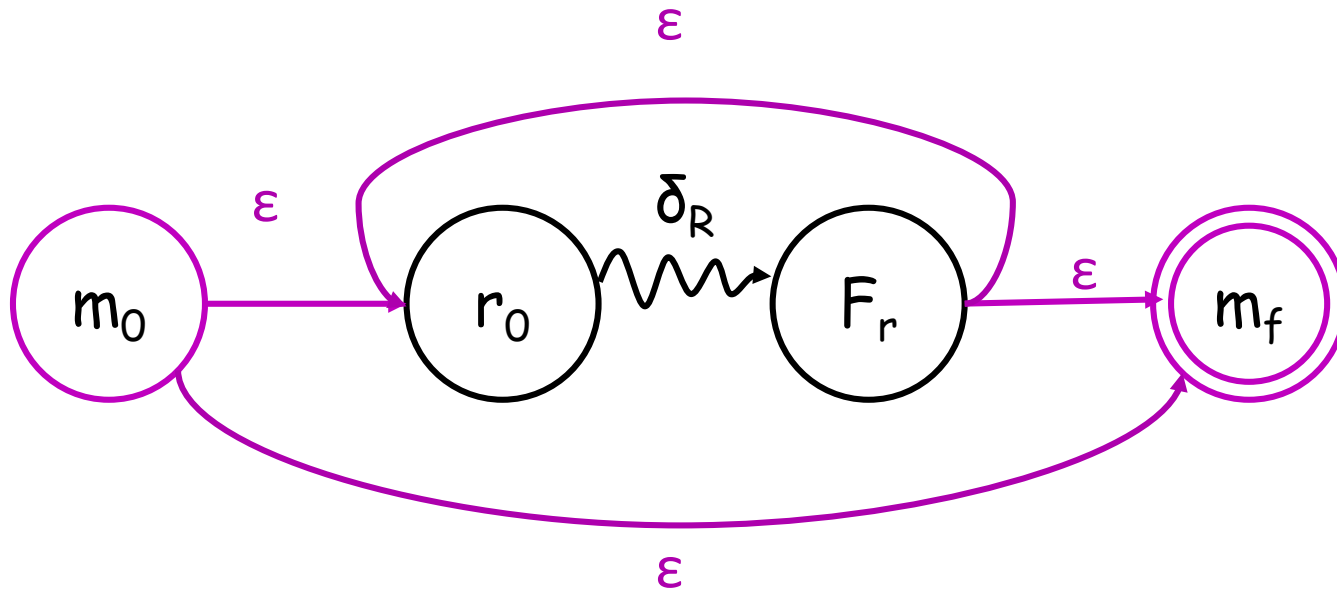
- for RE R with $M_R = \{\Sigma, s_R, \delta_R, r_0, F_r\}$



- $M_{R^*} = \{\Sigma, s_R \cup \{m_0, m_f\}, \delta_{R^*}, m_0, m_f\}$

R*

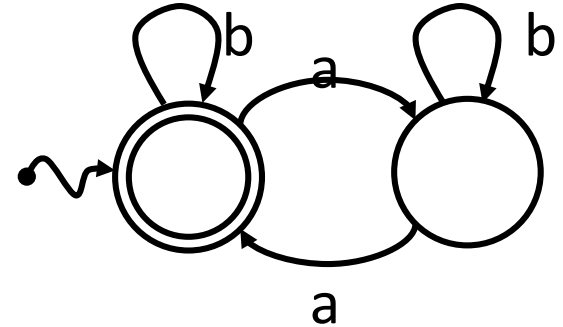
- for RE R with $M_R = \{\Sigma, s_R, \delta_R, r_0, F_r\}$



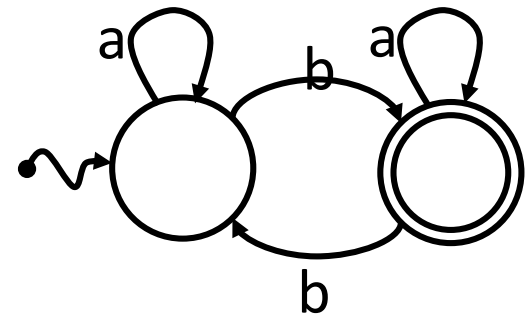
- $M_{R^*} = \{\Sigma, s_R \cup \{m_0, m_f\}, \delta_{R^*}, m_0, m_f\}$

Example REs \rightarrow FA

- even a's $b^* (a b^* a b^*)^*$



- odd b's $a^* b a^* (b a^* b a^*)^*$



- even a's or odd b's

$$R^A \mid R^B$$

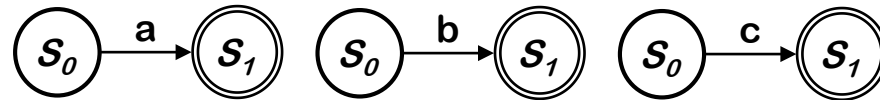
- even a's followed by odd b's

$$R^A R^B$$

Example of Thompson's Construction

Let's try $a (b | c)^*$

1. $a, b, \& c$



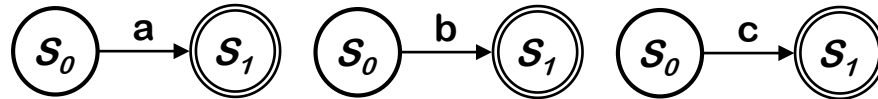
2. $b | c$

3. $(b | c)^*$

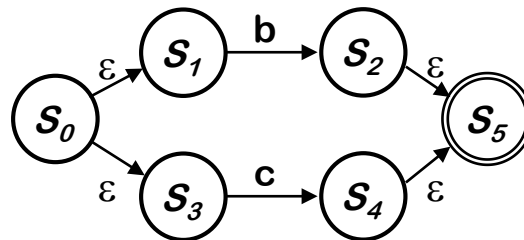
Example of Thompson's Construction

Let's try $a (b | c)^*$

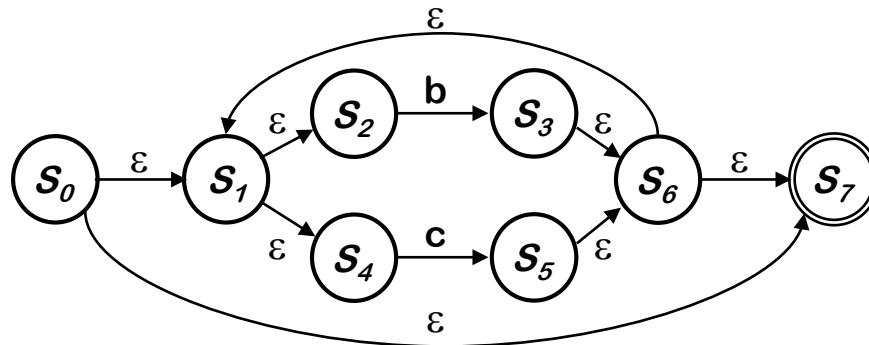
1. a, b, & c



2. b | c

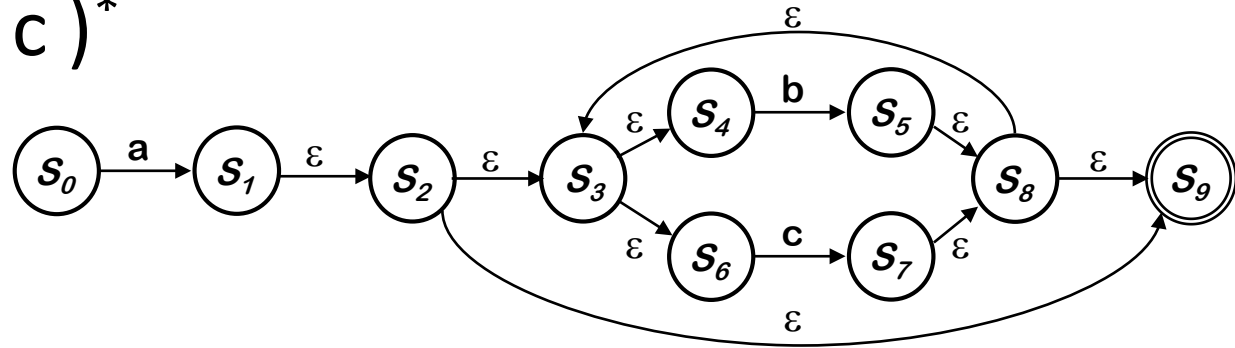


3. $(b | c)^*$



Example of Thompson's Construction

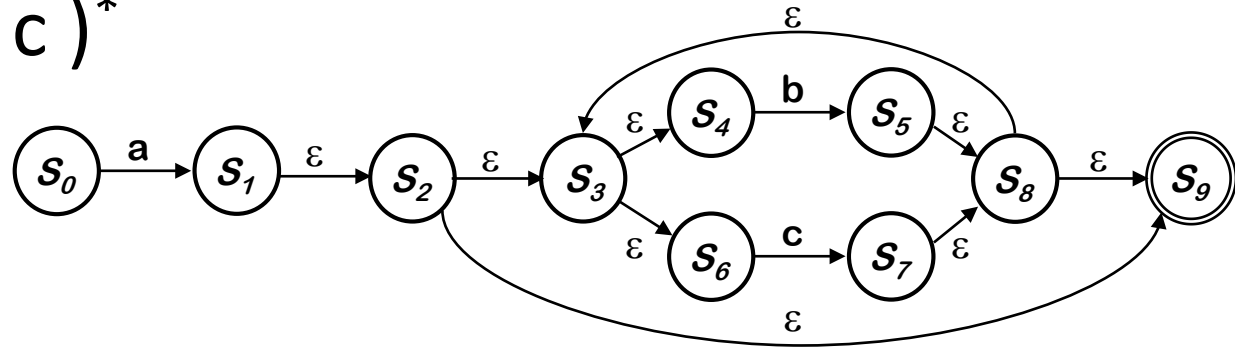
4. $a(b|c)^*$



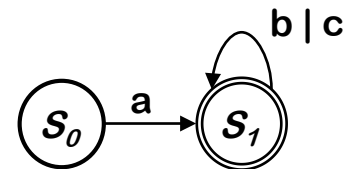
We could do a bit better. 😊

Example of Thompson's Construction


4. $a(b|c)^*$



We could do a bit better. 😊

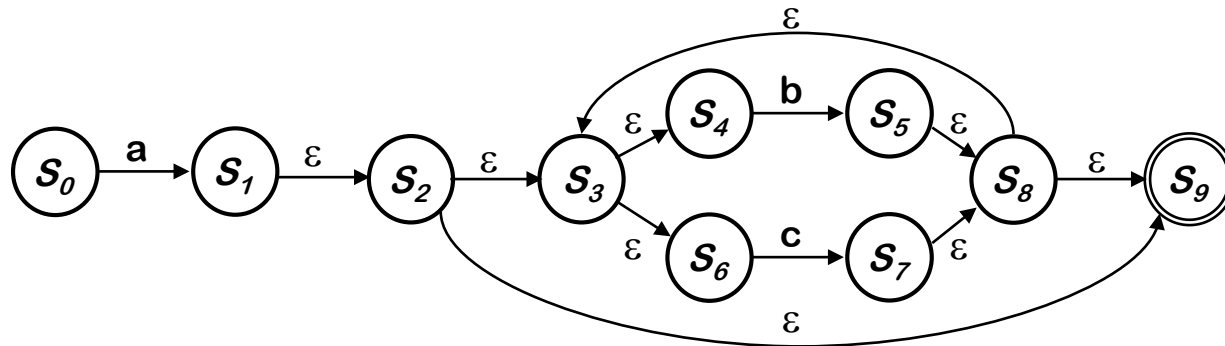


Today – part 1

- Lexing
 - Flex & other scanner generators
 - Regular Expressions
 - Finite Automata
 - RE \rightarrow NFA
 - **NFA \rightarrow DFA**
 - DFA \rightarrow Minimized DFA
 - Limits of Regular Languages
- 

