



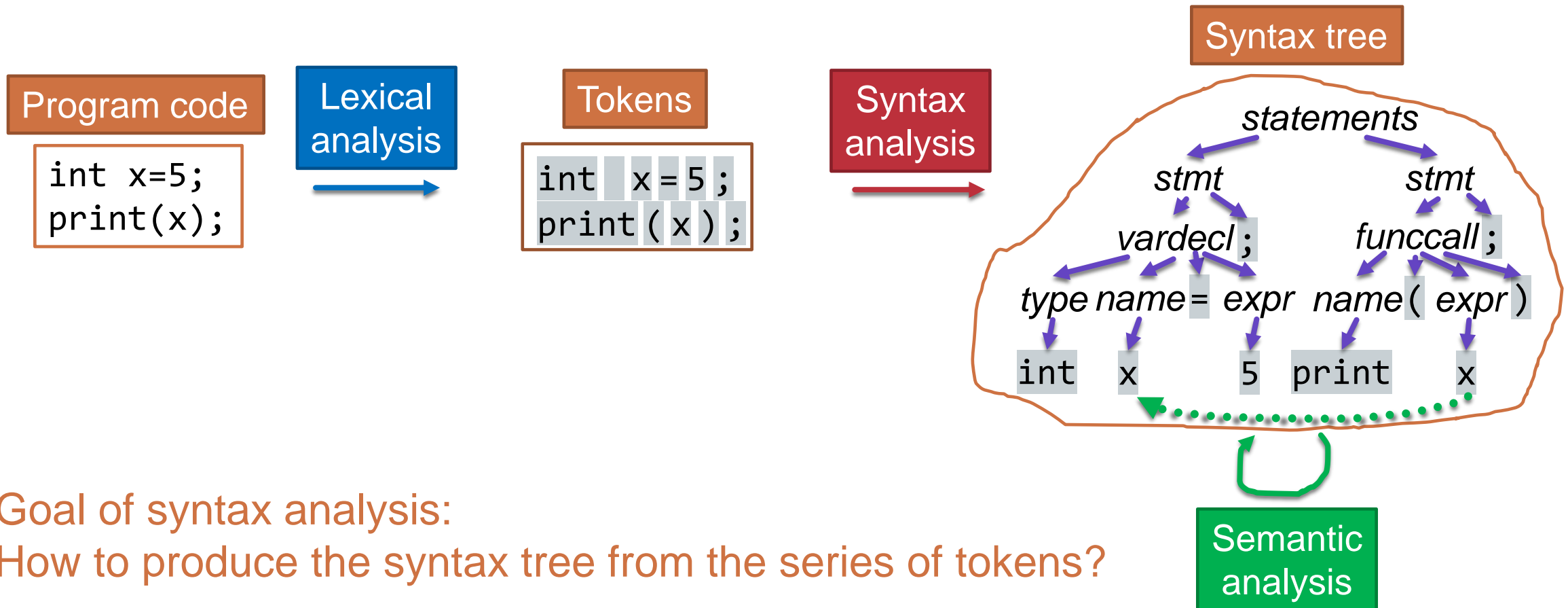
Model-based Software Development

Lecture 3

Syntax Analysis

Dr. Balázs Simon

Compiler front-end



This lecture: Syntax analysis

I. Context free (CF) grammars

II. Naïve methods (BFS, DFS)

III. Left analysis: $LL(k)$, $LL(*)$

IV. Right analysis: $LR(k)$, LALR

V. Error handling



Context free (CF) grammars

Non-terminal

$A \rightarrow \alpha \mid \beta \mid \dots$

A series of terminal and non-terminal symbols

Notation:

- Uppercase letters: non-terminal symbols
- Lowercase letters, other characters: terminals
- Greek letters: any terminal - non-terminal series

Parse tree:

- Root: start symbol
- Internal node: non-terminal symbol
- Leaf: terminal symbol (token)

CF grammar example

Production rules:

$E \rightarrow x|y|z|E \text{ Op} E|(E)$
 $Op \rightarrow +|-|*|/$

Program code:

$x+y*z$

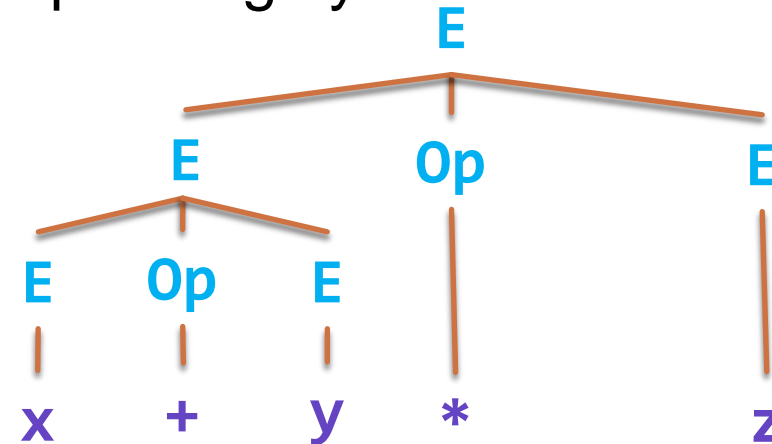
Lexical analysis (tokens/terminals):

