

## Resolving syntactic ambiguity

---

$\text{Int Expr} ::= \text{Int Lit}$   
 $\quad \quad \quad | \text{Int Expr} + \text{Int Expr}$

$"3 + 4 + 5" \longrightarrow (3 + 4) + 5$   
 $\quad \quad \quad \text{or}$   
 $\quad \quad \quad 3 + (4 + 5)$

---

$\text{Int Expr} ::= \text{Int Lit}$   
 $\quad \quad \quad | \text{Int Expr} + \text{Int Expr}$   
 $\quad \quad \quad | \text{Int Expr} * \text{Int Expr}$

---

$"3 + 4 * 5" \longrightarrow (3 + 4) * 5$   
 $\quad \quad \quad \text{or}$   
 $\quad \quad \quad 3 + (4 * 5)$

---

In the case of  $+$  &  $-$

$\text{Int Expr} ::= \text{Int Expr} + \text{Int Lit}$

$\text{Int Expr} ::= \text{Int Expr} + \text{Int Lit}$

"3 + (4 + 5)"

|  $\text{Int Expr} + (\text{Int Expr})$

---

Back to + & \* :

Suggestion: Right side of \* is  
always an integer literal

~~$\text{Int Expr} ::= \text{Int Expr} * \text{Int Expr}$~~

$\text{Int Expr} ::= \text{Int Expr} * \text{Int Lit}$

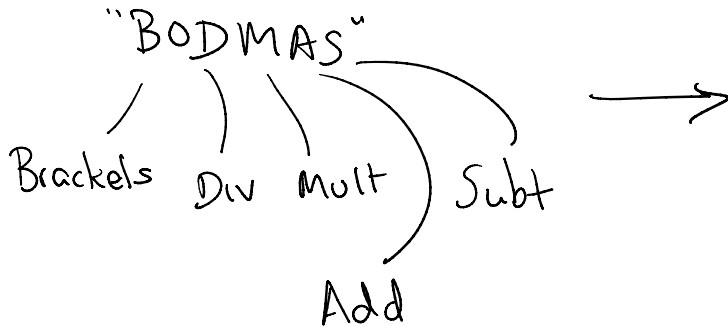
What about "3 \* (4 + 5)" ?

|  $\text{Int Expr} * (\text{Int Expr})$

---

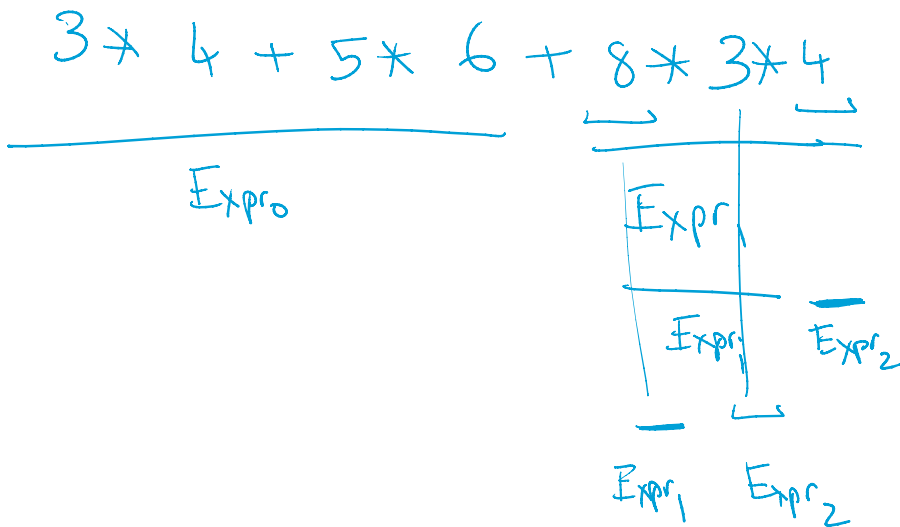
This suggests: Start with precedence

This suggests: Start with precedence rules.