RE \rightarrow NFA \rightarrow DFA

- Can't directly execute Non-deterministic FA
- Need to convert NFA to DFA
- Essentially, we will build a DFA that **simulates** the NFA

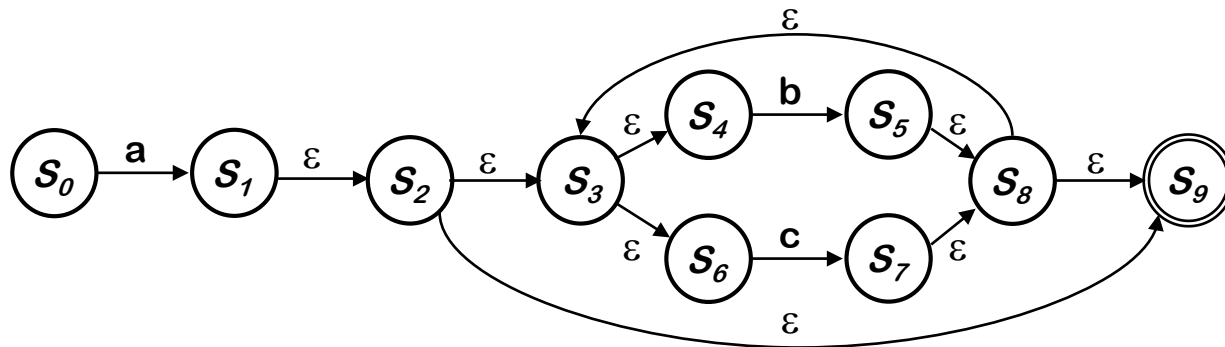


- Key idea: Keep track of all possible NFA states we could be in at each step:
the **set of all possible NFA states**
becomes the DFA state



Subset construction

- start in state $\{s_0\}$.
- For each edge create a set of all states that can be reached. Continue until done.
- All sets that contain an NFA accepting state are accepting.



Lets first deal with ϵ edges

- ϵ -closure: all states that can be reached only along ϵ -edges:
- Computing ϵ -closure(s) for $s \in S$:
 - initialize all ϵ -closure(s) = { s }
 - while some ϵ -closure(s) changed
 - foreach $s \in S$:
 - foreach $q \in \epsilon$ -closure(s) :
 - ϵ -closure(s) = ϵ -closure(s) \cup $\delta(q, \epsilon)$
- Terminates?

Subset Construction

- NFA: $\{\Sigma, Q, \delta, q_0, F\} \rightarrow$ DFA: $\{\Sigma, S, \Delta, s_0, F'\}$

$s_0 \leftarrow \varepsilon\text{-closure}(q_0)$

while \exists unmarked $s \in S$:

mark s

foreach $a \in \Sigma$

$t \leftarrow \varepsilon\text{-closure}(\text{Move}(s, a))$

if $t \notin S$:

add t to S

$\Delta(s, a) \leftarrow t$

Subset Construction

- NFA: $\{\Sigma, Q, \delta, q_0, F\} \rightarrow$ DFA: $\{\Sigma, S, \Delta, s_0, F'\}$

$s_0 \leftarrow \varepsilon\text{-closure}(q_0)$

while \exists unmarked $s \in S$:

mark s

foreach $a \in \Sigma$

$t \leftarrow \varepsilon\text{-closure}(\text{Move}(s, a))$

if $t \notin S$:

add t to S

$\Delta(s, a) \leftarrow t$

Move(s, a)

the set of states

reachable from s by a

Subset Construction

- NFA: $\{\Sigma, Q, \delta, q_0, F\} \rightarrow$ DFA: $\{\Sigma, S, \Delta, s_0, F'\}$

$s_0 \leftarrow \varepsilon\text{-closure}(q_0)$

while \exists unmarked $s \in S$:

mark s

foreach $a \in \Sigma$

$t \leftarrow \varepsilon\text{-closure}(\text{Move}(s, a))$

if $t \notin S$:

add t to S

$\Delta(s, a) \leftarrow t$

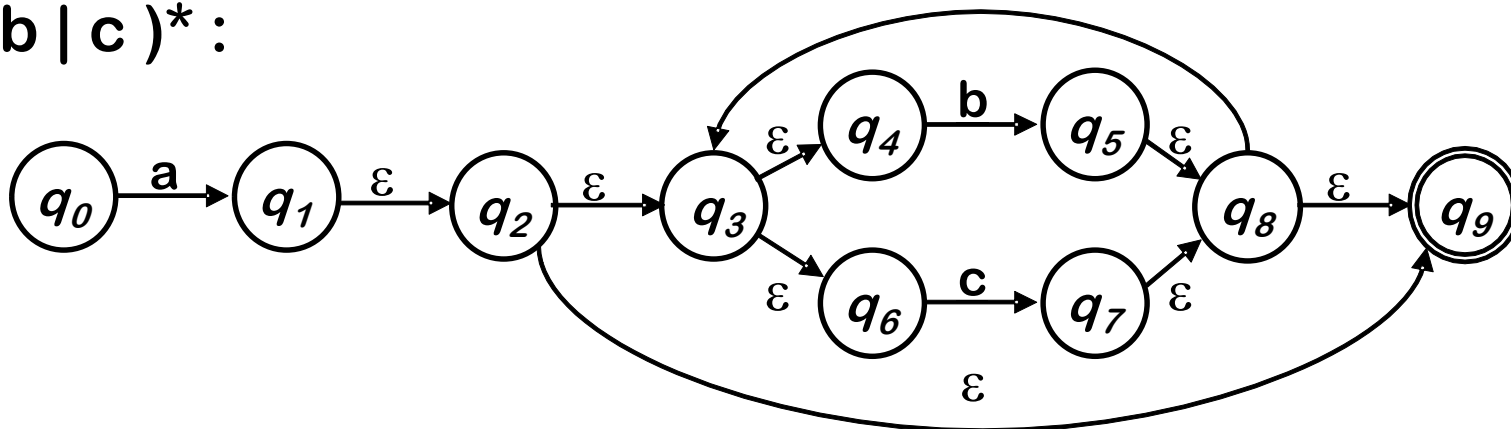
Why does this terminate?

Subset Construction

- NFA: $\{\Sigma, Q, \delta, q_0, F\} \rightarrow$ DFA: $\{\Sigma, S, \Delta, s_0, F'\}$
- Example of a fixed point computation
 - S is finite, at most ?
 - Always add to S , i.e., while loop is monotone
 - no duplicates in S
 - stop when S stops changing
- Other fixed point computations:
 - Constructing LR(1) items
 - Many Dataflow analysis (e.g., liveness)

example of subset construction

$a(b|c)^*$:

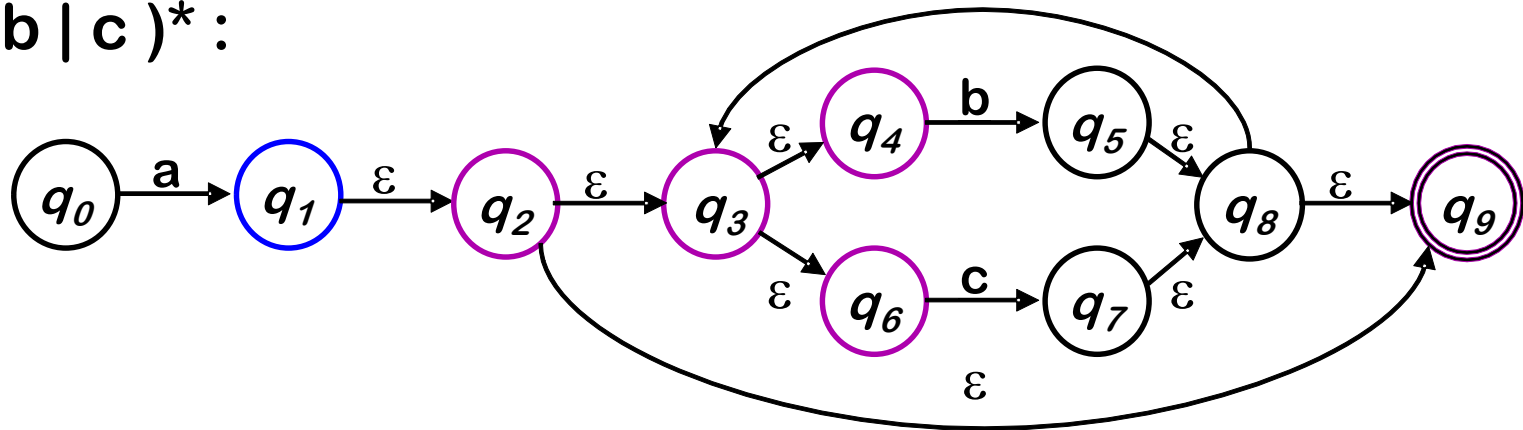


DFA States	NFA States	a	b	c
s_0	0			

Move(s_0, a)?

example of subset construction

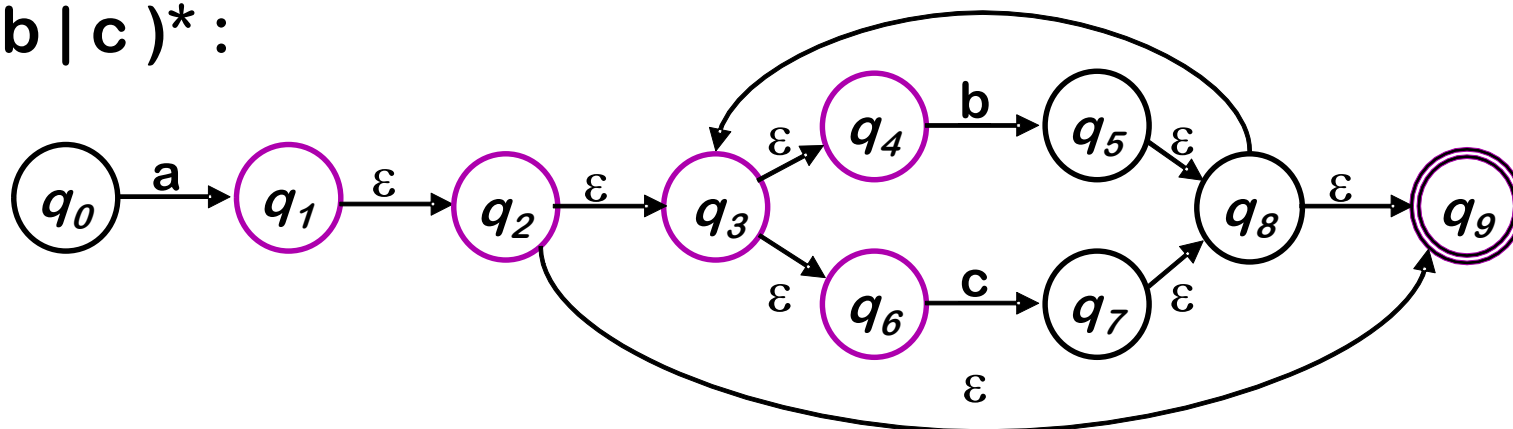
$a(b|c)^*$:



DFA States	NFA States	a	b	c
s_0	0	1, 2, 3, 4, 6, 9	-	-
s_1				

example of subset construction

$a(b|c)^*$:

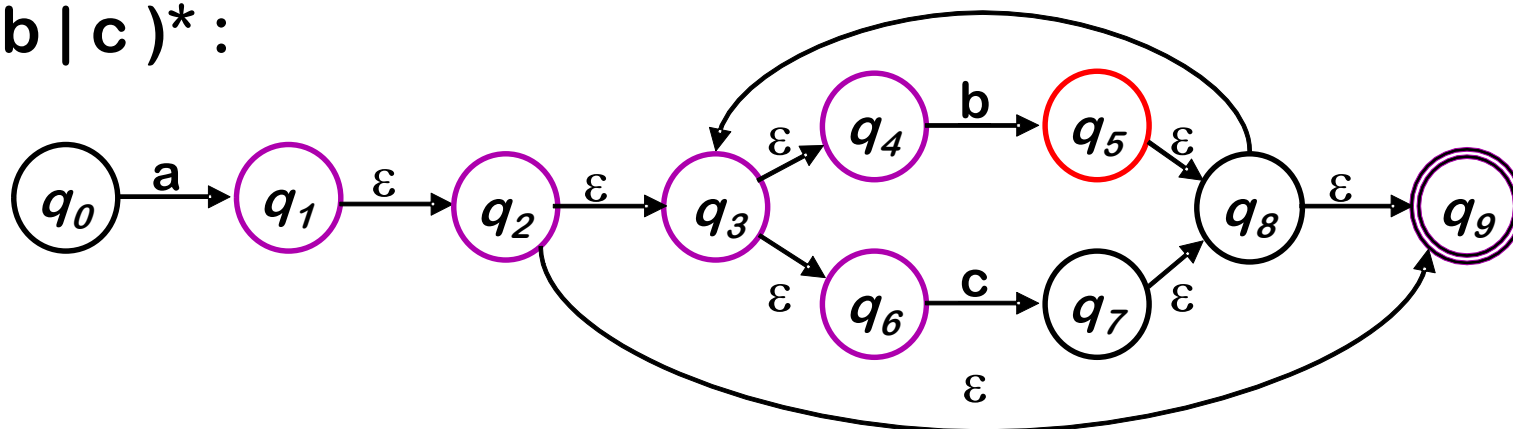


DFA States	NFA States	a	b	c
s_0	0	1, 2, 3, 4, 6, 9	-	-
s_1	1, 2, 3, 4, 6, 9	-		

b?

example of subset construction

$a(b|c)^*$:

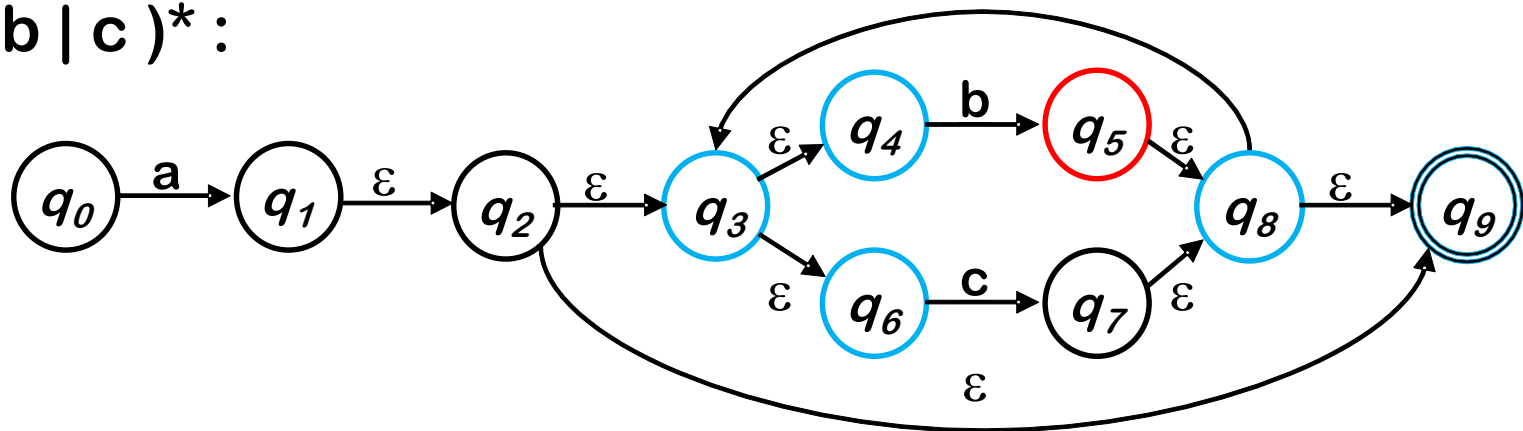


DFA States	NFA States	a	b	c
s_0	0	1, 2, 3, 4, 6, 9	-	-
s_1	1, 2, 3, 4, 6, 9	-	5	

ϵ -closure?

example of subset construction

$a(b|c)^*$:

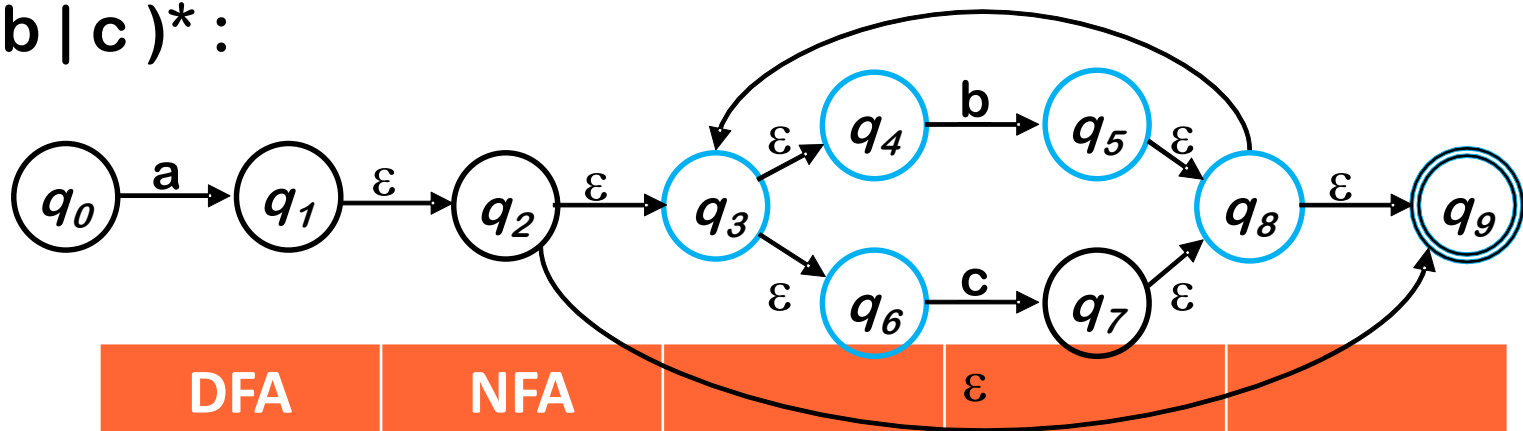


DFA States	NFA States	a	b	c
s_0	0	1, 2, 3, 4, 6, 9	-	-
s_1	1, 2, 3, 4, 6, 9	-	5, 3, 4, 6, 8, 9	

c?

example of subset construction

$a(b|c)^*$:

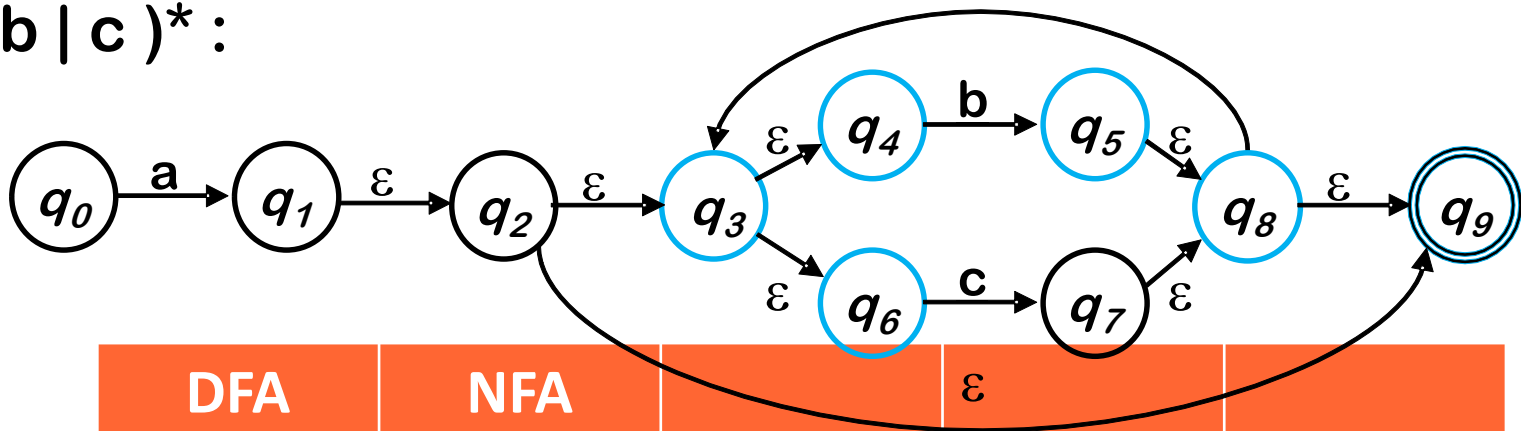


DFA States	NFA States	a	ε	b	c
s_0	0	1, 2, 3, 4, 6, 9		-	-
s_1	1, 2, 3, 4, 6, 9	-		5, 3, 4, 6, 8, 9	7, 3, 4, 6, 8, 9
s_2	5, 3, 4, 6, 8, 9				
s_3	7, 3, 4, 6, 8, 9				

$s_2, a?$

example of subset construction

$a(b|c)^*$:

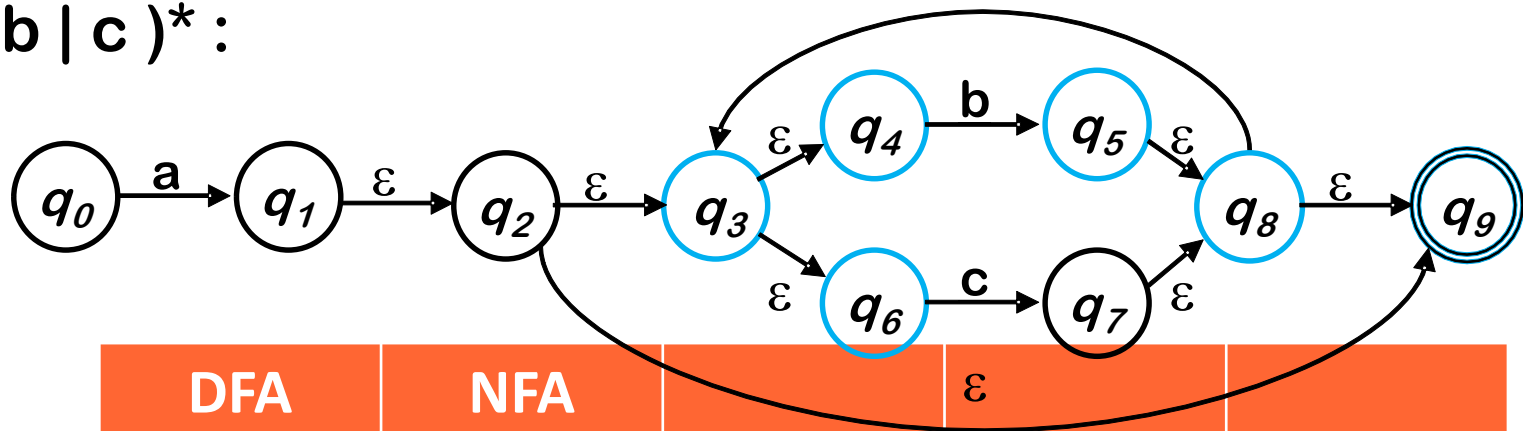


DFA States	NFA States	a	b	c
s_0	0	1, 2, 3, 4, 6, 9	-	-
s_1	1, 2, 3, 4, 6, 9	-	5, 3, 4, 6, 8, 9	7, 3, 4, 6, 8, 9
s_2	5, 3, 4, 6, 8, 9	-		
s_3	7, 3, 4, 6, 8, 9			