$x + y * z$

One possible derivation:

$E \xrightarrow{4} E \text{ Op} E \xrightarrow{4} E \text{ Op} E \text{ Op} E \rightarrow$
 $\xrightarrow{2} E \text{ Op} y \text{ Op} E \xrightarrow{6} E + y \text{ Op} E \rightarrow$
 $\xrightarrow{1} x + y \text{ Op} E \xrightarrow{8} x + y * E \rightarrow$
 $\xrightarrow{3} x + y * z$

Corresponding syntax tree:



meaning:
 $(x+y)*z$

Conclusions

- There are many possible derivations
 - > ambiguous grammar
- The meaning is determined by the structure of the syntax tree
 - > depends on the order of the applied rules
 - > precedence of operators
- Unambiguous grammar:
 - > syntax tree is always unambiguous

Unambiguous grammar

Production rules:

$E \rightarrow E+T \mid E-T \mid T$
 $T \rightarrow T * F \mid T / F \mid F$
 $F \rightarrow x \mid y \mid z \mid (E)$

Program code:

$x+y*z$

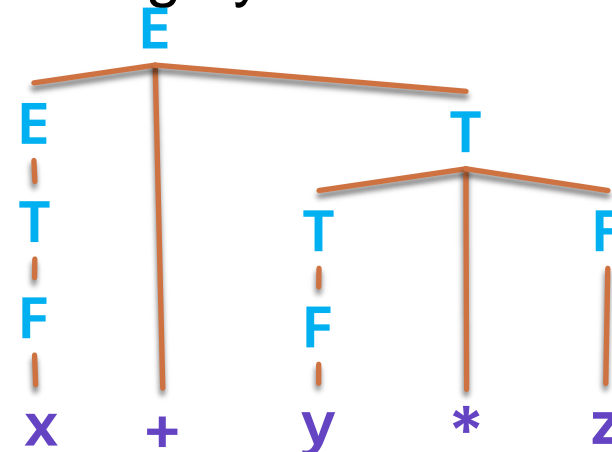
Lexical analysis (tokens/terminals):

$x + y * z$

One possible derivation:

$E \xrightarrow{1} E + T \xrightarrow{4} E + T * F \xrightarrow{6} E + F * F \xrightarrow{8} E + y * F \xrightarrow{3} T + y * F \xrightarrow{6} F + y * F \xrightarrow{7} x + y * F \xrightarrow{9} x + y * z$

Corresponding syntax tree:



meaning:
 $x+(y*z)$

Questions

- How can we make the analysis automated?
 - > there are many different algorithms
- How can we detect if a grammar is unambiguous?
 - > depends on the analyzer: either when we construct the analyzer or when we run the analyzer
- Can we make a grammar unambiguous?
 - > usually yes: by rewriting the production rules
 - > caution: the structure may change

This lecture: Syntax analysis

I. Context free (CF) grammars

II. Naïve methods (BFS, DFS)

III. Left analysis: $LL(k)$, $LL(*)$

IV. Right analysis: $LR(k)$, LALR

V. Error handling

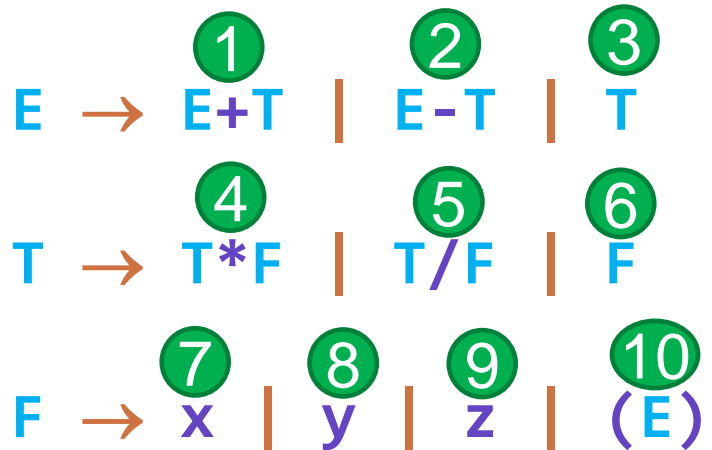


Naïve methods

- Idea:
 - > graph search algorithms (BFS, DFS)
 - > backtracking
- Start with the start symbol
- Nodes:
 - > symbol series derivable from the start symbol in one or more steps
- Edges:
 - > between α and β if β can be directly (in a single step) derived from α

Example

Production rules:



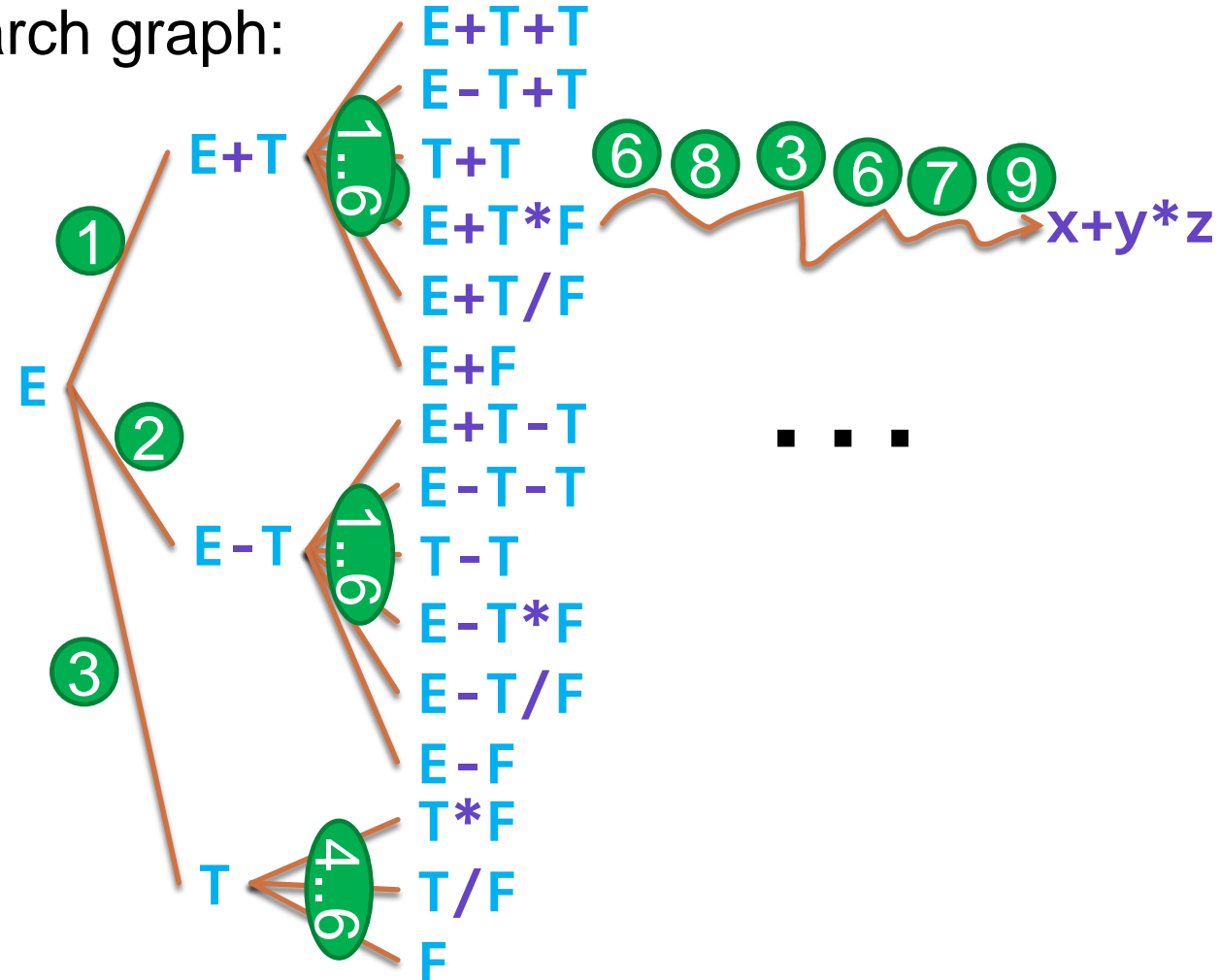
Program code:

$x+y*z$

Lexical analysis (tokens/terminals):

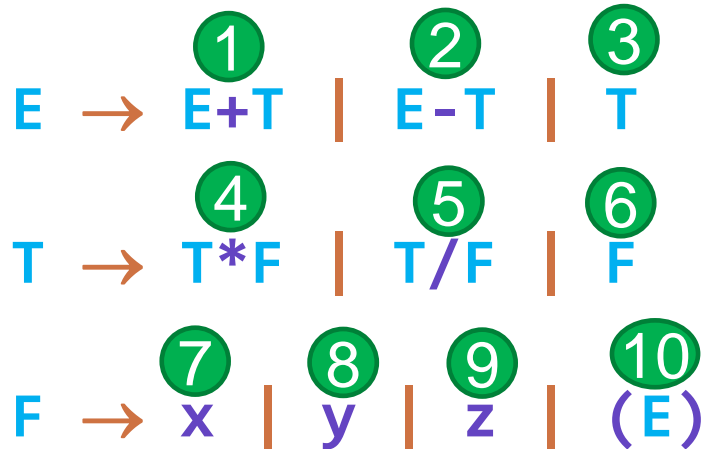
$x + y * z$

Search graph:



Example: Breadth-First Search (BFS)

Production rules:



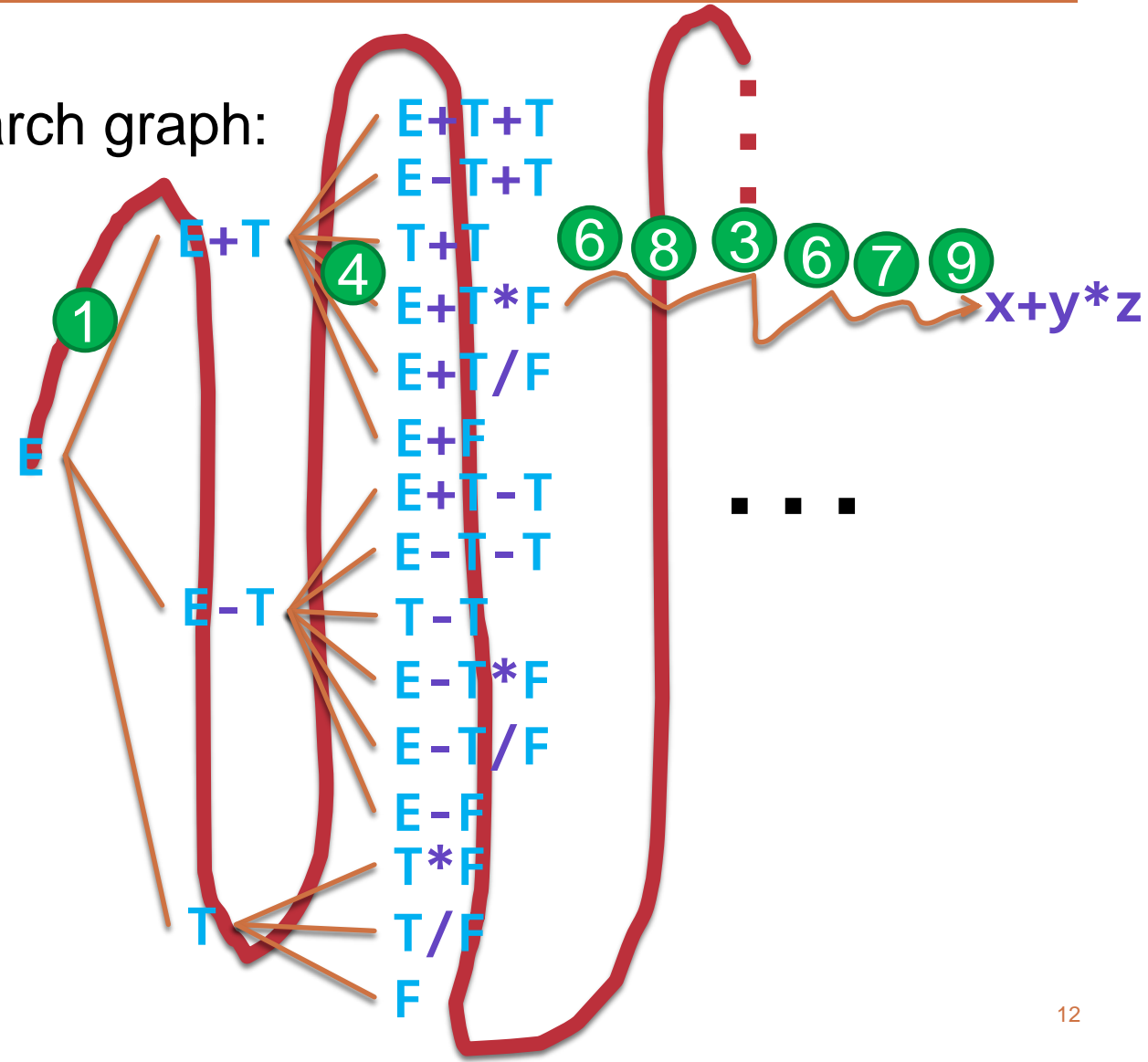
Program code:

$x+y*z$

Lexical analysis (tokens/terminals):

$x + y * z$

Search graph:



Example: Depth-First Search (DFS)

Production rules:

The diagram illustrates the derivation of the expression $(E + T) * (E - T) / (F * x) / (y / z)$ using the grammar rules:

- $E \rightarrow E + T \mid E - T \mid T$
- $T \rightarrow T * F \mid T / F \mid F$
- $F \rightarrow x \mid y \mid z \mid (E)$

The derivation steps are numbered 1 through 10:

- $E \rightarrow E + T$
- $E \rightarrow E - T$
- $E \rightarrow T$
- $T \rightarrow T * F$
- $T \rightarrow T / F$
- $T \rightarrow F$
- $F \rightarrow x$
- $F \rightarrow y$
- $F \rightarrow z$
- $F \rightarrow (E)$