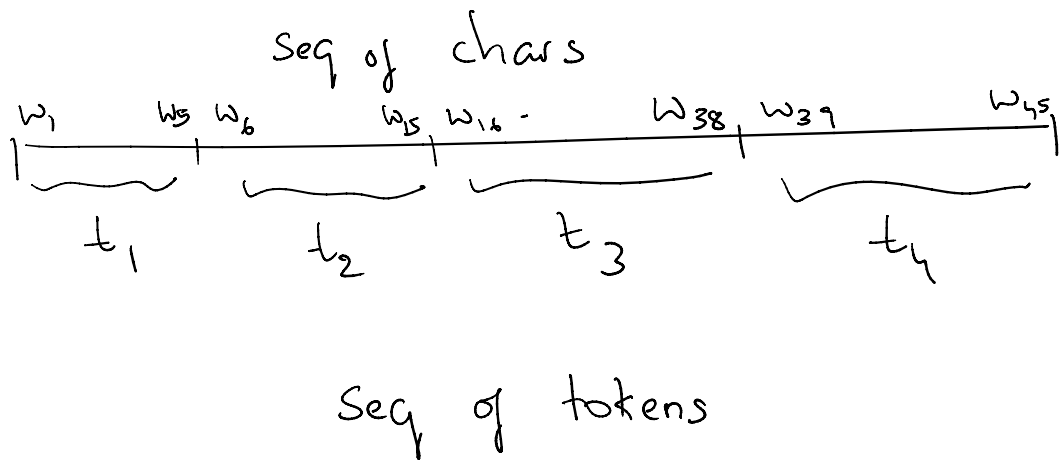
Brackets · Mult · Add

$Expr_0 ::= Expr_0 + \underline{Expr_1} (\tau)$   
 $Expr_1 ::= Expr_1 * \underline{Expr_2}$   
 $Expr_2 ::= \text{Int Lit}$   
 $Expr_2 ::= (Expr_0)$



Algorithmics

Lexical analyzer : seq of chars  $\longrightarrow$  seq of "tokens"



Lexical analyzer : set of pattern matching rules to identify tokens

---

Int Lit :  $[ '0' - '9' ]^+$

An integer literal is a sequence of ASCII digits

"9483 + 32"  $\longrightarrow$  Int Lit 9483 ; PLUS ; Int Lit 32

$\longrightarrow$  Int Lit 948 ; Int Lit 3 ; PLUS ; Int Lit 32

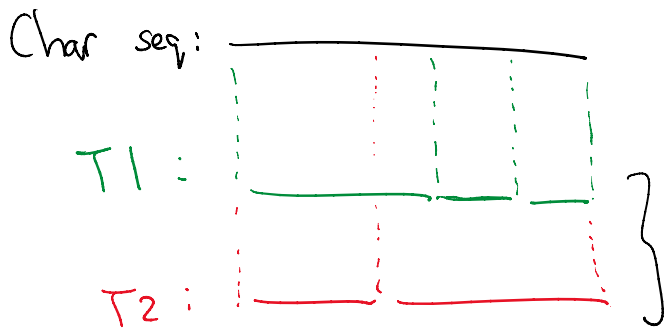
$\longrightarrow$  Int Lit 94 ; Int Lit 83 ; PLUS ; Int Lit 32

→ Int Lit 94; Int Lit 83; Plus; Int Lit 32

---

Ambiguity resolution using the Longest Match Principle.

Char seq:



Two alternative tokenizations  
Swallow as many characters  
as possible to emit first  
token.

"Longest match principle"

---

Pattern matching in the lexical analyzer

- Set of rules
- Each rule checks if a sequence of characters is a token or not.

characters is a token or not.

"Regular expressions": Way to describe patterns over strings.

Examples: Word =  $\left( \underbrace{('a' - 'z')}_{\textcircled{3} \text{ lowercase letter}} \cup \underbrace{('A' - 'Z')}_{\text{an } \textcircled{4} \text{ uppercase letter}} \right)^+$   
Non-empty sequence of chars  $\textcircled{1}$   
 $\textcircled{2}$  each of which is a

Email address:  $\left[ \begin{array}{l} ('a' - 'z') \cup ('A' - 'Z') \cup ('0' - '9') \\ \cdot '@' \cdot \left[ ('a' - 'z') \cup ('A' - 'Z') \cup ('0' - '9') \right]^+ \cdot \dots \end{array} \right]^+$

Word · '@' · Word · '.' · Word

---

Regular Expression,  $r ::=$  Char Range  $\left\{ \begin{array}{l} \text{Digit} \\ \text{Upper case letter} \\ \text{Vowel} \\ \vdots \end{array} \right.$   
 $| r_1 \cdot r_2 \dots$

|  $r_1 \cdot r_2$  \ vowel

String contains two parts

Left side matches  $r_1$

Right side matches  $r_2$

|  $r_1 \cup r_2$

String matches  $r_1$  or

String matches  $r_2$

|  $r^*$

Kleene \*



String can be decomposed

(possibly empty) into a sequence.

Each piece matches  $r$

|  $r^+$

String can be decomposed into

a sequence of one or more pieces, each of which matches  $r$ .

$$r^+ = r \cdot r^* = r^* \cdot r$$