$s_2, b?$

example of subset construction

$a(b|c)^*$:

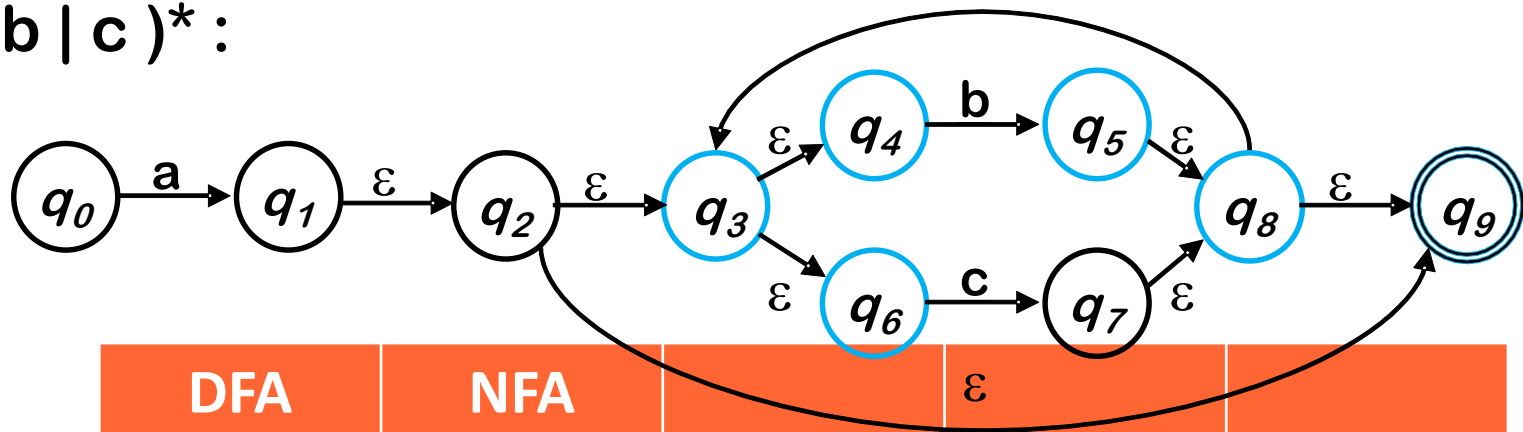


DFA States	NFA States	a	ε	b	c
s_0	0	1, 2, 3, 4, 6, 9		-	-
s_1	1, 2, 3, 4, 6, 9	-		5, 3, 4, 6, 8, 9	7, 3, 4, 6, 8, 9
s_2	5, 3, 4, 6, 8, 9	-		s_2	
s_3	7, 3, 4, 6, 8, 9				

$s_2, c?$

example of subset construction

$a(b|c)^*$:

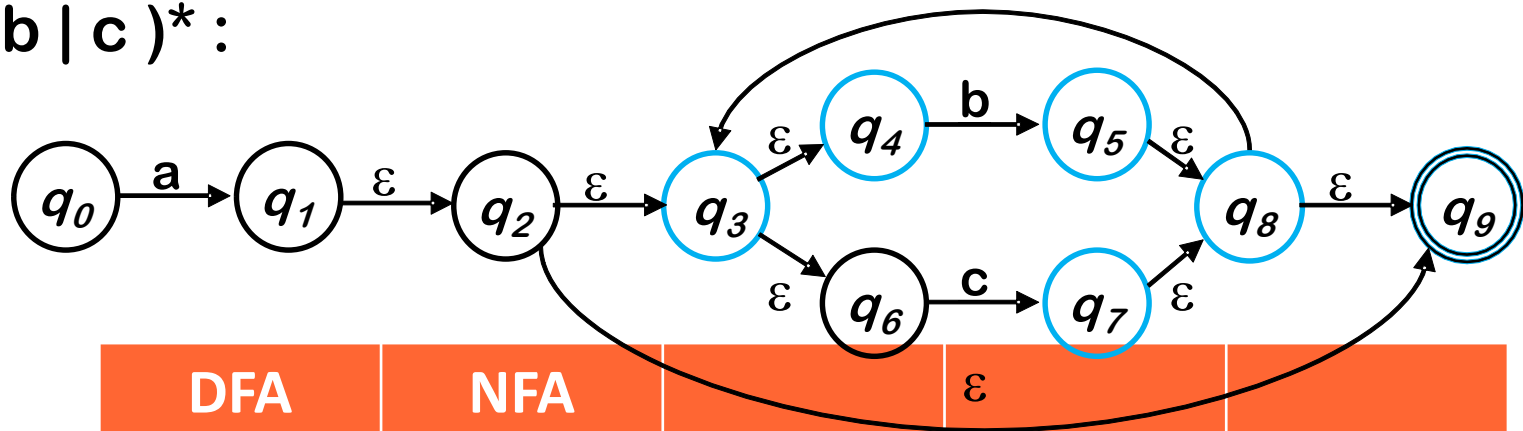


DFA States	NFA States	a	b	c
s_0	0	1, 2, 3, 4, 6, 9	-	-
s_1	1, 2, 3, 4, 6, 9	-	5, 3, 4, 6, 8, 9	7, 3, 4, 6, 8, 9
s_2	5, 3, 4, 6, 8, 9	-	s_2	s_3
s_3	7, 3, 4, 6, 8, 9			

rest?

example of subset construction

$a(b|c)^*$:

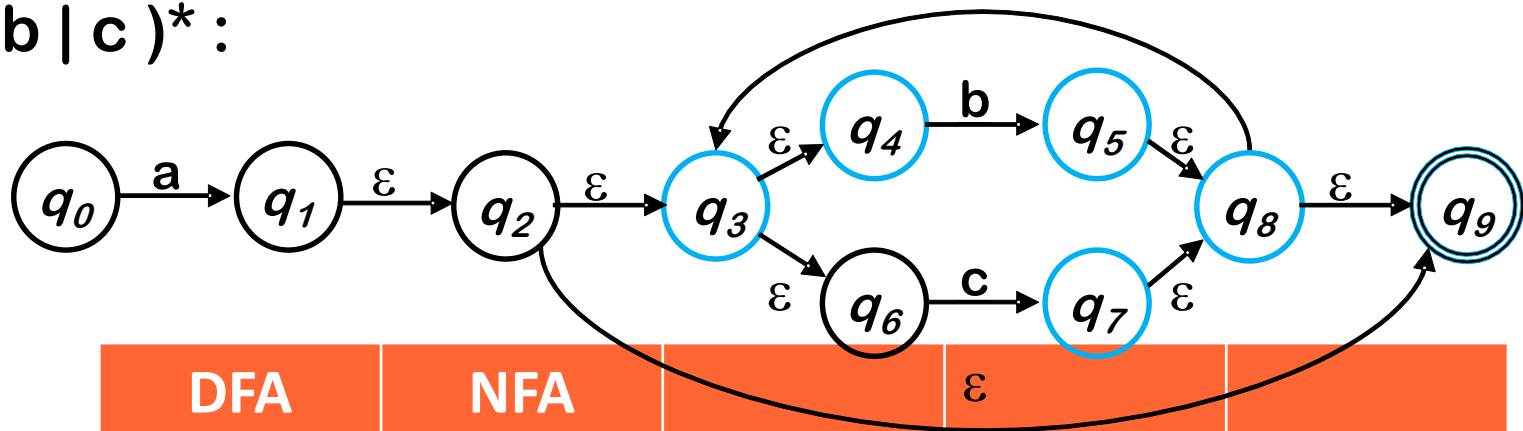


DFA States	NFA States	a	ε	b	c
s_0	0	1, 2, 3, 4, 6, 9		-	-
s_1	1, 2, 3, 4, 6, 9	-		5, 3, 4, 6, 8, 9	7, 3, 4, 6, 8, 9
s_2	5, 3, 4, 6, 8, 9	-		s_2	s_3
s_3	7, 3, 4, 6, 8, 9	-		s_2	s_3

Final?

example of subset construction

$a(b|c)^*$:

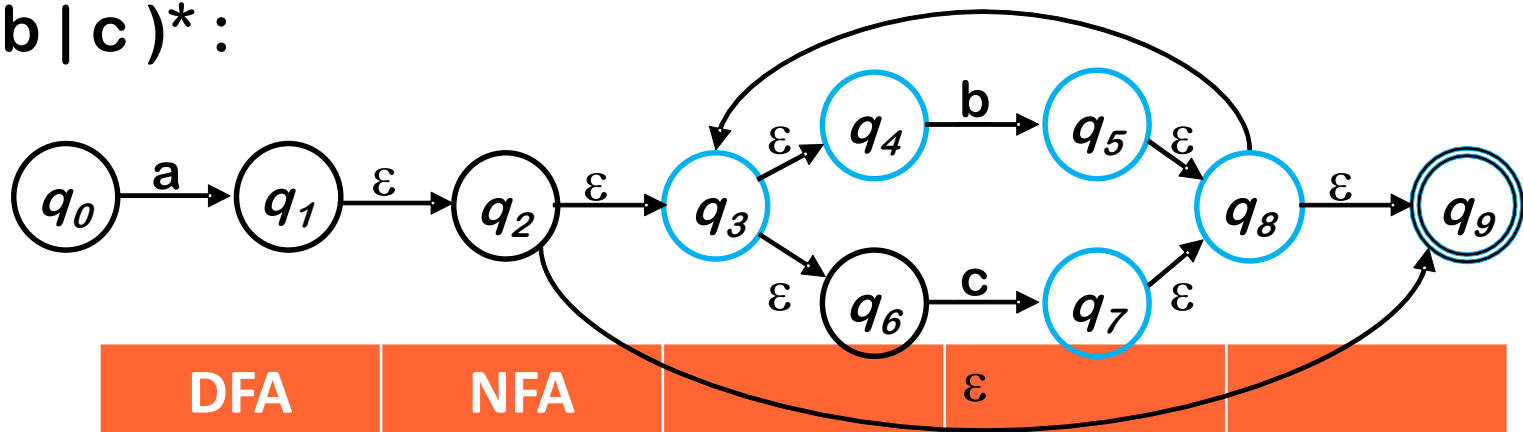


DFA States	NFA States	a	b	c
s_0	0	1, 2, 3, 4, 6, 9	-	-
s_1	1, 2, 3, 4, 6, 9	-	5, 3, 4, 6, 8, 9	7, 3, 4, 6, 8, 9
s_2	5, 3, 4, 6, 8, 9	-	s_2	s_3
s_3	7, 3, 4, 6, 8, 9	-	s_2	s_3

Final?

