

Parsing

Sujit Kumar Chakrabarti

IITB

Parsing

- Process of determining if the input belongs to the language of the grammar
- Builds a *parse tree*

Parsing

Example

12:10:45

Specification: A number followed by a colon followed by a number followed by a colon followed by a number.

Regular expression: num COLON num COLON num

Parsing

Example

12:10:45

11:09:22

...

Specification: A sequence of a number followed by a colon followed by a number followed by a colon followed by a number.

Regular expression: (num COLON num COLON num)+

Parsing

Example

()
() ()
(()) ()

Specification: A language of balanced parentheses

Regular expression: ?

Parsing

Example

()
(() ())
((())) ()

Specification: A language of balanced parentheses

Regular expression: ?

Parsing algorithm:

```
procedure BALANCED-PARENTHESES(buffer)  
  level  $\leftarrow$  0  
  while buffer has more characters do  
    c  $\leftarrow$  NEXTCHAR(buffer)  
    if c = LPAREN then  
      level  $\leftarrow$  level + 1  
    else if c = RPAREN then  
      level  $\leftarrow$  level - 1  
    if level < 0 then  
      return false  
  if level = 0 then  
    return true  
  else  
    return false
```

Parsing

Examples of structures that can't be expressed using regular expressions

```
(* ... (* ... (* ... *)*)*)
```

```
let ... in  
let ... in  
let in e
```

```
if(...) {  
  if(...) {  
    ...  
  }  
}
```

Context Free Grammar¹

Example: Balanced parentheses

¹Backus Naur Form (BNF)

Context Free Grammar¹

Example: Balanced parentheses

$$\begin{array}{lcl} S & \rightarrow & \epsilon \\ S & \rightarrow & (S) \\ S & \rightarrow & S S \end{array}$$

- **Components of a grammar:** Rules/productions, terminals, non-terminals, start symbol
- **Meta-language:** language in which the grammar is written; terminals, non-terminals are grammar-symbols/tokens in the meta-language.
- **Notational variances:** ::= instead of \rightarrow
- Could be rewritten as:

$$\begin{array}{lcl} S & \rightarrow & \epsilon \\ & | & (S) \\ & | & S S \end{array}$$

¹Backus Naur Form (BNF)

Parsing

Grammar

Example: Grammar for arithmetic expressions

Parsing

Grammar

Example: Grammar for arithmetic expressions

$$E \rightarrow E + E$$

$$E \rightarrow E * E$$

$$E \rightarrow (E)$$

$$E \rightarrow 0 \mid 1 \dots \mid 9$$

Grammar

Example: Grammar for function call

Grammar

Example: Grammar for function call

<i>fcall</i>	→	id (arglist)
<i>arglist</i>	→	<i>arg</i> *
<i>arg</i>	→	<i>exp</i>
<i>exp</i>	→	... SL ...

Note:

- This grammar contains multiple non-terminals.
- *arglist* → *arg** shorthand for

<i>arglist</i>	→	ε
<i>arglist</i>	→	<i>arg arglist</i>

Parsing

Derivation – Verifying $i \in L$

$((())) ()$

Parsing

Derivation – Verifying $i \in L$

$((())) ()$

- S
- $S S$
- $(S)(S)$
- $((S))(\epsilon)$
- $((((S))))()$
- $((((\epsilon))))()$
- $((()))()$

Parsing

Derivation

Parsing

Derivation

Example: Derivations for arithmetic expressions

- 5
- $1 + 2$
- $1 + 2 * 3$

Parsing

Parse Tree

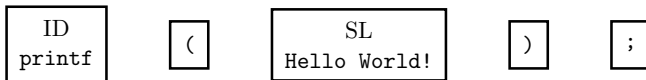
```
printf(" Hello World!");
```

Parsing

Parse Tree

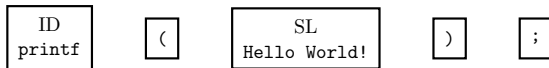
```
printf(" Hello World! ");
```

After lexical analysis:



