

Tokens

Token Type
Examples: ID, NUM, IF, EQUALS, ...

Lexeme
The characters actually matched.
Example:
... if x == **-12.30** then ...

Tokens

Token Type

Examples: ID, NUM, IF, EQUALS, ...

Lexeme

The characters actually matched.

Example:

... if x == **-12.30** then ...

How to describe/specify tokens?

Formal:

Regular Expressions

`Letter (Letter | Digit)*`

Informal:

“// through end of line”

Tokens

Token Type

Examples: ID, NUM, IF, EQUALS, ...

Lexeme

The characters actually matched.

Example:

... if x == **-12.30** then ...

How to describe/specify tokens?

Formal:

Regular Expressions

`Letter (Letter | Digit)*`

Informal:

“// through end of line”

Tokens will appear as TERMINALS in the grammar.

Stmt → while Expr do StmtList endwhile
→ ID “=” Expr “;”
→ ...

Lexical Errors

Most errors tend to be “typos”

Not noticed by the programmer

```
return 1.23;  
return 1,23;
```

... Still results in sequence of legal tokens

```
<ID, "return"> <INT, 1> <COMMA> <INT, 23> <SEMICOLON>
```

No lexical error, but problems during parsing!

Lexical Errors

Most errors tend to be “typos”

Not noticed by the programmer

```
return 1.23;  
return 1,23;
```

... Still results in sequence of legal tokens

```
<ID, "return"> <INT, 1> <COMMA> <INT, 23> <SEMICOLON>
```

No lexical error, but problems during parsing!

Errors caught by lexer:

- EOF within a String / missing ”
- Invalid ASCII character in file
- String / ID exceeds maximum length
- Numerical overflow
- etc...

Lexical Errors

Most errors tend to be “typos”

Not noticed by the programmer

```
return 1.23;  
return 1,23;
```

... Still results in sequence of legal tokens

```
<ID, "return"> <INT, 1> <COMMA> <INT, 23> <SEMICOLON>
```

No lexical error, but problems during parsing!

Errors caught by lexer:

- EOF within a String / missing ”
 - Invalid ASCII character in file
 - String / ID exceeds maximum length
 - Numerical overflow
- etc...

Lexer must keep going!

Always return a valid token.

Skip characters, if necessary.

May confuse the parser

The parser will detect syntax errors and get straightened out (hopefully!)

Managing Input Buffers

Option 1: Read one char from OS at a time.

Option 2: Read N characters per system call

e.g., N = 4096

Manage input buffers in Lexer

More efficient

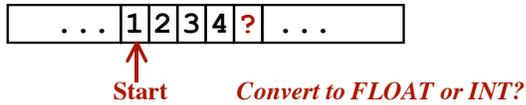
Managing Input Buffers

Option 1: Read one char from OS at a time.

Option 2: Read N characters per system call
e.g., N = 4096

Manage input buffers in Lexer
More efficient

Often, we need to look ahead



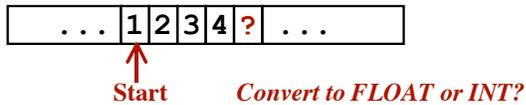
Managing Input Buffers

Option 1: Read one char from OS at a time.

Option 2: Read N characters per system call
e.g., N = 4096

Manage input buffers in Lexer
More efficient

Often, we need to look ahead



Token could overlap / span buffer boundaries.
⇒ need 2 buffers

Code:

```
if (ptr at end of buffer1) or (ptr at end of buffer2) then ...
```

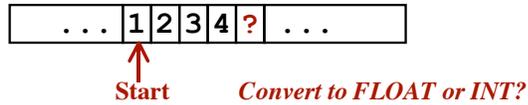
Managing Input Buffers

Option 1: Read one char from OS at a time.

Option 2: Read N characters per system call
e.g., N = 4096

Manage input buffers in Lexer
More efficient

Often, we need to look ahead



Token could overlap / span buffer boundaries.
⇒ need 2 buffers

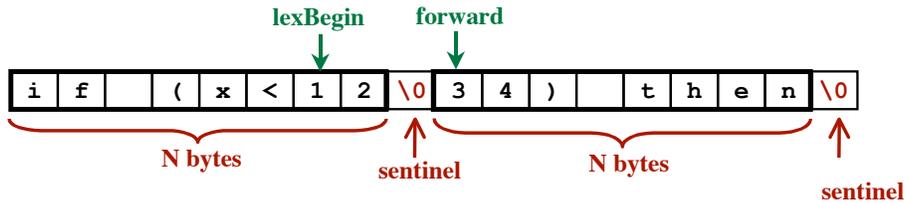
Code:

```
if (ptr at end of buffer1) or (ptr at end of buffer2) then ...
```

Technique: Use “Sentinels” to reduce testing

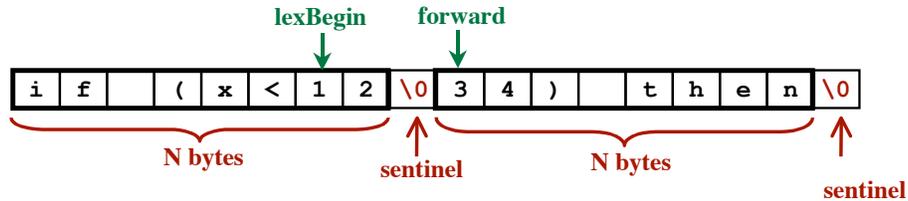
Choose some character that occurs rarely in most inputs

'\0'



Goal: Advance **forward** pointer to next character
...and reload buffer if necessary.

Lexical Analysis - Part 1



Goal: Advance **forward** pointer to next character
...and reload buffer if necessary.