example of subset construction

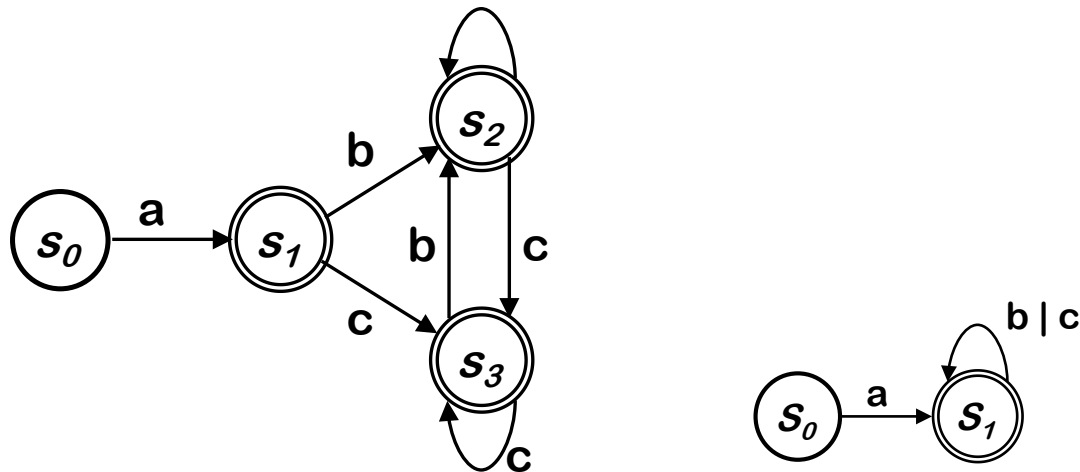
$a(b|c)^*$:



DFA States	NFA States	a	b	c
s_0	0	1, 2, 3, 4, 6, 9	-	-
s_1	1, 2, 3, 4, 6, 9	-	5, 3, 4, 6, 8, 9	7, 3, 4, 6, 8, 9
s_2	5, 3, 4, 6, 8, 9	-	s_2	s_3
s_3	7, 3, 4, 6, 8, 9	-	s_2	s_3

example of subset construction

$a(b|c)^*$:



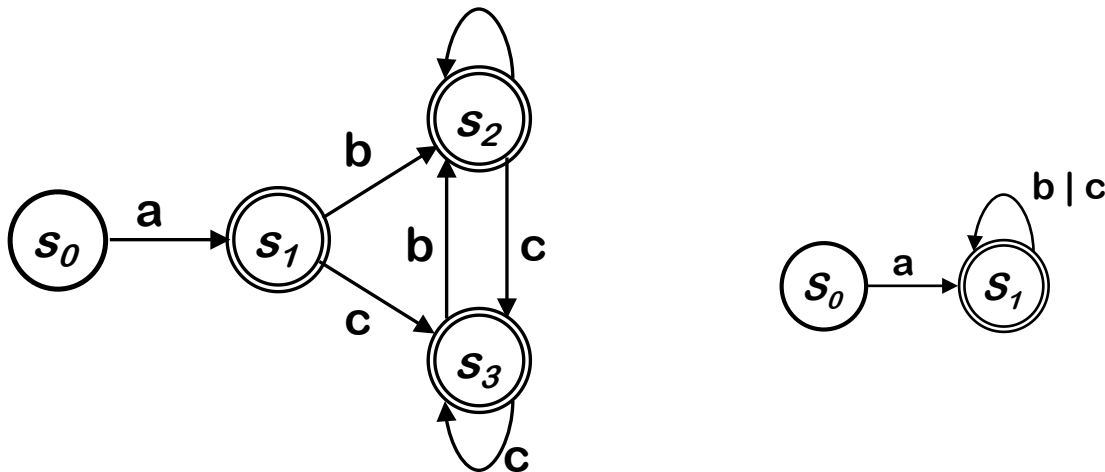
DFA States	NFA States	a	b	c
s_0	0	1, 2, 3, 4, 6, 9	-	-
s_1	1, 2, 3, 4, 6, 9	-	5, 3, 4, 6, 8, 9	7, 3, 4, 6, 8, 9
s_2	5, 3, 4, 6, 8, 9	-	s_2	s_3
s_3	7, 3, 4, 6, 8, 9	-	s_2	s_3

Today – part 1

- Lexing
- Flex & other scanner generators
- Regular Expressions
- Finite Automata
- RE \rightarrow NFA
- NFA \rightarrow DFA
- **DFA \rightarrow Minimized DFA**
- Limits of Regular Languages

DFA Minimization

- Partition states into equivalent sets
- Two states are equivalent iff:
 - paths entering them are the same
 - $\forall a \in \Sigma$, transitions lead to equivalent states
- transition on a to different sets \Rightarrow different states.



DFA Minimization

- Plan:
 - start with maximal sets: $\{ Q \}$ and $\{ Q - F \}$
 - partition sets for each $a \in \Sigma$ until no change
 - partitions become new states of minimized DFA
- Partitioning a set on “ α ”
 - Assume $q_a, \& q_b \in s$, and $\delta(q_a, \alpha) = q_x \& \delta(q_b, \alpha) = q_y$
 - If $q_x \& q_y$ are not in the same set, then s must be split
(q_a has transition on α , q_b does not $\Rightarrow \alpha$ splits s)
- One state in the final DFA cannot have two transitions on α

DFA Minimization

$P \leftarrow \{ F, \{Q-F\} \}$

while (P is still changing)

$T \leftarrow \{ \}$

for each set $S \in P$

for each $\alpha \in \Sigma$

partition S by α into S_1, S_2, \dots, S_k

$T \leftarrow T \cup S_1 \cup S_2 \cup \dots \cup S_k$

if $T \neq P$ then

$P \leftarrow T$

DFA Minimization

$P \leftarrow \{ F, \{Q-F\} \}$

while (P is still changing)

$T \leftarrow \{ \}$

 for each set $S \in P$

 for each $\alpha \in \Sigma$

 partition S by α into S_1, S_2, \dots, S_k

$T \leftarrow T \cup S_1 \cup S_2 \cup \dots \cup S_k$

 if $T \neq P$ then

$P \leftarrow T$

Another Fixed Point Alg

Terminates:

- maximum of $2^{|Q|}$ sets
- Always adding to P
- Never combining sets in P

Initial partition ensures that final states remain final.

Hopcroft's worklist algorithm is efficient.

Today – part 1

- Lexing
- Flex & other scanner generators
- Regular Expressions
- Finite Automata
- RE \rightarrow NFA
- NFA \rightarrow DFA
- DFA \rightarrow Minimized DFA
- **Limits of Regular Languages**

Regular Languages

- Regular Expressions are great
 - concise notation
 - automatic scanner generation
 - lots of useful languages
- But, ...
 - Not all languages are regular
 - Context Free Languages
 - Context Sensitive Languages
 - Even simple things like balanced parenthesis, e.g., $L = \{ A^k B^k \}$ (or nested comments!)
 - RL can't count

Not all Scanning is easy

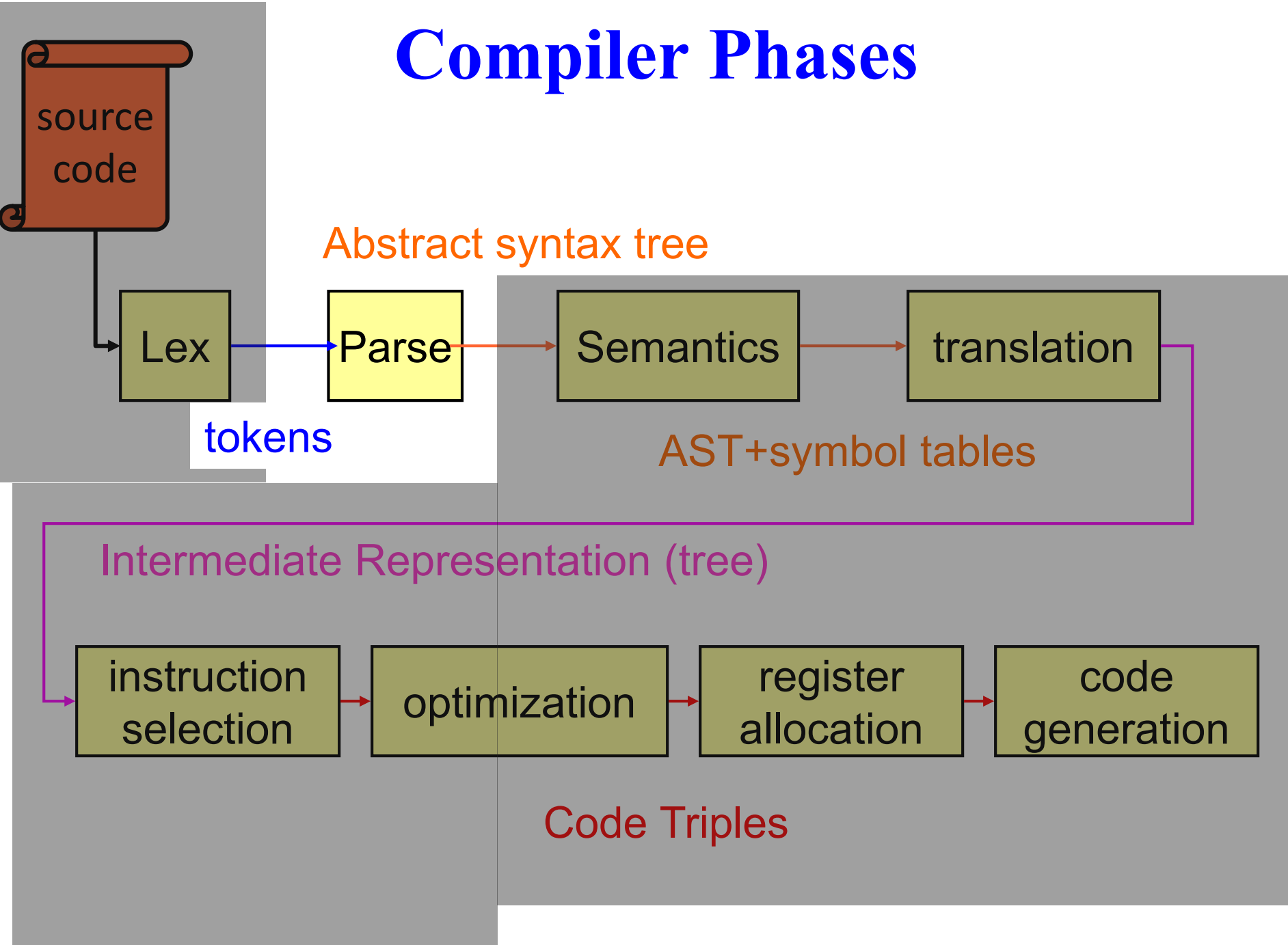
- Language design should start with lexemes
 - My favorite example from PL/I
`if then then then = else; else else = then`
- blanks not important in Fortran
- nested comments in C
- limited identifier lengths in Fortran

Today – part 2

Parsing

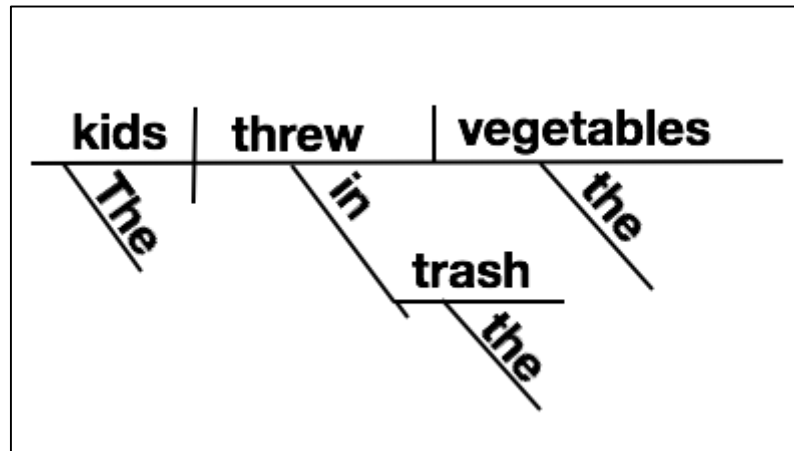
- Languages and Grammars
- Context Free Grammars
- Derivations & Parse Trees
- Ambiguity
- Top-down parsers
- FIRST, FOLLOW, and NULLABLE
- Bottom-up parsers