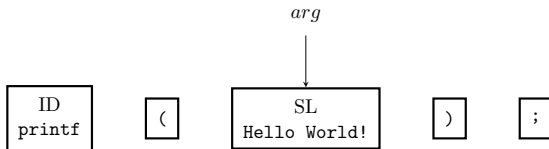
Parsing

Parse Tree



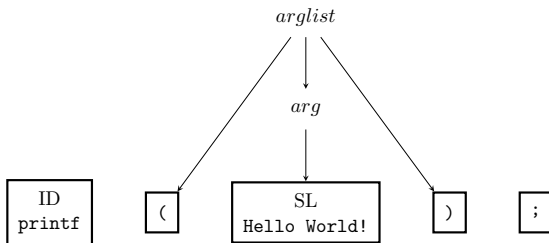
Parsing

Parse Tree



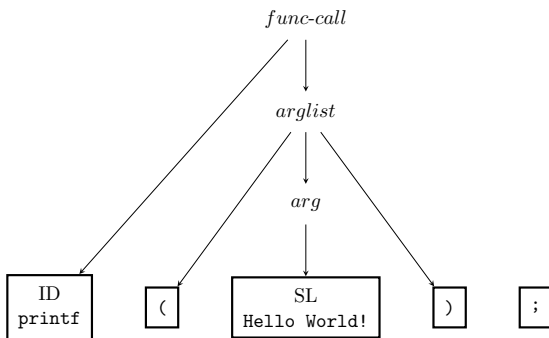
Parsing

Parse Tree



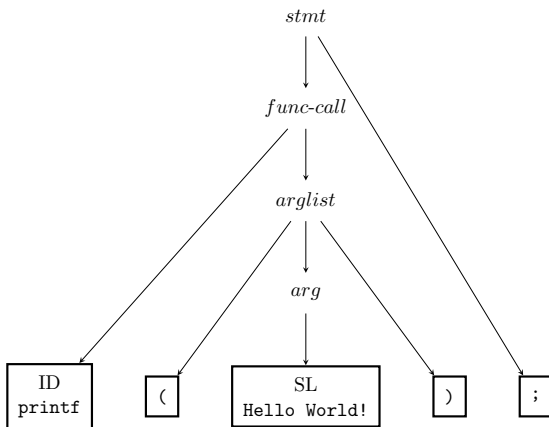
Parsing

Parse Tree



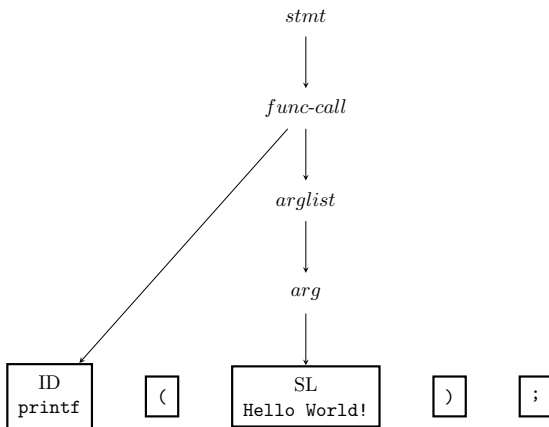
Parsing

Parse Tree



Parsing

Abstract Syntax Tree



Parsing

Parse Tree

- Grammar symbol \mapsto Nodes
- Starting symbol \mapsto Root node
- Non-terminals \mapsto internal nodes
- Terminals \mapsto leaves
- Productions \mapsto Edges

Parsing

Ambiguity

 $1 + 2 * 3$

$$E \rightarrow E + E$$

$$E \rightarrow E * E$$

$$E \rightarrow (E)$$

$$E \rightarrow 0 \mid 1 \dots \mid 9$$

Parsing

Ambiguity

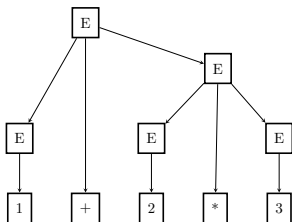
$1 + 2 * 3$

$$E \rightarrow E + E$$

$$E \rightarrow E * E$$

$$E \rightarrow (E)$$

$$E \rightarrow 0 \mid 1 \mid \dots \mid 9$$



Parsing

Ambiguity

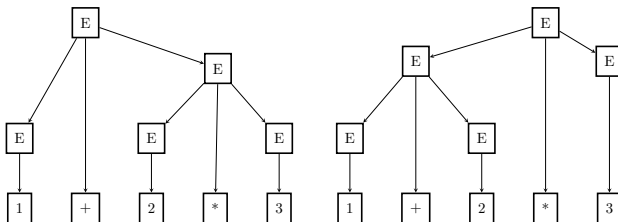
$1 + 2 * 3$

$$E \rightarrow E + E$$

$$E \rightarrow E * E$$

$$E \rightarrow (E)$$

$$E \rightarrow 0 \mid 1 \mid \dots \mid 9$$



Parsing

Ambiguity

- Language processors can't deal with ambiguity.
- Ambiguous grammars are common.
- Methods of dealing with ambiguity:
 - Fixing the grammar
 - Associativity
 - Operator precedence

Parsing

Handling Ambiguity – Fixing the grammar

Dangling else

$$\begin{array}{lcl} stmt & \rightarrow & \text{if } expr \text{ then } stmt \text{ else } stmt \\ & | & \text{if } expr \text{ then } stmt \end{array}$$

```
if C1 then
  if C2 then
    S1
  else
    S2
```

Parsing

Handling Ambiguity – Fixing the grammar

Dangling else

```

stmt    →  if expr then stmt else stmt
           |  if expr then stmt

```

```

stmt           →  matched_stmt
                  |  open_stmt
matched_stmt  →  if expr then matched_stmt else matched_stmt
                  |  other
open_stmt     →  if expr then stmt
                  |  if expr then matched_stmt else open_stmt

```


Parsing

Handling Ambiguity – Operator Precedence

$1 + 2 * 3$

$$\begin{array}{lll} E & \rightarrow & E + E \\ E & \rightarrow & E * E \\ E & \rightarrow & (E) \\ E & \rightarrow & 0 \mid 1 \dots \mid 9 \end{array}$$

Parsing

Handling Ambiguity – Operator Precedence

$1 + 2 * 3$

E	\rightarrow	$E + E$
E	\rightarrow	$E * E$
E	\rightarrow	(E)
E	\rightarrow	$0 \mid 1 \dots \mid 9$

* has higher precedence than +.