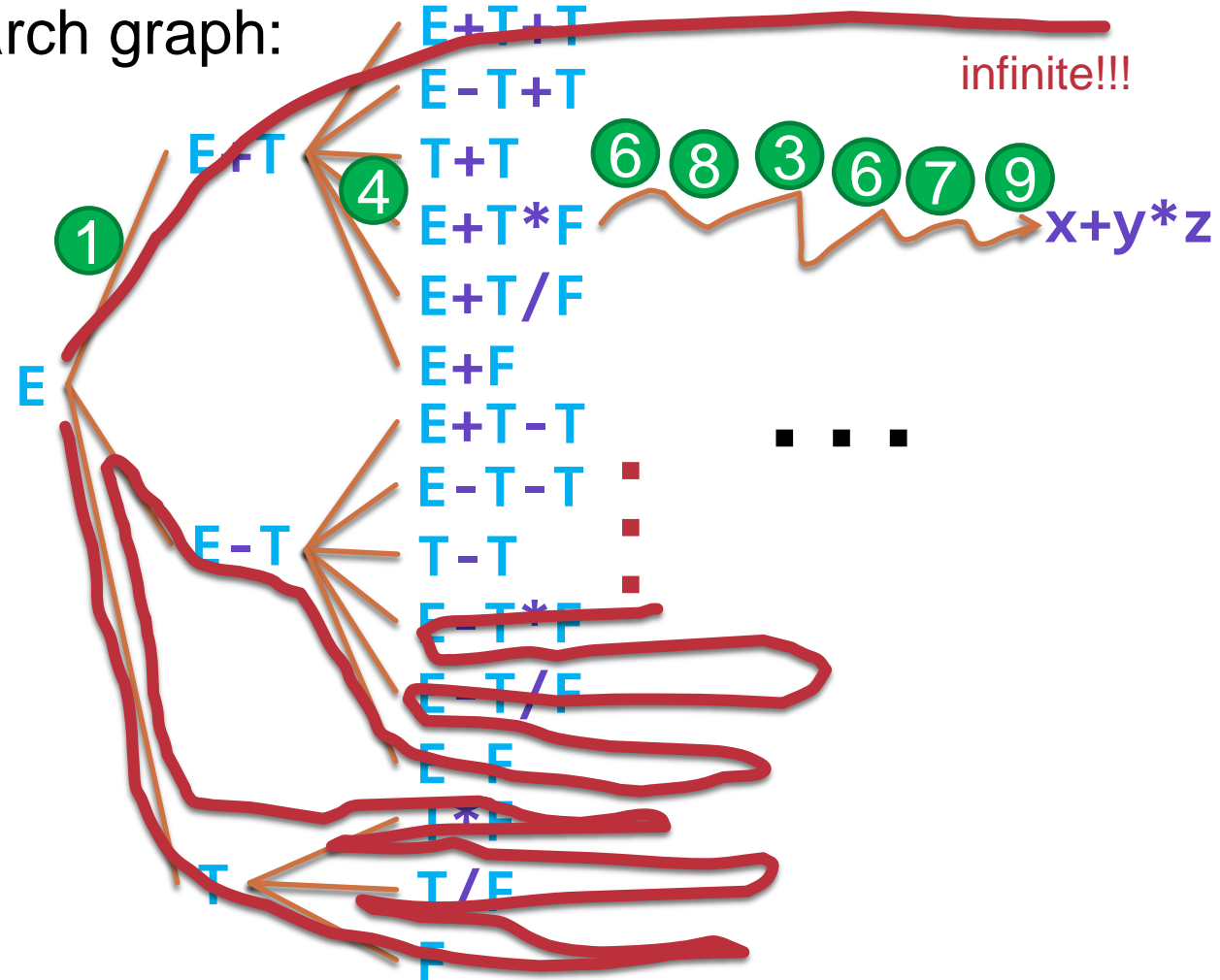
Program code:

$$x + y * z$$

Lexical analysis (tokens/terminals):

x + y * z

Search graph:



Conclusions

- Search graph \neq syntax tree
 - > syntax tree can be constructed based on the edges of the search graph
- Problem: large state space, many branches
 - > BFS: exponential memory, exponential time
 - > DFS: smaller memory, but can be infinite in time
 - problem: left recursion ($E \rightarrow E+T \rightarrow E+E+T \rightarrow E+E+E+T \rightarrow \dots$)
 - solution: later...
- State space must somehow be reduced
- Rarely used in practice

This lecture: Syntax analysis

I. Context free (CF) grammars

II. Naïve methods (BFS, DFS)

III. Left analysis: $LL(k)$, $LL(*)$

IV. Right analysis: $LR(k)$, LALR

V. Error handling



Left analysis: LL(k)

- Prediction instead of backtracking
- Idea:
 - > L: read from left to right
 - > L: always rewrite the leftmost non-terminal (leftmost derivation)
 - > decision between alternatives: based on the lookahead of **k** terminals (token)
- Errors recognized automatically: comparing the produced and predicted terminals
 - > extraneous and missing tokens
- Problem: left recursion

Left analysis example

Production rules:

$S \rightarrow Sa \mid b$

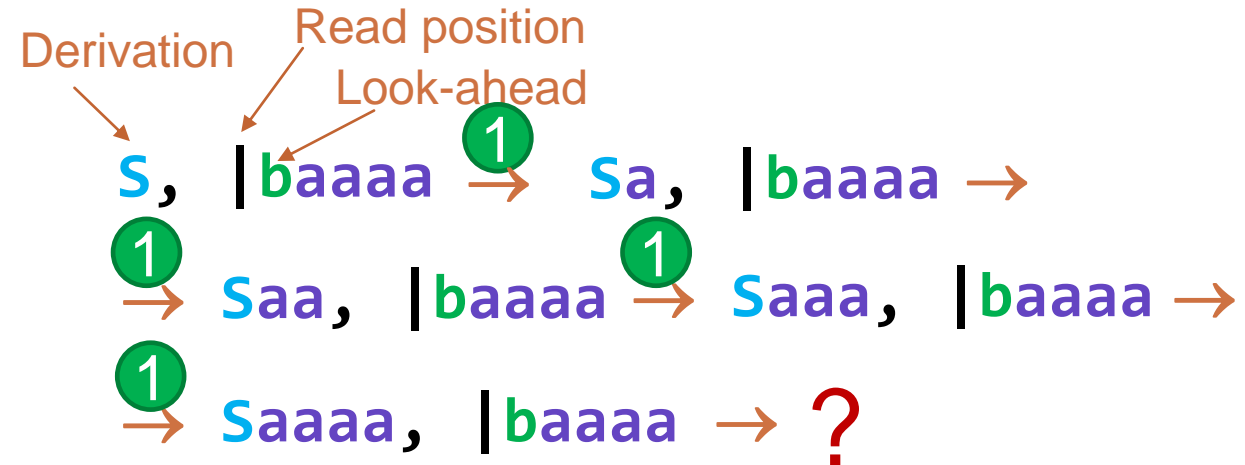
Program code:

baaaa

Lexical analysis (tokens/terminals):

b a a a a

LL(1) analysis:



We should stop here using rule 2, however, we can't possibly know this by the look-ahead!

Resolving left recursion: by transforming the grammar

Production rules:

$S \rightarrow b\hat{S}$ (1)
 $\hat{S} \rightarrow a\hat{S} \mid \epsilon$ (2) (3)

Program code:

`baaaa`

Lexical analysis (tokens/terminals):

`b a a a a`

LL(1) analysis:

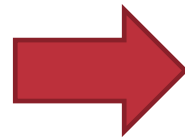
Derivation Read position Look-ahead

$S, \mid baaaa \xrightarrow{(1)} b\hat{S}, b \mid aaaa \rightarrow$
 $\xrightarrow{(2)} ba\hat{S}, ba \mid aaa \xrightarrow{(2)} baa\hat{S}, baa \mid aa \rightarrow$
 $\xrightarrow{(2)} baaa\hat{S}, baaa \mid a \rightarrow$
 $\xrightarrow{(2)} baaaa\hat{S}, baaaa \mid \epsilon \rightarrow$
 $\xrightarrow{(3)} baaaa, baaaa \mid \epsilon$

Left recursion

- Solution 1: transforming the grammar
 - > disadvantage: the structure of the syntax tree changes
- Solution 2: priority between different rules
 - > ANTLR4: direct left recursion is allowed (but indirect is not!)

Solution 1: transforming the grammar

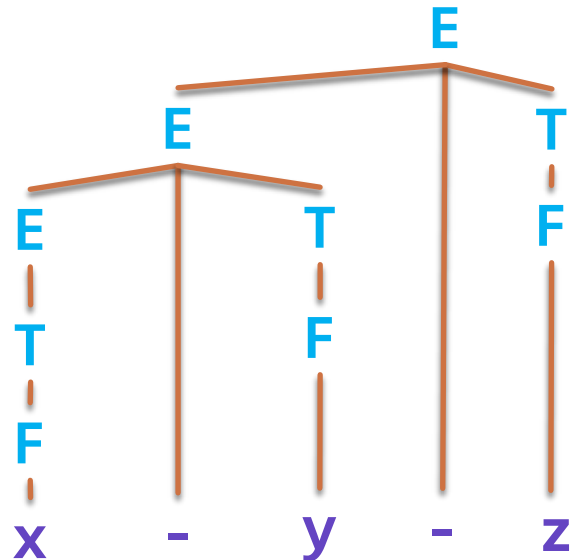
$$\begin{aligned} E &\rightarrow E+T \mid E-T \mid T \\ T &\rightarrow T*F \mid T/F \mid F \\ F &\rightarrow x \mid y \mid z \mid (E) \end{aligned}$$


Question: why is this a bad solution?

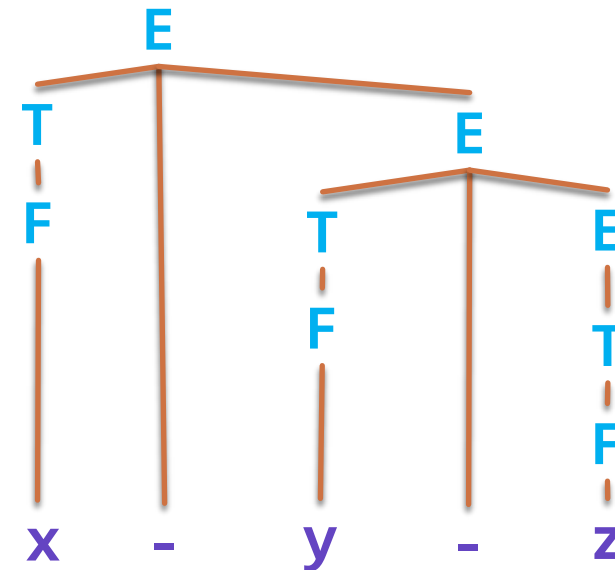
$$\begin{aligned} E &\rightarrow T+E \mid T-E \mid T \\ T &\rightarrow F*T \mid F/T \mid F \\ F &\rightarrow x \mid y \mid z \mid (E) \end{aligned}$$

Answer: associativity of the operators is important!

meaning:
 $(x-y)-z$



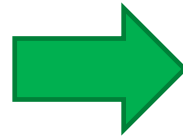
vs.