Compiler Phases



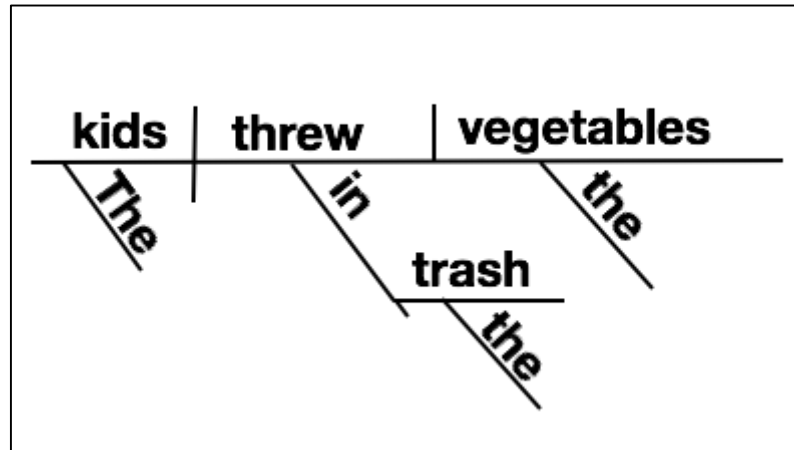
Languages

- Compiler translates from sequence of characters to an executable.
- A series of language transformations
- lexing: characters → tokens
- parsing: tokens → “sentences”



Languages

- Compiler translates from sequence of characters to an executable.
- A series of language transformations
- lexing: characters \rightarrow tokens
- parsing: tokens \rightarrow parse trees



Grammars and Languages

- A grammar, G , recognizes a language, $L(G)$
 - Σ set of terminal symbols
 - A set of non-terminals
 - S the start symbol, a non-terminal
 - P a set of productions
- Usually,
 - $\alpha, \beta, \gamma, \dots$ strings of terminals and/or non-terminals
 - A, B, C, \dots are non-terminals
 - a, b, c, \dots are terminals
- General form of a production is: $\alpha \rightarrow \beta$

Derivation

- A sequence of applying productions starting with S and ending with w

$$S \rightarrow \gamma_1 \rightarrow \gamma_2 \dots \rightarrow \gamma_{n-1} \rightarrow w$$

$$S \rightarrow^* w$$

- $L(G)$ are all the w that can be derived from S

Regular Grammar (NFA)

- Regular expressions and NFAs can be described by a regular grammar
- E.G., a^*bc^*

$$S \rightarrow aS$$

$$S \rightarrow bA$$

$$A \rightarrow \epsilon$$

$$A \rightarrow cA$$

- An example derivation of $aabc$:

$$S \rightarrow aS$$

Regular Grammar (NFA)

- Regular expressions and NFAs can be described by a regular grammar
- E.G., a^*bc^*

$$S \rightarrow aS$$

$$S \rightarrow bA$$

$$A \rightarrow \epsilon$$

$$A \rightarrow cA$$

- An example derivation of $aabc$:

$$S \rightarrow aS \rightarrow aaS$$

Regular Grammar (NFA)

- Regular expressions and NFAs can be described by a regular grammar
- E.G., a^*bc^*

$$S \rightarrow aS$$

$$S \rightarrow bA$$

$$A \rightarrow \epsilon$$

$$A \rightarrow cA$$

- An example derivation of $aabc$:

$$S \rightarrow aS \rightarrow aaS \rightarrow aabA$$

Regular Grammar (NFA)

- Regular expressions and NFAs can be described by a regular grammar
- E.G., a^*bc^*

$$S \rightarrow aS$$

$$S \rightarrow bA$$

$$A \rightarrow \epsilon$$

$$A \rightarrow cA$$

- An example derivation of $aabc$:

$$S \rightarrow aS \rightarrow aaS \rightarrow aabA \rightarrow aabcA$$

Regular Grammar (NFA)

- Regular expressions and NFAs can be described by a regular grammar
- E.G., a^*bc^*

$$S \rightarrow aS$$

$$S \rightarrow bA$$

$$A \rightarrow \epsilon$$

$$A \rightarrow cA$$

- An example derivation of $aabc$:

$$S \rightarrow aS \rightarrow aaS \rightarrow aabA \rightarrow aabcA \rightarrow aabc$$

Regular Grammar (NFA)

- Regular expressions and NFAs can be described by a regular grammar

- E.G., a^*bc^*

$$S \rightarrow aS$$

$$S \rightarrow bA$$

$$A \rightarrow \epsilon$$

$$A \rightarrow cA$$

- Above is a right-regular grammar

- All rules are of form:
 $A \rightarrow a$
 $A \rightarrow aB$
 $A \rightarrow \epsilon$

Regular Grammar (NFA)

- Regular expressions and NFAs can be described by a regular grammar
- right regular grammar:
 - $A \rightarrow a$
 - $A \rightarrow aB$
 - $A \rightarrow \epsilon$
- left regular grammar:
 - $A \rightarrow a$
 - $A \rightarrow Ba$
 - $A \rightarrow \epsilon$
- Regular grammars are either right-regular or left-regular.

Expressiveness

- Restrictions on production rules limit expressiveness of grammars.
- No restrictions allow a grammar to recognize all recursively enumerable languages
- A bit too expressive for our uses 😊
- Regular grammars cannot recognize $a^n b^n$
- We need something more expressive

Chomsky Hierarchy

Class	Language	Automaton	Form	“word” problem	Example
0	Recursively Enumerable	Turing Machine	any	undecidable	Post’s Corresp. problem
1	Context Sensitive	Linear-Bounded TM	$\alpha A \beta \rightarrow \alpha \gamma \beta$	PSPACE-complete	$a^n b^n c^n$
2	Context Free	Pushdown Automata	$A \rightarrow \alpha$	cubic	$a^n b^n$
3	Regular	NFA	$A \rightarrow a$ $A \rightarrow aB$	linear	$a^* b^*$

Today – part 2

- Languages and Grammars
- Context Free Grammars
- Derivations & Parse Trees
- Ambiguity
- Top-down parsers
- FIRST, FOLLOW, and NULLABLE
- Bottom-up parsers

Context-Free Grammar

- A context-free grammar, G , is described by:
 - Σ , a **set of terminals** (which are just the set of possible tokens from the lexer)
e.g., **if, then, while, id, int, string, ...**
 - A , a **set of non-terminals**.
Non-terminals are syntactic variables which define sets of strings in the language
e.g., **stmt, expr, term, factor, vardecl, ...**
 - S
 - P

Context-Free Grammar

- A context-free grammar, G , is described by:
 - Σ , a **set of terminals** ...
 - A , a **set of non-terminals**.
 - S , $S \in A$, the **start symbol**
The set of strings derived from S are the valid string in the language.
 - P , set of **productions** that specify how terminals and non-terminals combine to form strings in the language
a production, p , has the form: $A \rightarrow \alpha$

Context-Free Grammar

- A context-free grammar, G , is described by:
 - Σ , a **set of terminals** ...
 - A , a **set of non-terminals**.
 - S , $S \in A$, the **start symbol**
 - P , set of **productions** ...

a production, p , has the form: $A \rightarrow \alpha$

– E.g.,: $S := E$

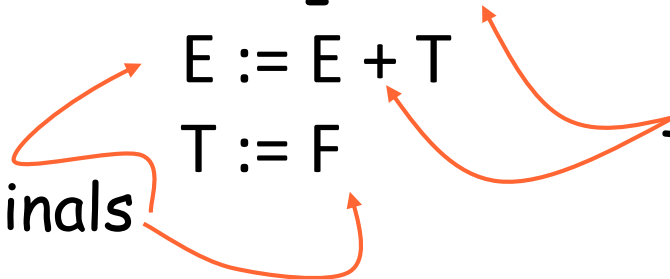
$S := \mathbf{print} E$

$E := E + T$

$T := F$

non-terminals

terminals



What makes a grammar CF?

- Only one NT on left-hand side \rightarrow context-free
- What makes a grammar context-sensitive?
- $\alpha A \beta \rightarrow \alpha \gamma \beta$ where
 - α or β may be empty,
 - but γ is not-empty
- Are context-sensitive grammars useful for compiler writers?

Simple Grammar of Expressions

S := Exp

Exp := Exp + Exp

Exp := Exp - Exp

Exp := Exp * Exp

Exp := Exp / Exp

Exp := **id**

Exp := **int**

Describes a language of expressions. e.g.: $2+3*x$

Derivations

- A sequence of steps in which a non-terminal is replaced by its right-hand side.

1 $S \Rightarrow \text{Exp}$ S

2 $\text{Exp} \Rightarrow \text{Exp} + \text{Exp}$ Exp

3 $\text{Exp} \Rightarrow \text{int}_2 + \text{Exp} * \text{id}_x$ Exp

4 $\text{Exp} := \text{Exp} * \text{Exp}$ Exp by 3 $\Rightarrow \text{Exp} + \text{id}_x$

5 $\text{Exp} := \text{Exp} / \text{Exp}$ Exp by 2 $\Rightarrow \text{Exp} + \text{Exp} * \text{id}_x$

6 $\text{Exp} := \text{id}$ Exp by 7 $\Rightarrow \text{int}_2 + \text{Exp} * \text{id}_x$

7 $\text{Exp} := \text{int}$ Exp by 7 $\Rightarrow \text{int}_2 + \text{int}_3 * \text{id}_x$

There are possibly many derivations determined by the NT chosen to expand.

Leftmost Derivations

- Leftmost derivation: leftmost NT always chosen

1	$S := \text{Exp}$	S
2	$\text{Exp} := \text{Exp} + \text{Exp}$	by 1 $\Rightarrow \text{Exp}$
3	$\text{Exp} := \text{Exp} - \text{Exp}$	by 4 $\Rightarrow \text{Exp} * \text{Exp}$
4	$\text{Exp} := \text{Exp} * \text{Exp}$	by 2 $\Rightarrow \text{Exp} + \text{Exp} * \text{Exp}$
5	$\text{Exp} := \text{Exp} / \text{Exp}$	by 7 $\Rightarrow \text{int}_2 + \text{Exp} * \text{Exp}$
6	$\text{Exp} := \text{id}$	by 7 $\Rightarrow \text{int}_2 + \text{int}_3 * \text{Exp}$
7	$\text{Exp} := \text{int}$	by 6 $\Rightarrow \text{int}_2 + \text{int}_3 * \text{id}_x$

Rightmost Derivations

- Rightmost derivation: rightmost NT always chosen

1	$S := \text{Exp}$	S
2	$\text{Exp} := \text{Exp} + \text{Exp}$	by 1 $\Rightarrow \text{Exp}$
3	$\text{Exp} := \text{Exp} - \text{Exp}$	by 4 $\Rightarrow \text{Exp} * \text{Exp}$
4	$\text{Exp} := \text{Exp} * \text{Exp}$	by 6 $\Rightarrow \text{Exp} * \text{id}_x$
5	$\text{Exp} := \text{Exp} / \text{Exp}$	by 2 $\Rightarrow \text{Exp} + \text{Exp} * \text{id}_x$
6	$\text{Exp} := \text{id}$	by 7 $\Rightarrow \text{Exp} + \text{int}_3 * \text{id}_x$
7	$\text{Exp} := \text{int}$	by 7 $\Rightarrow \text{int}_2 + \text{int}_3 * \text{id}_x$