Code :

```
forward++;  
if *forward == '\0' then  
  if forward at end of buffer #1 then  
    Read next N bytes into buffer #2;  
    forward = address of first char of buffer #2;  
  elseif forward at end of buffer #2 then  
    Read next N bytes into buffer #1;  
    forward = address of first char of buffer #1;  
  else  
    // do nothing; a real \0 occurs in the input  
  endif  
endif
```

*One fast test
...which usually fails*

Lexical Analysis - Part 1

“Alphabet” (Σ)

A set of symbols (“characters”)

Examples: $\Sigma = \{ a, b, c, d \}$

$\Sigma =$ ASCII character set

Lexical Analysis - Part 1

“Alphabet” (Σ)

A set of symbols (“characters”)

Examples: $\Sigma = \{ a, b, c, d \}$

$\Sigma =$ ASCII character set

“String” (or “Sentence”)

Sequence of symbols

Finite in length

Example: `abbadc` Length of $s = |s|$

Lexical Analysis - Part 1

“Alphabet” (Σ)

A set of symbols (“characters”)

Examples: $\Sigma = \{ a, b, c, d \}$

$\Sigma =$ ASCII character set

“String” (or “Sentence”)

Sequence of symbols

Finite in length

Example: `abbadc` Length of $s = |s|$

“Empty String” (ϵ , “epsilon”)

It is a string

$|\epsilon| = 0$

Lexical Analysis - Part 1

“Alphabet” (Σ)

A set of symbols (“characters”)

Examples: $\Sigma = \{ a, b, c, d \}$

$\Sigma =$ ASCII character set

“String” (or “Sentence”)

Sequence of symbols

Finite in length

Example: `abbadc` Length of $s = |s|$

“Empty String” (ϵ , “epsilon”)

It is a string

$|\epsilon| = 0$

“Language”

A set of strings

Examples: $L_1 = \{ a, baa, bccb \}$

$L_2 = \{ \}$

$L_3 = \{ \epsilon \}$

$L_4 = \{ \epsilon, ab, abab, ababab, abababab, \dots \}$

$L_5 = \{ s \mid s \text{ can be interpreted as an English sentence}$

$\text{making a true statement about mathematics} \}$

*Each string is finite in length,
but the set may have an infinite
number of elements.*

Lexical Analysis - Part 1

“Prefix” ...of string s

$s =$ `hello`

Prefixes: `he`
`hello`
 `ϵ`

Lexical Analysis - Part 1

“Concatenation”

Strings: x, y

Concatenation: xy

Example:

$x = \mathbf{abb}$

$y = \mathbf{cdc}$

$xy = \mathbf{abbcdc}$

$yx = \mathbf{cdcabb}$

Other notations: $x \parallel y$

$x + y$

$x \mathbf{++} y$

$x \cdot y$

Lexical Analysis - Part 1

“Concatenation”

Strings: x, y

Concatenation: xy

Example:

$x = \mathbf{abb}$

$y = \mathbf{cdc}$

$xy = \mathbf{abbcdc}$

$yx = \mathbf{cdcabb}$

What is the “identity” for concatenation?

$\epsilon x = x\epsilon = x$

Multiplication \Leftrightarrow Concatenation

Exponentiation $\Leftrightarrow ?$

Other notations: $x \parallel y$

$x + y$

$x \mathbf{++} y$

$x \cdot y$

Lexical Analysis - Part 1

“Concatenation”
Strings: x, y
Concatenation: xy
Example:
 $x = \mathbf{abb}$
 $y = \mathbf{cdc}$
 $xy = \mathbf{abbc dc}$
 $yx = \mathbf{cdcabb}$

What is the “identity” for concatenation?
 $\epsilon x = x\epsilon = x$

Multiplication \Leftrightarrow Concatenation
Exponentiation $\Leftrightarrow ?$

Define $s^0 = \epsilon$
 $s^N = s^{N-1}s$

Example $x = \mathbf{ab}$
 $x^0 = \epsilon$
 $x^1 = x = \mathbf{ab}$
 $x^2 = xx = \mathbf{abab}$
 $x^3 = xxx = \mathbf{ababab}$
...etc...

Other notations: $x \parallel y$
 $x + y$
 $x ++ y$
 $x \cdot y$

© Harry H. Porter, 2005 25

Lexical Analysis - Part 1

“Concatenation”
Strings: x, y
Concatenation: xy
Example:
 $x = \mathbf{abb}$
 $y = \mathbf{cdc}$
 $xy = \mathbf{abbc dc}$
 $yx = \mathbf{cdcabb}$

What is the “identity” for concatenation?
 $\epsilon x = x\epsilon = x$

Multiplication \Leftrightarrow Concatenation
Exponentiation $\Leftrightarrow ?$

Define $s^0 = \epsilon$
 $s^N = s^{N-1}s$

Example $x = \mathbf{ab}$
 $x^0 = \epsilon$
 $x^1 = x = \mathbf{ab}$
 $x^2 = xx = \mathbf{abab}$
 $x^3 = xxx = \mathbf{ababab}$
...etc...
 $x^* = x^\infty = \mathbf{abababababab\dots}$

Infinite sequence of symbols!
Technically, this is not a “string”

© Harry H. Porter, 2005 26

Lexical Analysis - Part 1

“Language”
A set of strings
 $L = \{ \dots \}$
 $M = \{ \dots \}$

Generally, these are infinite sets.

Lexical Analysis - Part 1

“Language”
A set of strings
 $L = \{ \dots \}$
 $M = \{ \dots \}$

Generally, these are infinite sets.

“Union” of two languages
 $L \cup M = \{ s \mid s \text{ is in } L \text{ or is in } M \}$

Example:
 $L = \{ a, ab \}$
 $M = \{ c, dd \}$
 $L \cup M = \{ a, ab, c, dd \}$

“Language”
 A set of strings
 $L = \{ \dots \}$
 $M = \{ \dots \}$

Generally, these are infinite sets.