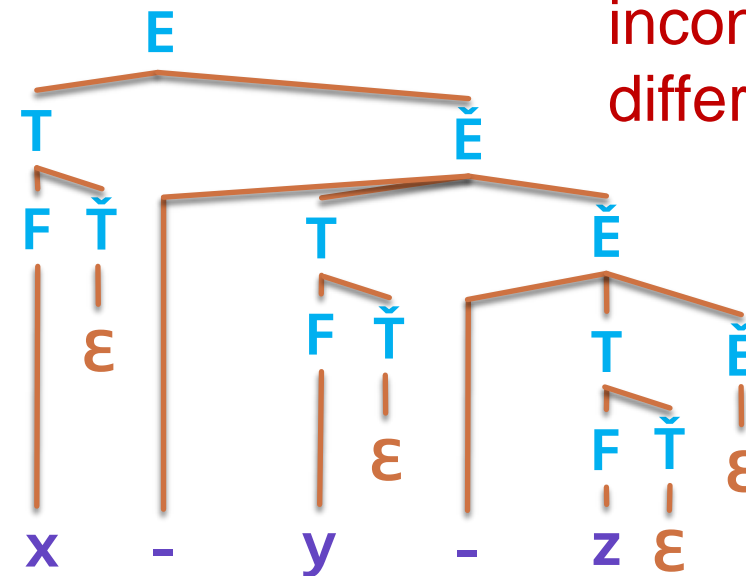
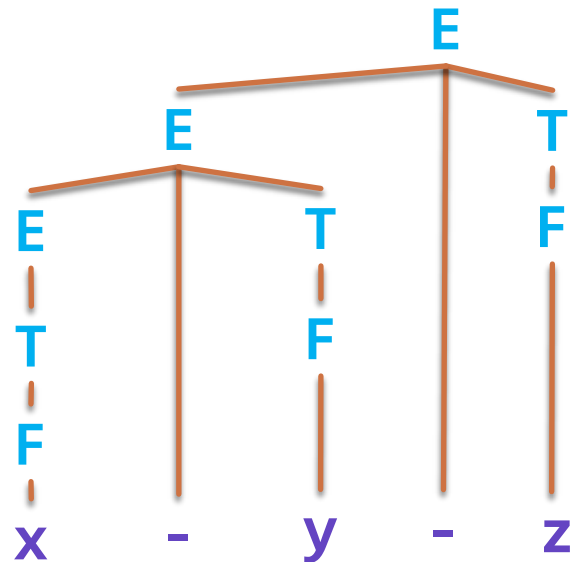
meaning:
 $x-(y-z)$

Solution 1: correct transformation of the grammar

$E \rightarrow E+T \mid E-T \mid T$
 $T \rightarrow T * F \mid T / F \mid F$
 $F \rightarrow x \mid y \mid z \mid (E)$



$E \rightarrow T\check{E}$
 $\check{E} \rightarrow +T\check{E} \mid -T\check{E} \mid \varepsilon$
 $T \rightarrow F\check{T}$
 $\check{T} \rightarrow *F\check{T} \mid /F\check{T} \mid \varepsilon$
 $F \rightarrow x \mid y \mid z \mid (E)$



inconvenience:
different structure

LL(1) analysis

Production rules:

$$\begin{aligned} E &\rightarrow T\check{E} \\ \check{E} &\rightarrow +T\check{E} \mid -T\check{E} \mid \epsilon \\ T &\rightarrow F\check{T} \\ \check{T} &\rightarrow *F\check{T} \mid /F\check{T} \mid \epsilon \\ F &\rightarrow x \mid y \mid z \mid (E) \end{aligned}$$