Parse Trees

- symbols in rhs are children of NT being rewritten

S

by 1 \Rightarrow Exp

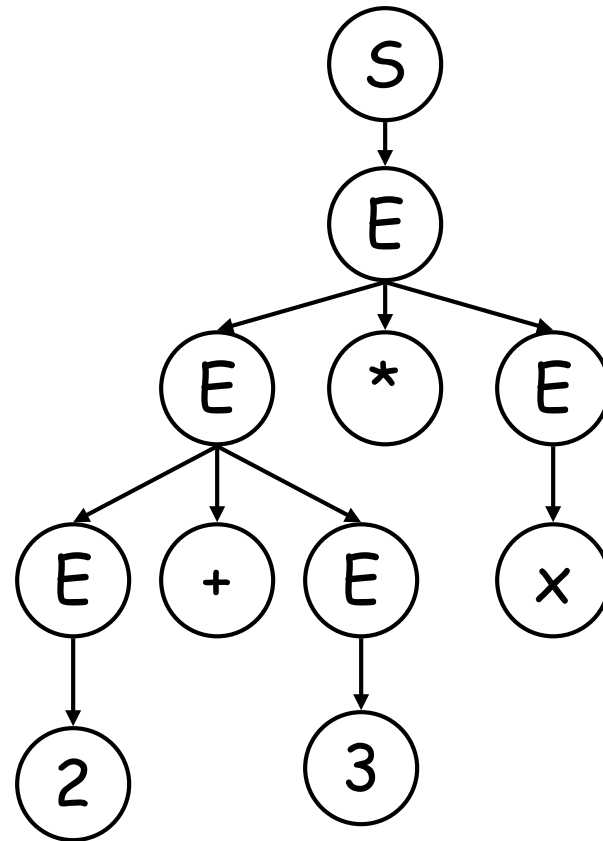
by 4 \Rightarrow $Exp * Exp$

by 2 \Rightarrow $Exp + Exp * Exp$

by 7 \Rightarrow $int_2 + Exp * Exp$

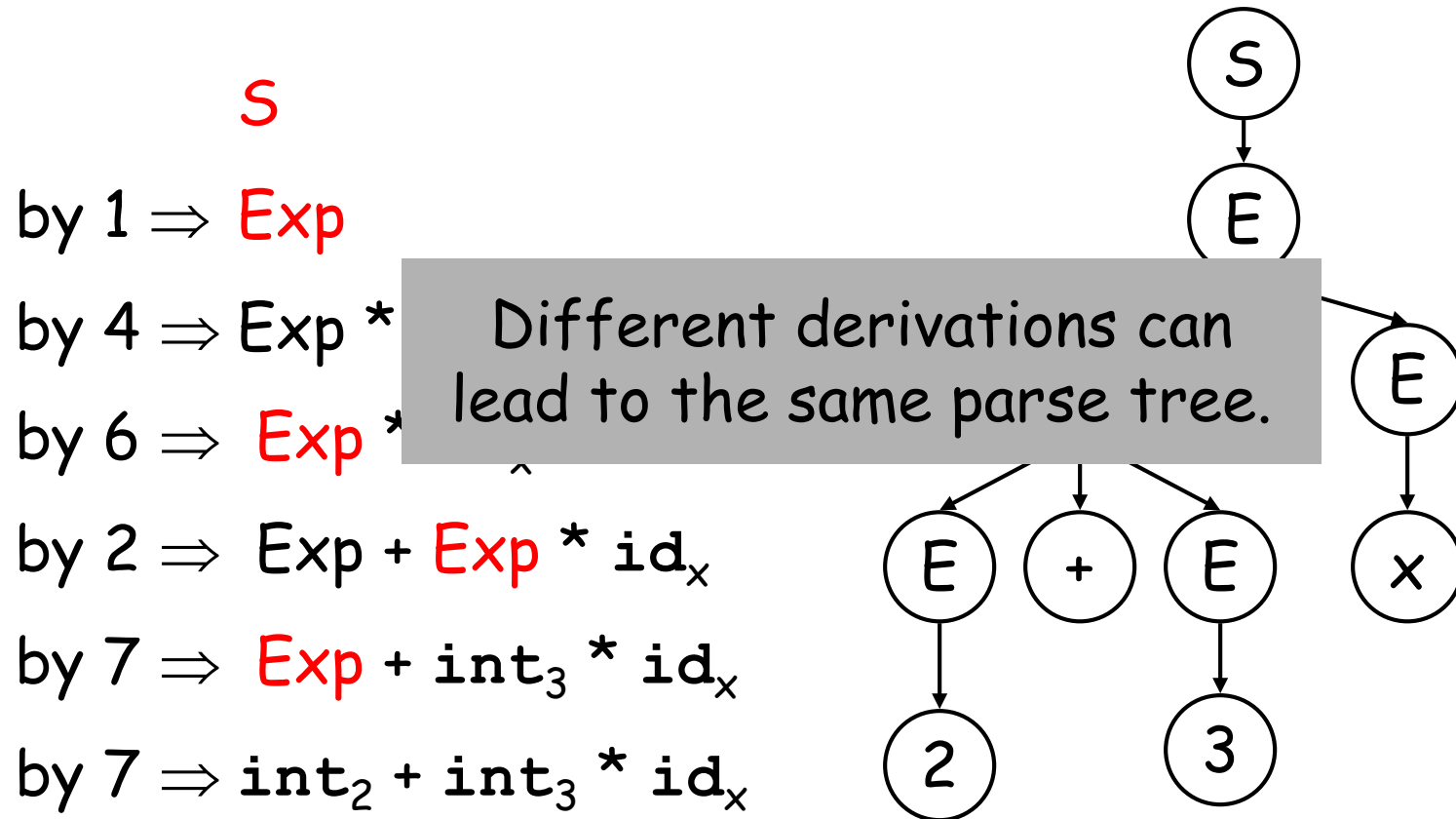
by 7 \Rightarrow $int_2 + int_3 * Exp$

by 6 \Rightarrow $int_2 + int_3 * id_x$



Parse Trees

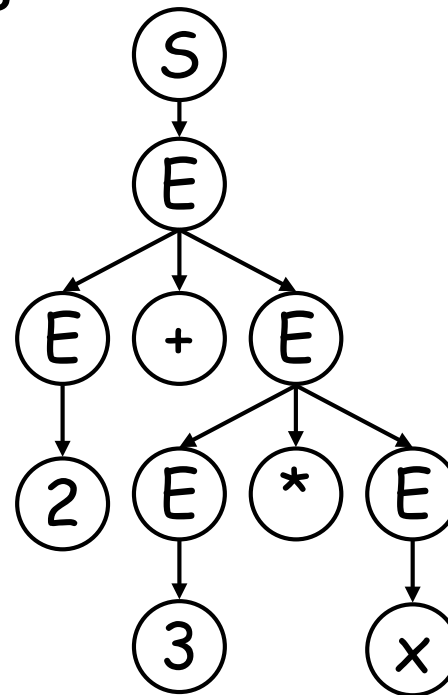
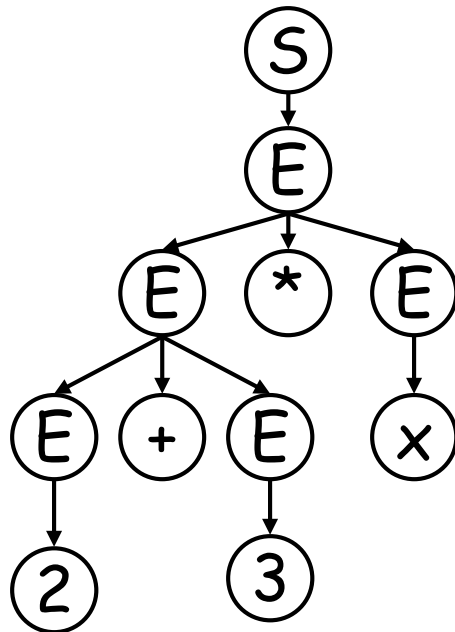
- parse tree for rightmost derivation



What about different parse trees for same sentence?

Ambiguous Grammars

- A grammar is ambiguous if it has a sentence with >1 parse trees. or,
- If a grammar has >1 leftmost (rightmost) derivations it is ambiguous



Converting Expression Grammar

- Adding precedence with more non-terminals
- One for each level of precedence:
 - (+, -) exp
 - (*, /) term
 - (**id**, **int**) factor
 - Make sure parse derives sentences that respect the precedence
 - Make sure that extra levels of precedence can be bypassed, i.e., “x” is still legal

A Better Exp Grammar

1	S	$:= \text{Exp}$	S
2	Exp	$:= \text{Exp} + \text{Term}$	by 1 $\Rightarrow \text{Exp}$
3	Exp	$:= \text{Exp} - \text{Term}$	by 2 $\Rightarrow \text{Exp} + \text{Term}$
4	Exp	$:= \text{Term}$	by 4 $\Rightarrow \text{Term} + \text{Term}$
5	Term	$:= \text{Term} * \text{Factor}$	by 7 $\Rightarrow \text{Factor} + \text{Term}$
6	Term	$:= \text{Term} / \text{Factor}$	by 9 $\Rightarrow \text{int}_2 + \text{Term}$
7	Term	$:= \text{Factor}$	by 5 $\Rightarrow \text{int}_2 + \text{Term} * \text{Factor}$
8	Factor	$:= \text{id}$	by 7 $\Rightarrow \text{int}_2 + \text{Factor} * \text{Factor}$
9	Factor	$:= \text{int}$	by 9 $\Rightarrow \text{int}_2 + \text{int}_3 * \text{Factor}$
			by 8 $\Rightarrow \text{int}_2 + \text{int}_3 * \text{id}_x$

What is the parse tree?

Another Ambiguous Grammar

S := if E then S
| if E then S else S
| other

- What is the parse tree for:
if E then if E then S else S?
- What is the language designers intention?
- Is there a context-free solution?

Dangling Else Grammar

S := matchedS
 | unmatchedS

unmatchedS := **if** E **then** S
 | if E then matchedS else unmatchedS

matchedS := **if** E **then** matchedS **else** matchedS
 | **other**

- Is this clearer?
- What is parse tree for: **if** E **then** **if** E **then** S **else** S?

Parser generators provide a better way

Parsing a CFG

- Top-Down
 - start at root of parse-tree
 - pick a production and expand to match input
 - may require backtracking
 - if no backtracking required, predictive
- Bottom-up
 - start at leaves of tree
 - recognize valid prefixes of productions
 - consume input and change state to match
 - use stack to track state

Top-down Parsers

- Starts at root of parse tree and recursively expands children that match the input
- In general case, may require backtracking
- Such a parser uses recursive descent.
- When a grammar does not require backtracking a **predictive parser** can be built.

A Predictive Parser

S := B S F
|
B := b
F := f

Idea is for parser to do something besides recognize legal sentences.

```
S() {  
    if match('b') -> B(); S(); F(); action();  
    else return;  
}  
  
B() { mustMatch('b'); action(); return;}  
F() { mustMatch('f'); action(); return;}
```

Top-Down parsing

- Start with root of tree, i.e., S
- Repeat until entire input matched:
 - pick a non-terminal, A , and pick a production $A \rightarrow \gamma$ that can match input, and expand tree
 - if no such rule applies, backtrack
- Key is obviously selecting the right production

Top-down for Exp Grammar

1	$S := E$
2	$E := E + T$
3	$E := E - T$
4	$E := T$
5	$T := T * F$
6	$T := T / F$
7	$T := F$
8	$F := id$
9	$F := int$

S
by 1 $\Rightarrow E$

| $int_2 - int_3 * id_x$

| $int_2 - int_3 * id_x$

Top-down for Exp Grammar

1	$S := E$
2	$E := E + T$
3	$E := E - T$
4	$E := T$
5	$T := T * F$
6	$T := T / F$
7	$T := F$
8	$F := id$
9	$F := int$

S
 by 1 $\Rightarrow E$

 by 2 $\Rightarrow E + T$
 by 4 $\Rightarrow T + T$
 by 7 $\Rightarrow F + T$
 by 9 $\Rightarrow int_2 + T$

$| int_2 - int_3 * id_x$
 $| int_2 - int_3 * id_x$

 $| int_2 - int_3 * id_x$
 $| int_2 - int_3 * id_x$
 $| int_2 - int_3 * id_x$
 $int_2 | - int_3 * id_x$

Must backtrack here!