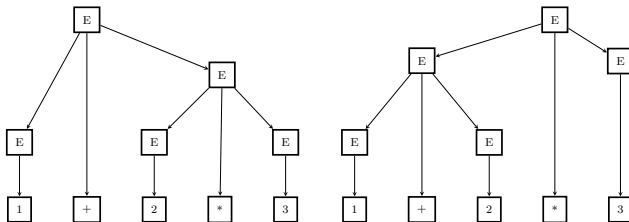
Parsing

Handling Ambiguity – Operator Precedence

$1 + 2 * 3$

E	\rightarrow	$E + E$
E	\rightarrow	$E * E$
E	\rightarrow	(E)
E	\rightarrow	$0 \mid 1 \dots \mid 9$

$*$ has higher precedence than $+$.



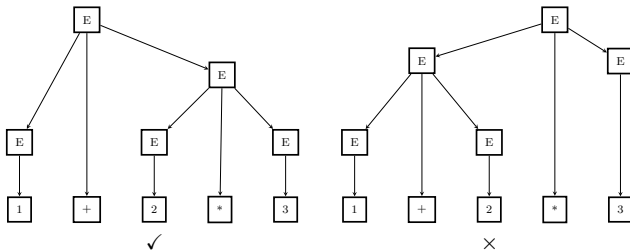
Parsing

Handling Ambiguity – Operator Precedence

$1 + 2 * 3$

E	\rightarrow	$E + E$
E	\rightarrow	$E * E$
E	\rightarrow	(E)
E	\rightarrow	$0 \mid 1 \dots \mid 9$

$*$ has higher precedence than $+$.



Parsing

Handling Ambiguity – Associativity

4 - 2 - 1

E	\rightarrow	$E + E$
E	\rightarrow	$E - E$
E	\rightarrow	(E)
E	\rightarrow	$0 \mid 1 \dots \mid 9$

Parsing

Handling Ambiguity – Associativity

4 - 2 - 1

$$\begin{array}{lll} E & \rightarrow & E + E \\ E & \rightarrow & E - E \\ E & \rightarrow & (E) \\ E & \rightarrow & 0 \mid 1 \dots \mid 9 \end{array}$$

- is left associative.

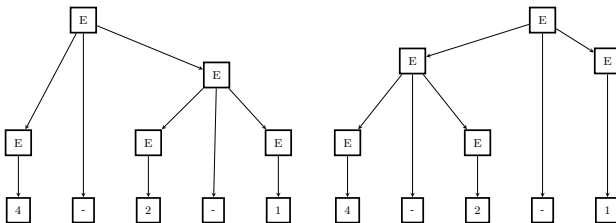
Parsing

Handling Ambiguity – Associativity

4 - 2 - 1

E	\rightarrow	$E + E$
E	\rightarrow	$E - E$
E	\rightarrow	(E)
E	\rightarrow	$0 \mid 1 \dots \mid 9$

- is left associative.



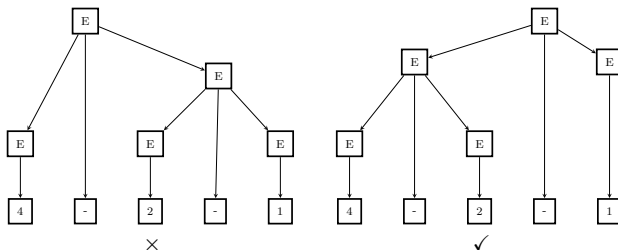
Parsing

Handling Ambiguity – Associativity

4 - 2 - 1

E	\rightarrow	$E + E$
E	\rightarrow	$E - E$
E	\rightarrow	(E)
E	\rightarrow	$0 \mid 1 \dots \mid 9$

- is left associative.



Parsing

Ambiguity

- Language processors can't deal with ambiguity.
- Ambiguous grammars are common.
- Methods of dealing with ambiguity:
 - Fixing the grammar
 - Associativity
 - Operator precedence
- Non-trivial
- Not covered

Parsing

Ambiguity

- Language processors can't deal with ambiguity.
- Ambiguous grammars are common.
- Methods of dealing with ambiguity:
 - Fixing the grammar
 - Associativity
 - Operator precedence
- Non-trivial
- Not covered



Next

Recursive Descent Parsing