“Union” of two languages
 $L \cup M = \{ s \mid s \text{ is in } L \text{ or is in } M \}$

Example:
 $L = \{ a, ab \}$
 $M = \{ c, dd \}$
 $L \cup M = \{ a, ab, c, dd \}$

“Concatenation” of two languages
 $LM = \{ st \mid s \in L \text{ and } t \in M \}$

Example:
 $L = \{ a, ab \}$
 $M = \{ c, dd \}$
 $LM = \{ ac, add, abc, abdd \}$

Repeated Concatenation

Let: $L = \{ a, bc \}$

Example: $L^0 = \{ \epsilon \}$
 $L^1 = L = \{ a, bc \}$
 $L^2 = LL = \{ aa, abc, bca, bcbc \}$
 $L^3 = LLL = \{ aaa, aabc, abca, abcbc, bcaa, bcabc, bcbca, bcbcbc \}$
 ...etc...
 $L^N = L^{N-1}L = LL^{N-1}$

Kleene Closure

Let: $L = \{ a, bc \}$

Example: $L^0 = \{ \epsilon \}$

$L^1 = L = \{ a, bc \}$

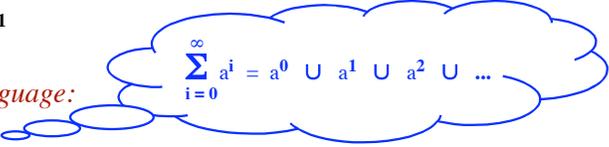
$L^2 = LL = \{ aa, abc, bca, bcbc \}$

$L^3 = LLL = \{ aaa, aabc, abca, abcbc, bcaa, bcabc, bcbca, bcbcbc \}$

...etc...

$L^N = L^{N-1}L = LL^{N-1}$

The "Kleene Closure" of a language:



$$\sum_{i=0}^{\infty} a^i = a^0 \cup a^1 \cup a^2 \cup \dots$$

$$L^* = \bigcup_{i=0}^{\infty} L^i = L^0 \cup L^1 \cup L^2 \cup L^3 \cup \dots$$

Example:

$$L^* = \{ \underbrace{\epsilon}_{L^0}, \underbrace{a, bc}_{L^1}, \underbrace{aa, abc, bca, bcbc}_{L^2}, \underbrace{aaa, aabc, abca, abcbc, \dots}_{L^3}, \dots \}$$

Positive Closure

Let: $L = \{ a, bc \}$

Example: $L^0 = \{ \epsilon \}$

$L^1 = L = \{ a, bc \}$

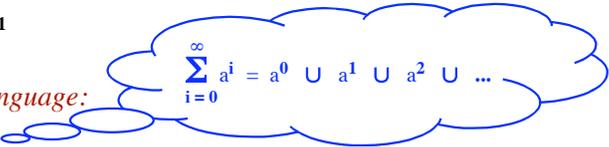
$L^2 = LL = \{ aa, abc, bca, bcbc \}$

$L^3 = LLL = \{ aaa, aabc, abca, abcbc, bcaa, bcabc, bcbca, bcbcbc \}$

...etc...

$L^N = L^{N-1}L = LL^{N-1}$

The "Positive Closure" of a language:



$$\sum_{i=1}^{\infty} a^i = a^1 \cup a^2 \cup \dots$$

$$L^+ = \bigcup_{i=1}^{\infty} L^i = L^1 \cup L^2 \cup L^3 \cup \dots$$

Example:

$$L^+ = \{ \underbrace{a, bc}_{L^1}, \underbrace{aa, abc, bca, bcbc}_{L^2}, \underbrace{aaa, aabc, abca, abcbc, \dots}_{L^3}, \dots \}$$

Positive Closure

Let: $L = \{ a, bc \}$

Example: $L^0 = \{ \epsilon \}$

$L^1 = L = \{ a, bc \}$

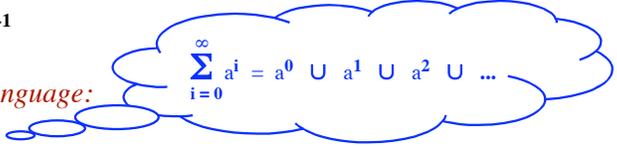
$L^2 = LL = \{ aa, abc, bca, bcbc \}$

$L^3 = LLL = \{ aaa, aabc, abca, abcbc, bcaa, bcabc, bcbca, bcbcbc \}$

...etc...

$L^N = L^{N-1}L = LL^{N-1}$

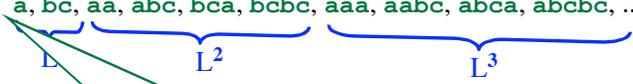
The "Positive Closure" of a language:



$$L^+ = \bigcup_{i=1}^{\infty} L^i = L^1 \cup L^2 \cup L^3 \cup \dots$$

Example:

$L^+ = \{ a, bc, aa, abc, bca, bcbc, aaa, aabc, abca, abcbc, \dots \}$



Note that ϵ is not included
UNLESS it is in L to start with

Examples

Let: $L = \{ a, b, c, \dots, z \}$

$D = \{ 0, 1, 2, \dots, 9 \}$

$D^+ =$

Examples

Let: $L = \{ a, b, c, \dots, z \}$
 $D = \{ 0, 1, 2, \dots, 9 \}$

$D^+ =$
 “The set of strings with one or more digits”

$L \cup D =$

Examples

Let: $L = \{ a, b, c, \dots, z \}$
 $D = \{ 0, 1, 2, \dots, 9 \}$

$D^+ =$
 “The set of strings with one or more digits”

$L \cup D =$
 “The set of alphanumeric characters”
 $\{ a, b, c, \dots, z, 0, 1, 2, \dots, 9 \}$

$(L \cup D)^* =$

Examples

Let: $L = \{ \mathbf{a, b, c, \dots, z} \}$
 $D = \{ \mathbf{0, 1, 2, \dots, 9} \}$

$D^+ =$
 “The set of strings with one or more digits”