Top-down for Exp Grammar

1	$S := E$
2	$E := E + T$
3	$E := E - T$
4	$E := T$
5	$T := T * F$
6	$T := T / F$
7	$T := F$
8	$F := id$
9	$F := int$

S | int₂ - int₃ * id_x

by 1 $\Rightarrow E$ | int₂ - int₃ * id_x

by 2 $\Rightarrow E + T$ | int₂ - int₃ * id_x

by 4 $\Rightarrow T + T$ | int₂ - int₃ * id_x

by 7 $\Rightarrow F + T$ | int₂ - int₃ * id_x

by 9 $\Rightarrow int_2 + T$ int₂ | - int₃ * id_x

by 3 $\Rightarrow E - T$ | int₂ - int₃ * id_x

by 4 $\Rightarrow T - T$ | int₂ - int₃ * id_x

by 7 $\Rightarrow F - T$ | int₂ - int₃ * id_x

by 9 $\Rightarrow int_2 - T$ int₂ | - int₃ * id_x

by 5 $\Rightarrow int_2 - T * F$ int₂ - | int₃ * id_x

Top-down for Exp Grammar

1	$S := E$
2	$E := E + T$
3	$E := E - T$
4	$E := T$
5	$T := T * F$
6	$T := T / F$
7	$T := F$
8	$F := id$
9	$F := int$

	S	int ₂ - int ₃ * id _x
by 1 ⇒	E	int ₂ - int ₃ * id _x
<hr/>		
by 2 ⇒	$E + T$	int ₂ - int ₃ * id _x
by 4 ⇒	$T + T$	int ₂ - int ₃ * id _x
by 7 ⇒	$F + T$	int ₂ - int ₃ * id _x
by 9 ⇒	int ₂ + T	int ₂ - int ₃ * id _x
<hr/>		
by 3 ⇒	$E - T$	int ₂ - int ₃ * id _x
by 4 ⇒	$T - T$	int ₂ - int ₃ * id _x
by 7 ⇒	$F - T$	int ₂ - int ₃ * id _x
by 9 ⇒	int ₂ - T	int ₂ - int ₃ * id _x

What kind of derivation is this parsing?

int₂ - int₃ * id_x

Top-down for Exp Grammar

1	$S := E$
2	$E := E + T$
3	$E := E - T$
4	$E := T$
5	$T := T * F$
6	$T := T / F$
7	$T := F$
8	$F := id$
9	$F := int$

S
by 1 $\Rightarrow E$
by 2 $\Rightarrow E + T$
by 2 $\Rightarrow E + E + T$
by 2 $\Rightarrow E + E + E + T$

$| int_2 - int_3 * id_x$
 $| int_2 - int_3 * id_x$
 $| int_2 - int_3 * id_x$
 $| int_2 - int_3 * id_x$
 $| int_2 - int_3 * id_x$

Will not terminate! Why?

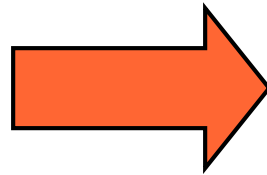
grammar is left-recursive

What should we do about it?

Eliminate left-recursion

Does this work?

```
1 S := E
2 E := E + T
3 E := E - T
4 E := T
5 T := T * F
6 T := T / F
7 T := F
8 F := id
9 F := int
```



```
1 S := E
2 E := T + E
3 E := T - E
4 E := T
5 T := F * T
6 T := F / T
7 T := F
8 F := id
9 F := int
```

It is right recursive, but also right associative!

Eliminating Left-Recursion

- Given 2 productions:

$$A := A \alpha \mid \beta$$

Where neither α nor β start with A

(e.g., For example, $E := E + T \mid T$)
 α β

- Make it right-recursive:

$A := \beta R$
$R := \alpha R$

R is right recursive

- Extends to general case.

Rewriting Exp Grammar

```
1 S := E
2 E := E + T
3 E := E - T
4 E := T
5 T := T * F
6 T := T / F
7 T := F
8 F := id
9 F := int
```

```
1 S := E
2' E' := + T E'
3' E' := - T E'
4' E' :=
5' T' := * F T'
6' T' := / F T'
7' T' :=
8 F := id
9 F := int
```

```
2 E := T E'
5 T := F T'
```

Is this legible?

Try again

```

1  S := E
2  E := T E'
2' E' := + T E'
3' E' := - T E'
4' E' :=
5  T := F T'
5' T' := * F T'
6' T' := / F T'
7' T' :=
8  F := id
9  F := int
    
```

```

          S
by 1 ⇒ E
by 2 ⇒ T E'
by 5 ⇒ F T' E'
by 9 ⇒ 2 T' E'
by 7' ⇒ 2 E'
by 3' ⇒ 2 - T E'
by 5 ⇒ 2 - F T' E'
by 9 ⇒ 2 - 3 T' E'
by 5' ⇒ 2 - 3 * F T' E'
    
```

```

●int2 - int3 * idx
●int2 - int3 * idx
●int2 - int3 * idx
●int2 - int3 * idx
int2 ● - int3 * idx
int2 ● - int3 * idx
int2 - ●int3 * idx
int2 - ●int3 * idx
int2 - int3 ● * idx
int2 - int3 * ●idx
int3 * idx ●
int3 * idx ●
int3 * idx ●
    
```

Unlike previous time we tried this, it appears that only one production applies at a time. I.e., no backtracking needed. Why?

Lookahead

- How to pick right production?
- Lookahead in input stream for guidance
- General case: arbitrary lookahead required
- Luckily, many context-free grammars can be parsed with limited lookahead
- If we have $A \rightarrow \alpha \mid \beta$, then we want to correctly choose either $A \rightarrow \alpha$ or $A \rightarrow \beta$
- define $\text{FIRST}(\alpha)$ as the set of tokens that can be first symbol of α , i.e.,
$$a \in \text{FIRST}(\alpha) \text{ iff } \alpha \rightarrow^* a\gamma \text{ for some } \gamma$$

Lookahead

- How to pick right production?
- If we have $A \rightarrow \alpha \mid \beta$, then we want to correctly choose either $A \rightarrow \alpha$ or $A \rightarrow \beta$
- define $\text{FIRST}(\alpha)$ as the set of tokens that can be first symbol of α , i.e.,
$$a \in \text{FIRST}(\alpha) \text{ iff } \alpha \rightarrow^* a\gamma \text{ for some } \gamma$$
- If $A \rightarrow \alpha \mid \beta$ we want:
$$\text{FIRST}(\alpha) \cap \text{FIRST}(\beta) = \emptyset$$
- If that is always true, we can build a predictive parser.

Computing FIRST(α)

- Given $X := A B C$, $\text{FIRST}(X) = \text{FIRST}(A B C)$
- Can we ignore B or C?
- Consider:

A := a

|

B := b

| A

C := c

Computing FIRST(α)

- Given $X := A B C$, $\text{FIRST}(X) = \text{FIRST}(A B C)$
- Can we ignore B or C?

- Consider:

A := a

|

B := b

| A

C := c

- $\text{FIRST}(X)$ must also include $\text{FIRST}(C)$
- IOW:
 - Must keep track of NTs that are nullable
 - For nullable NTs, determine $\text{FOLLOWS}(\text{NT})$

nullable(A)

- nullable(A) is
 - true if A can derive the empty string
 - false otherwise
- For example:

B := X Y b

X := x

| Y Y

Y :=

In this case, nullable(X) = nullable(Y) = true
nullable(B) = false

FOLLOW(A)

- FOLLOW(A) is the set of terminals that can immediately follow A in a sentential form.
- I.e.,
 $a \in \text{FOLLOW}(A)$ iff $S \Rightarrow^* \alpha A a \beta$ for some α and β

Building a Predictive Parser

- We want to know for each non-terminal which production to choose based on the next input character.
- Build a table with rows labeled by non-terminals, A , and columns labeled by terminals, a . We will put the production, $A := \alpha$, in (A, a) iff
 - $\text{FIRST}(\alpha)$ contains a or
 - $\text{nullable}(\alpha)$ and $\text{FOLLOW}(A)$ contains a



The table for the robot

S := B S F

|

B := b

F := f

	FIRST	FOLLOW	nullable
S	b	\$	yes
B	b	b,f	no
F	f	f,\$	no

	b	f	\$
S			
B			
F			

The table for the robot

$S := B S F$

|

$B := b$

F $\text{FIRST}(BSF) = b$

	FIRST	FOLLOW	nullable
S	b	\$	yes
B	b	b,f	no
F	f	f,\$	no

	b	f	\$
S	$S := BSF$		$S :=$
B	$B := b$		
F		$F := f$	

nullable(ϵ)=true
and
FOLLOW(S) = \$

Table 1

- 1 $S := E$
- 2 $E := TE'$
- 2' $E' := +TE'$
- 3' $E' := -TE'$
- 4' $E' :=$
- 5 $T := FT'$
- 5' $T' := *FT'$
- 6' $T' := /FT'$
- 7' $T' :=$
- 8 $F := id$
- 9 $F := int$

	FIRST	FOLLOW	nullable
S	id, int	\$	
E	id, int	\$	
E'	+, -	\$	yes
T	id, int	+, -, \$	
T'	/, *	+, -, \$	yes
F	id, int	/, *, \$	

	+	-	*	/	id	int	\$
S							
E							
E'							
T							
T'							
F							

Table 1

- 1 $S := E$
- 2 $E := TE'$
- 2' $E' := +TE'$
- 3' $E' := -TE'$
- 4' $E' :=$
- 5 $T := FT'$
- 5' $T' := *FT'$
- 6' $T' := /FT'$
- 7' $T' :=$
- 8 $F := id$
- 9 $F := int$

	FIRST	FOLLOW	nullable
S	id, int	\$	
E	id, int	\$	
E'	+, -	\$	yes
T	id, int	+, -, \$	
T'	/, *	+, -, \$	yes
F	id, int	/, *, \$	

	+	-	*	/	id	int	\$
S					:=E	:=E	
E					:=TE'	:=TE'	
E'	:=+TE'	:= -TE'					:=
T					:=FT'	:=FT'	
T'	:=	:=	:=*FT'	:=/FT'			:=
F					:=id	:=int	

Using the Table

- Each row in the table becomes a function
- For each input token with an entry:
Create a series of invocations that implement the production, where
 - a non-terminal is eaten
 - a terminal becomes a recursive call
- For the blank cells implement errors

Example function

	+	-	*	/	id	int	\$
S					:=E	:=E	
E					:=TE'	:=TE'	
E'	:=+TE'	:= -TE'			:=TE'	:=TE'	:=
T							
T'	:=	:=	:=*FT				
F					:=id	:=int	

How to handle errors?

```

Eprime () {
    switch (token) {
        case PLUS:    eat (PLUS) ; T () ; Eprime () ; break ;
        case MINUS:   eat (MINUS) ; T () ; Eprime () ; break ;
        case ID:      T () ; Eprime () ;
        case INT:     T () ; Eprime () ;
        default:      error () ;
    }
}
    
```

Left-Factoring

- Predictive parsers need to make a choice based on the next terminal.

- Consider:

```
S := if E then S else S
    | if E then S
```

- When looking at **if**, can't decide
- so **left-factor** the grammar

```
S := if E then S X
X := else S
    |
```

Top-Down Parsing

- Can be constructed by hand
- LL(k) grammars can be parsed
 - Left-to-right
 - Leftmost-derivation
 - with k symbols lookahead
- Often requires
 - left-factoring
 - Elimination of left-recursion