$L \cup D =$
 “The set of alphanumeric characters”
 $\{ \mathbf{a, b, c, \dots, z, 0, 1, 2, \dots, 9} \}$

$(L \cup D)^* =$
 “Sequences of zero or more letters and digits”

$L (L \cup D)^* =$

Examples

Let: $L = \{ \mathbf{a, b, c, \dots, z} \}$
 $D = \{ \mathbf{0, 1, 2, \dots, 9} \}$

$D^+ =$
 “The set of strings with one or more digits”

$L \cup D =$
 “The set of alphanumeric characters”
 $\{ \mathbf{a, b, c, \dots, z, 0, 1, 2, \dots, 9} \}$

$(L \cup D)^* =$
 “Sequences of zero or more letters and digits”

$L ((L \cup D)^*) =$

Examples

Let: $L = \{ a, b, c, \dots, z \}$
 $D = \{ 0, 1, 2, \dots, 9 \}$

$D^+ =$
 “The set of strings with one or more digits”

$L \cup D =$
 “The set of alphanumeric characters”
 $\{ a, b, c, \dots, z, 0, 1, 2, \dots, 9 \}$

$(L \cup D)^* =$
 “Sequences of zero or more letters and digits”

$L((L \cup D)^*) =$
 “Set of strings that start with a letter, followed by zero or more letters and and digits.”

Regular Expressions

Assume the alphabet is given... *e.g.*, $\Sigma = \{ a, b, c, \dots, z \}$

Example: $a (b | c) d^* e$

A regular expression describes a language.

Notation:

r = regular expression
 L(r) = the corresponding language

Example:

$r = a (b | c) d^* e$
 $L(r) = \{$
 $abe,$
 $abde,$
 $abdde,$
 $abddde,$
 $\dots,$
 $ace,$
 $acde,$
 $acdde,$
 $acddee,$
 $\dots \}$

Meta Symbols:

()
|
*
 ϵ

How to “Parse” Regular Expressions

- * has highest precedence.
- Concatenation as middle precedence.
- | has lowest precedence.
- Use parentheses to override these rules.

How to “Parse” Regular Expressions

- * has highest precedence.
- Concatenation as middle precedence.
- | has lowest precedence.
- Use parentheses to override these rules.

Examples:

a b* =

How to “Parse” Regular Expressions

* has highest precedence.
Concatenation as middle precedence.
| has lowest precedence.
Use parentheses to override these rules.

Examples:

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

If you want $(a b)^*$ you must use parentheses.

How to “Parse” Regular Expressions

* has highest precedence.
Concatenation as middle precedence.
| has lowest precedence.
Use parentheses to override these rules.

Examples:

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

If you want $(a b)^*$ you must use parentheses.

$a | b c =$

How to “Parse” Regular Expressions

* has highest precedence.
Concatenation as middle precedence.
| has lowest precedence.
Use parentheses to override these rules.

Examples:

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

If you want $(a b)^*$ you must use parentheses.

$a | b c = a | (b c)$

If you want $(a | b) c$ you must use parentheses.

How to “Parse” Regular Expressions

* has highest precedence.
Concatenation as middle precedence.
| has lowest precedence.
Use parentheses to override these rules.

Examples:

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

If you want $(a b)^*$ you must use parentheses.

$a | b c = a | (b c)$

If you want $(a | b) c$ you must use parentheses.

Concatenation and | are associative.

$(a b) c = a (b c) = a b c$

$(a | b) | c = a | (b | c) = a | b | c$

How to “Parse” Regular Expressions

* has highest precedence.
 Concatenation as middle precedence.
 | has lowest precedence.
 Use parentheses to override these rules.

Examples:

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

If you want $(a b)^*$ you must use parentheses.

$a | b c = a | (b c)$

If you want $(a | b) c$ you must use parentheses.

Concatenation and | are associative.

$(a b) c = a (b c) = a b c$

$(a | b) | c = a | (b | c) = a | b | c$

Example:

$b d | e f^* | g a =$

How to “Parse” Regular Expressions

* has highest precedence.
 Concatenation as middle precedence.
 | has lowest precedence.
 Use parentheses to override these rules.

Examples:

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

If you want $(a b)^*$ you must use parentheses.

$a | b c = a | (b c)$

If you want $(a | b) c$ you must use parentheses.

Concatenation and | are associative.

$(a b) c = a (b c) = a b c$

$(a | b) | c = a | (b | c) = a | b | c$

Example:

$b d | e f^* | g a = b d | e (f^*) | g a$

How to “Parse” Regular Expressions

* has highest precedence.
 Concatenation as middle precedence.
 | has lowest precedence.
 Use parentheses to override these rules.

Examples:

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

If you want $(a b)^*$ you must use parentheses.

$$a | b c = a | (b c)$$

If you want $(a | b) c$ you must use parentheses.

Concatenation and | are associative.

$$(a b) c = a (b c) = a b c$$

$$(a | b) | c = a | (b | c) = a | b | c$$

Example:

$$b d | e f^* | g a = (b d) | (e (f^*)) | (g a)$$

How to “Parse” Regular Expressions

* has highest precedence.
 Concatenation as middle precedence.
 | has lowest precedence.
 Use parentheses to override these rules.

Examples:

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

If you want $(a b)^*$ you must use parentheses.

$$a | b c = a | (b c)$$

If you want $(a | b) c$ you must use parentheses.

Concatenation and | are associative.

$$(a b) c = a (b c) = a b c$$

$$(a | b) | c = a | (b | c) = a | b | c$$

Example:

$$b d | e f^* | g a = ((b d) | (e (f^*))) | (g a)$$

Fully parenthesized

Definition: Regular Expressions

(Over alphabet Σ)

- ϵ is a regular expression.
- If a is a symbol (i.e., if $a \in \Sigma$), then a is a regular expression.
- If R and S are regular expressions, then $R|S$ is a regular expression.
- If R and S are regular expressions, then RS is a regular expression.
- If R is a regular expression, then R^* is a regular expression.
- If R is a regular expression, then (R) is a regular expression.

Definition: Regular Expressions

(Over alphabet Σ)

And, given a regular expression R , what is $L(R)$?

- ϵ is a regular expression.
- If a is a symbol (i.e., if $a \in \Sigma$), then a is a regular expression.
- If R and S are regular expressions, then $R|S$ is a regular expression.
- If R and S are regular expressions, then RS is a regular expression.
- If R is a regular expression, then R^* is a regular expression.
- If R is a regular expression, then (R) is a regular expression.

Definition: Regular Expressions

(Over alphabet Σ)

And, given a regular expression R , what is $L(R)$?

- ϵ is a regular expression.
 $L(\epsilon) = \{ \epsilon \}$
- If a is a symbol (i.e., if $a \in \Sigma$), then a is a regular expression.
- If R and S are regular expressions, then $R|S$ is a regular expression.
- If R and S are regular expressions, then RS is a regular expression.
- If R is a regular expression, then R^* is a regular expression.
- If R is a regular expression, then (R) is a regular expression.

Definition: Regular Expressions

(Over alphabet Σ)

And, given a regular expression R , what is $L(R)$?

- ϵ is a regular expression.
 $L(\epsilon) = \{ \epsilon \}$
- If a is a symbol (i.e., if $a \in \Sigma$), then a is a regular expression.
 $L(a) = \{ a \}$
- If R and S are regular expressions, then $R|S$ is a regular expression.
- If R and S are regular expressions, then RS is a regular expression.
- If R is a regular expression, then R^* is a regular expression.
- If R is a regular expression, then (R) is a regular expression.

Definition: Regular Expressions

(Over alphabet Σ)

And, given a regular expression R , what is $L(R)$?

- ϵ is a regular expression.
 $L(\epsilon) = \{ \epsilon \}$
- If a is a symbol (i.e., if $a \in \Sigma$), then a is a regular expression.
 $L(a) = \{ a \}$
- If R and S are regular expressions, then $R|S$ is a regular expression.
 $L(R|S) = L(R) \cup L(S)$
- If R and S are regular expressions, then RS is a regular expression.
- If R is a regular expression, then R^* is a regular expression.
- If R is a regular expression, then (R) is a regular expression.

Definition: Regular Expressions

(Over alphabet Σ)

And, given a regular expression R , what is $L(R)$?

- ϵ is a regular expression.
 $L(\epsilon) = \{ \epsilon \}$
- If a is a symbol (i.e., if $a \in \Sigma$), then a is a regular expression.
 $L(a) = \{ a \}$
- If R and S are regular expressions, then $R|S$ is a regular expression.
 $L(R|S) = L(R) \cup L(S)$
- If R and S are regular expressions, then RS is a regular expression.
 $L(RS) = L(R) L(S)$
- If R is a regular expression, then R^* is a regular expression.
- If R is a regular expression, then (R) is a regular expression.

Definition: Regular Expressions

(Over alphabet Σ)

And, given a regular expression R , what is $L(R)$?

- ϵ is a regular expression.
 $L(\epsilon) = \{ \epsilon \}$
- If a is a symbol (i.e., if $a \in \Sigma$), then a is a regular expression.
 $L(a) = \{ a \}$
- If R and S are regular expressions, then $R|S$ is a regular expression.
 $L(R|S) = L(R) \cup L(S)$
- If R and S are regular expressions, then RS is a regular expression.
 $L(RS) = L(R) L(S)$
- If R is a regular expression, then R^* is a regular expression.
 $L(R^*) = (L(R))^*$
- If R is a regular expression, then (R) is a regular expression.

Definition: Regular Expressions

(Over alphabet Σ)

And, given a regular expression R , what is $L(R)$?

- ϵ is a regular expression.
 $L(\epsilon) = \{ \epsilon \}$
- If a is a symbol (i.e., if $a \in \Sigma$), then a is a regular expression.
 $L(a) = \{ a \}$
- If R and S are regular expressions, then $R|S$ is a regular expression.
 $L(R|S) = L(R) \cup L(S)$
- If R and S are regular expressions, then RS is a regular expression.
 $L(RS) = L(R) L(S)$
- If R is a regular expression, then R^* is a regular expression.
 $L(R^*) = (L(R))^*$
- If R is a regular expression, then (R) is a regular expression.
 $L((R)) = L(R)$

Regular Languages

Definition: “Regular Language” (or “Regular Set”)

... A language that can be described by a regular expression.

- Any finite language (i.e., finite set of strings) is a regular language.
- Regular languages are (usually) infinite.
- Regular languages are, in some sense, simple languages.

Regular Languages \subset Context-Free Languages

Examples:

$a \mid b \mid cab$	$\{a, b, cab\}$
b^*	$\{\epsilon, b, bb, bbb, \dots\}$
$a \mid b^*$	$\{a, \epsilon, b, bb, bbb, \dots\}$
$(a \mid b)^*$	$\{\epsilon, a, b, aa, ab, ba, bb, aaa, \dots\}$ “Set of all strings of a’s and b’s, including ϵ .”

Equality v. Equivalence

Are these regular expressions equal?

$$R = a a^* (b \mid c)$$

$$S = a^* a (c \mid b)$$

... No!

Yet, they describe the same language.

$$L(R) = L(S)$$

“Equivalence” of regular expressions

If $L(R) = L(S)$ then we say $R \approx S$

“R is equivalent to S”

“Syntactic” equality versus “deeper” equality...

Algebra:

$$\text{Does... } x(3+b) = 3x+bx \quad ?$$

From now on, we’ll just say $R = S$ to mean $R \approx S$

Notation:

Equality

=

Equivalence

=

≡

≈

≡

↔

Algebraic Laws of Regular Expressions

Let R, S, T be regular expressions...

I is commutative
 $RIS = SIR$

I is associative
 $RI(SIT) = (RIS)IT = RISIT$

Concatenation is associative
 $R(ST) = (RS)T = RST$

Concatenation distributes over I
 $R(SIT) = RSIRT$
 $(RIS)T = RTIST$

ϵ is the identity for concatenation
 $\epsilon R = R\epsilon = R$

* is idempotent
 $(R^*)^* = R^*$

Relation between * and ϵ
 $R^* = (R | \epsilon)^*$

Preferred

Preferred

Regular Definitions

Letter = a | b | c | ... | z

Digit = 0 | 1 | 2 | ... | 9

ID = Letter (Letter | Digit)*

Regular Definitions

Letter = a | b | c | ... | z
Digit = 0 | 1 | 2 | ... | 9
ID = Letter (Letter | Digit) *

Names (e.g., Letter) are underlined to distinguish from a sequence of symbols.

Letter (Letter | Digit) *
 = { "Letter", "LetterLetter", "LetterDigit", ... }

Regular Definitions

Letter = a | b | c | ... | z
Digit = 0 | 1 | 2 | ... | 9
ID = Letter (Letter | Digit) *

Names (e.g., Letter) are underlined to distinguish from a sequence of symbols.

Letter (Letter | Digit) *
 = { "Letter", "LetterLetter", "LetterDigit", ... }

Each definition may only use names *previously* defined.

- ⇒ No recursion
- Regular Sets = no recursion
- CFG = recursion

Addition Notation / Shorthand

One-or-more: +

$X^+ = X(X^*)$

Digit⁺ = Digit Digit^{*} = Digits

Addition Notation / Shorthand

One-or-more: +

$X^+ = X(X^*)$

Digit⁺ = Digit Digit^{*} = Digits

Optional (zero-or-one): ?

$X? = (X | \epsilon)$

Num = Digit⁺ (. Digit⁺)?

Addition Notation / Shorthand

One-or-more: +

$$X^+ = X(X^*)$$

$$\underline{\text{Digit}}^+ = \underline{\text{Digit}} \underline{\text{Digit}}^* = \underline{\text{Digits}}$$

Optional (zero-or-one): ?

$$X? = (X | \epsilon)$$

$$\underline{\text{Num}} = \underline{\text{Digit}}^+ (. \underline{\text{Digit}}^+)?$$

Character Classes: [FirstChar-LastChar]

Assumption: The underlying alphabet is known ...and is ordered.

$$\underline{\text{Digit}} = [0-9]$$

$$\underline{\text{Letter}} = [a-zA-Z] = [A-Za-z]$$

Addition Notation / Shorthand

One-or-more: +

$$X^+ = X(X^*)$$

$$\underline{\text{Digit}}^+ = \underline{\text{Digit}} \underline{\text{Digit}}^* = \underline{\text{Digits}}$$

Optional (zero-or-one): ?

$$X? = (X | \epsilon)$$

$$\underline{\text{Num}} = \underline{\text{Digit}}^+ (. \underline{\text{Digit}}^+)?$$

Character Classes: [FirstChar-LastChar]

Assumption: The underlying alphabet is known ...and is ordered.

$$\underline{\text{Digit}} = [0-9]$$

$$\underline{\text{Letter}} = [a-zA-Z] = [A-Za-z]$$

Variations:

Zero-or-more: $ab^*c = a\{b\}c = a\{b\}^*c$

One-or-more: $ab^+c = a\{b\}^+c$

Optional: $ab?c = a\{b\}c$

What does
ab...bc
mean?

Lexical Analysis - Part 1

Many sets of strings are not regular.
...no regular expression for them!

Lexical Analysis - Part 1

Many sets of strings are not regular.
...no regular expression for them!

The set of all strings in which parentheses are balanced.

`((((())))`

Must use a CFG!

Lexical Analysis - Part 1

Many sets of strings are not regular.
...no regular expression for them!

The set of all strings in which parentheses are balanced.

$(() (())$

Must use a CFG!

Strings with repeated substrings

$\{ XcX \mid X \text{ is a string of } a\text{'s and } b\text{'s} \}$

$abbbabcabbbab$

CFG is not even powerful enough.

Lexical Analysis - Part 1

Many sets of strings are not regular.
...no regular expression for them!

The set of all strings in which parentheses are balanced.

$(() (())$

Must use a CFG!

Strings with repeated substrings

$\{ XcX \mid X \text{ is a string of } a\text{'s and } b\text{'s} \}$

$abbbabcabbbab$

CFG is not even powerful enough.

The Problem?

In order to recognize a string,
these languages require *memory*!

Lexical Analysis - Part 1

Problem: How to describe tokens?

Solution: Regular Expressions

Problem: How to recognize tokens?

Approaches:

- Hand-coded routines
 Examples: E-Language, PCAT-Lexer
- Finite State Automata
- Scanner Generators (Java: JLex, C: Lex)

Scanner Generators

Input: Sequence of regular definitions

Output: A lexer (e.g., a program in Java or “C”)

Approach:

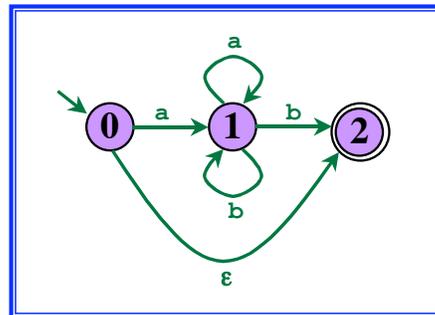
- Read in regular expressions
- Convert into a Finite State Automaton (FSA)
- Optimize the FSA
- Represent the FSA with tables / arrays
- Generate a table-driven lexer (Combine “canned” code with tables.)

Lexical Analysis - Part 1

Finite State Automata (FSAs)

(“Finite State Machines”, “Finite Automata”, “FA”)

- One start state
- Many final states
- Each state is labeled with a state name
- Directed edges, labeled with symbols
 - **Deterministic** (DFA)
 - No ϵ -edges
 - Each outgoing edge has different symbol
 - **Non-deterministic** (NFA)



Finite State Automata (FSAs)

Formalism: $\langle S, \Sigma, \delta, S_0, S_F \rangle$

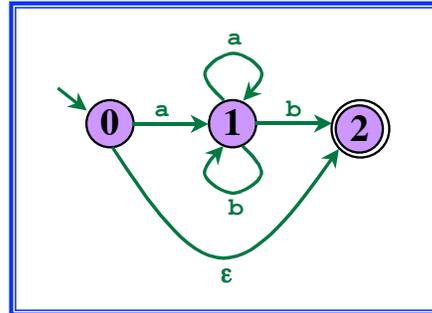
S = Set of states
 $S = \{s_0, s_1, \dots, s_N\}$

Σ = Input Alphabet
 $\Sigma = \text{ASCII Characters}$

δ = Transition Function
 $S \times \Sigma \rightarrow \text{States (deterministic)}$
 $S \times \Sigma \rightarrow \text{Sets of States (non-deterministic)}$

s_0 = Start State
 "Initial state"
 $s_0 \in S$

S_F = Set of final states
 "accepting states"
 $S_F \subseteq S$



Finite State Automata (FSAs)

Formalism: $\langle S, \Sigma, \delta, S_0, S_F \rangle$

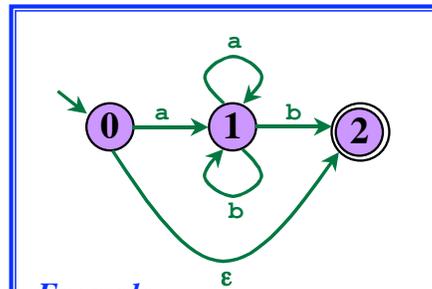
S = Set of states
 $S = \{s_0, s_1, \dots, s_N\}$

Σ = Input Alphabet
 $\Sigma = \text{ASCII Characters}$

δ = Transition Function
 $S \times \Sigma \rightarrow \text{States (deterministic)}$
 $S \times \Sigma \rightarrow \text{Sets of States (non-deterministic)}$

s_0 = Start State
 "Initial state"
 $s_0 \in S$

S_F = Set of final states
 "accepting states"
 $S_F \subseteq S$



Example:

$S = \{0, 1, 2\}$

$\Sigma = \{a, b\}$

$s_0 = 0$

$S_F = \{2\}$

$\delta =$

		Input Symbols		
		a	b	ε
States	0	{1}	{}	{2}
	1	{1}	{1,2}	{}
	2	{}	{}	{}

Finite State Automata (FSAs)

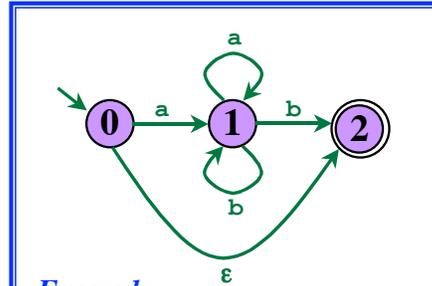
A string is “**accepted**”...
(a string is “**recognized**”...)

by a FSA if there is a path
from Start to any accepting state
where edge labels match the string.

Example:

This FSA accepts:

ϵ
aaab
abbb



Example:

$S = \{0, 1, 2\}$

$\Sigma = \{a, b\}$

$s_0 = 0$

$S_F = \{2\}$

$\delta =$

		<i>Input Symbols</i>		
		a	b	ϵ
{	<i>States</i> 0	{1}	{}	{2}
	1	{1}	{1,2}	{}
	2	{}	{}	{}

Deterministic Finite Automata (DFAs)

No ϵ -moves

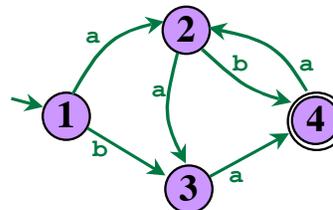
The transition function returns a single state

$\delta: S \times \Sigma \rightarrow S$

function Move (s:State, a:Symbol) returns State

$\delta =$

		<i>Input Symbols</i>	
		a	b
{	1	2	3
	2	3	4
	3	4	---
	4	2	---



Deterministic Finite Automata (DFAs)

No ϵ -moves

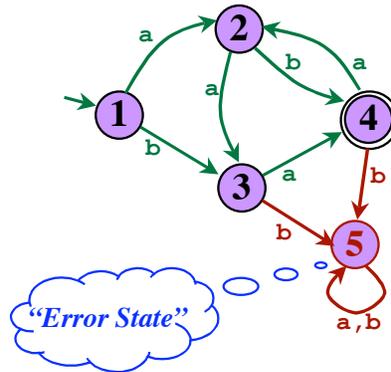
The transition function returns a single state

$$\delta: S \times \Sigma \rightarrow S$$

function Move (s:State, a:Symbol) returns State

$\delta =$

		Input Symbols	
		a	b
States	1	2	3
	2	3	4
	3	4	5
	4	2	5
	5	5	5



Deterministic Finite Automata (DFAs)

No ϵ -moves

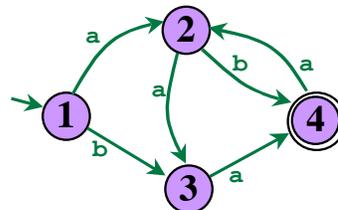
The transition function returns a single state

$$\delta: S \times \Sigma \rightarrow S$$

function Move (s:State, a:Symbol) returns State

$\delta =$

		Input Symbols	
		a	b
States	1	2	3
	2	3	4
	3	4	---
	4	2	---



Non-Deterministic Finite Automata (NFAs)

Allow ϵ -moves

The transition function returns a *set of* states

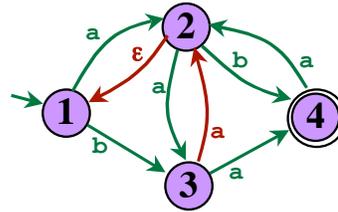
$$\delta: S \times \Sigma \rightarrow \text{Powerset}(S)$$

$$\delta: S \times \Sigma \rightarrow P(S)$$

function Move (s:State, a:Symbol) returns *set of* State

$\delta =$

	<i>Input Symbols</i>		
	a	b	ϵ
1	{2}	{3}	{}
2	{3}	{4}	{1}
3	{4,2}	{}	{}
4	{2}	{}	{}



Theoretical Results

- The set of strings recognized by an NFA can be described by a Regular Expression.

Theoretical Results

- The set of strings recognized by an NFA can be described by a Regular Expression.
- The set of strings described by a Regular Expression can be recognized by an NFA.

Theoretical Results

- The set of strings recognized by an NFA can be described by a Regular Expression.
- The set of strings described by a Regular Expression can be recognized by an NFA.
- The set of strings recognized by an DFA can be described by a Regular Expression.

Theoretical Results

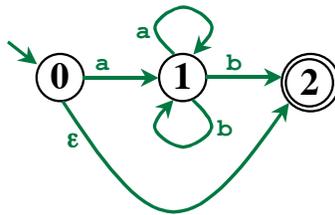
- The set of strings recognized by an NFA can be described by a Regular Expression.
- The set of strings described by a Regular Expression can be recognized by an NFA.
- The set of strings recognized by a DFA can be described by a Regular Expression.
- The set of strings described by a Regular Expression can be recognized by a DFA.

Theoretical Results

- The set of strings recognized by an NFA can be described by a Regular Expression.
- The set of strings described by a Regular Expression can be recognized by an NFA.
- The set of strings recognized by a DFA can be described by a Regular Expression.
- The set of strings described by a Regular Expression can be recognized by a DFA.
- DFAs, NFAs, and Regular Expressions all have the same “power”. They describe “Regular Sets” (“Regular Languages”)

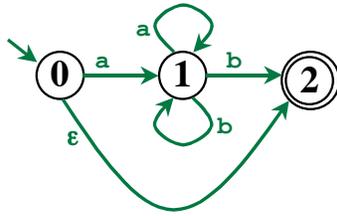
Theoretical Results

- The set of strings recognized by an NFA can be described by a Regular Expression.
- The set of strings described by a Regular Expression can be recognized by an NFA.
- The set of strings recognized by a DFA can be described by a Regular Expression.
- The set of strings described by a Regular Expression can be recognized by a DFA.
- DFAs, NFAs, and Regular Expressions all have the same “power”. They describe “Regular Sets” (“Regular Languages”).
- The DFA may have a lot more states than the NFA. (May have exponentially as many states, but...)



What is the regular expression?

Lexical Analysis - Part 1

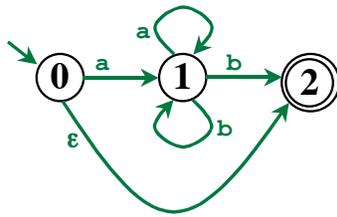


What is the regular expression?

$\epsilon \mid a (a|b)^* b$

What is an equivalent DFA?

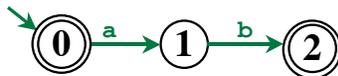
Lexical Analysis - Part 1



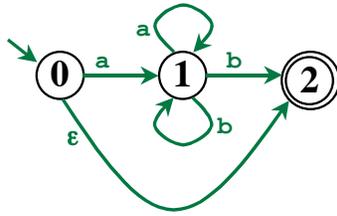
What is the regular expression?

$\epsilon \mid a (a|b)^* b$

What is an equivalent DFA?



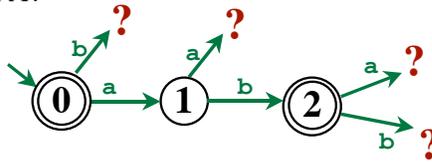
Lexical Analysis - Part 1



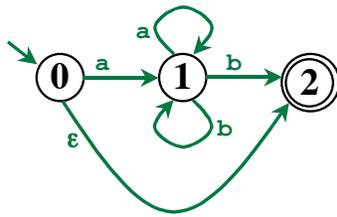
What is the regular expression?

$\epsilon \mid a (a|b)^* b$

What is an equivalent DFA?



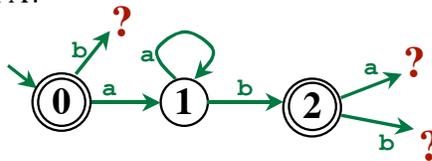
Lexical Analysis - Part 1



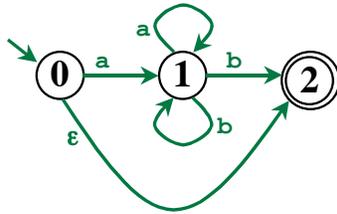
What is the regular expression?

$\epsilon \mid a (a|b)^* b$

What is an equivalent DFA?



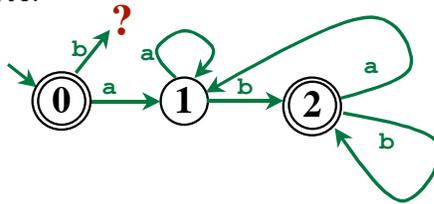
Lexical Analysis - Part 1



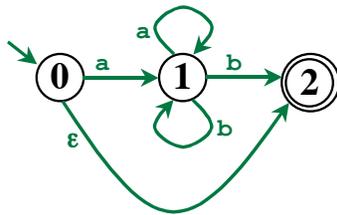
What is the regular expression?

$\epsilon \mid a (a|b)^* b$

What is an equivalent DFA?



Lexical Analysis - Part 1



What is the regular expression?

$\epsilon \mid a (a|b)^* b$

What is an equivalent DFA?