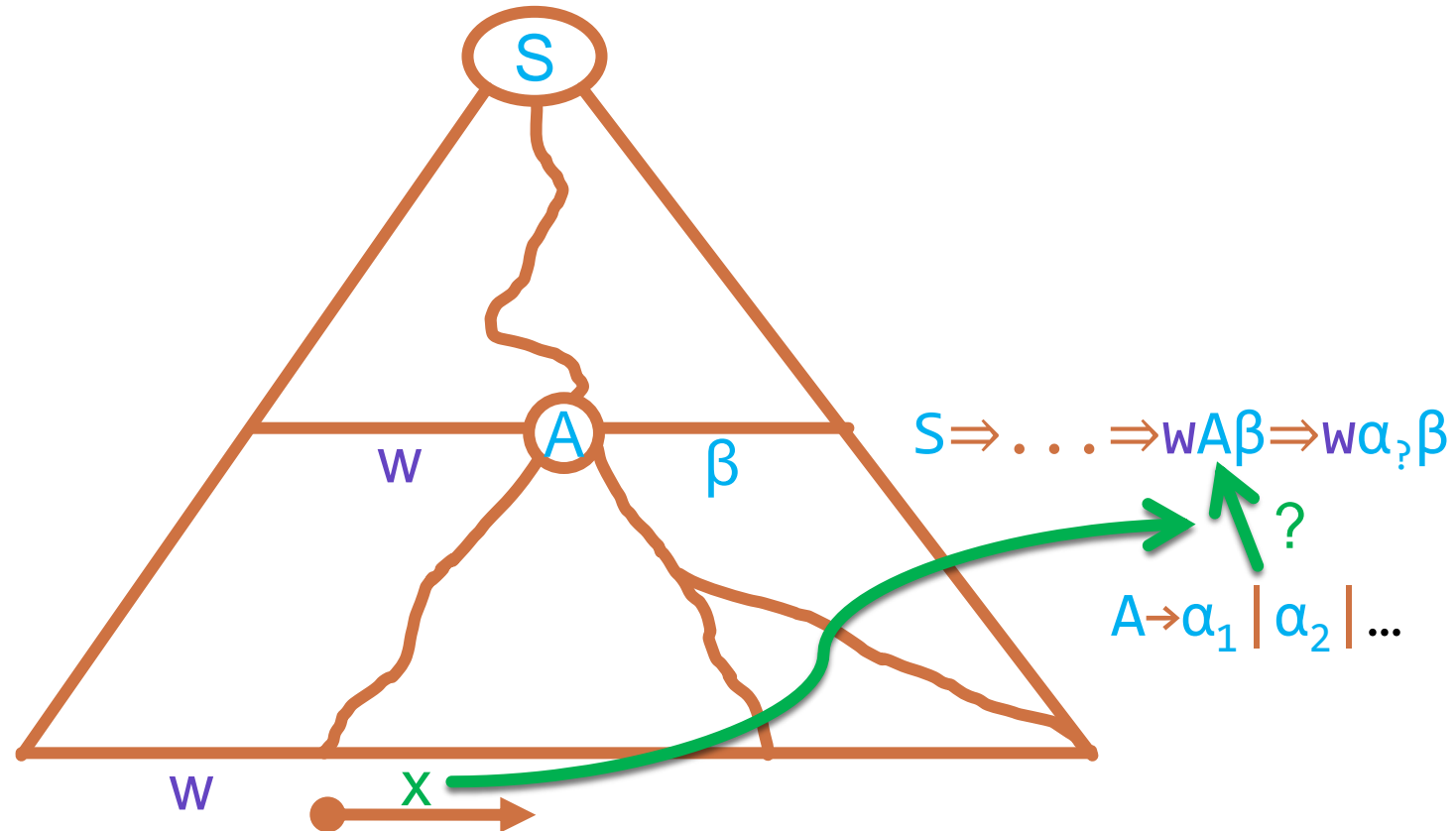
Program code: $x+y*z$

Lexical analysis (tokens): $x \ + \ y \ * \ z$

LL(1) analysis:

$$\begin{aligned} E, \mid x+y*z &\xrightarrow{1} T\check{E}, \mid x+y*z \rightarrow \\ &\xrightarrow{5} F\check{T}\check{E}, \mid x+y*z \xrightarrow{9} x\check{T}\check{E}, x \mid +y*z \rightarrow \\ &\xrightarrow{8} x\check{E}, x \mid +y*z \xrightarrow{2} x+T\check{E}, x+ \mid y*z \rightarrow \\ &\xrightarrow{5} x+F\check{T}\check{E}, x+ \mid y*z \xrightarrow{10} x+y\check{T}\check{E}, x+y \mid *z \rightarrow \\ &\xrightarrow{6} x+y*F\check{T}\check{E}, x+y* \mid z \rightarrow \\ &\xrightarrow{11} x+y*z\check{T}\check{E}, x+y*z \mid \epsilon \rightarrow \\ &\xrightarrow{8} x+y*z\check{E}, x+y*z \mid \epsilon \rightarrow \\ &\xrightarrow{4} x+y*z, x+y*z \mid \epsilon \end{aligned}$$

LL(k) overview



We have to take into account the characters generated by A and β , when we compute x of length k .

Advantages of LL(k)

- Simple
 - > can be programmed manually (e.g., Roslyn – C# compiler): Recursive Descent Parser
 - > analyzer can be easily generated: based on a table
- Fast
- Small memory footprint
- Easy to detect and signal errors
 - > smaller errors can be fixed, and the analysis can be continued
 - > e.g., skipping an unexpected token, inserting a missing token

Example (1/3): LL(k) analyzer written manually

Production rules:

$E \rightarrow T\check{E}$
 $\check{E} \rightarrow +T\check{E} \mid -T\check{E} \mid \epsilon$
 $T \rightarrow F\check{T}$
 $\check{T} \rightarrow *F\check{T} \mid /F\check{T} \mid \epsilon$
 $F \rightarrow x \mid y \mid z \mid (E)$

C# code:



```
E ParseE()  
{  
    var t = ParseT();  
    var ehat = ParseEHat();  
    return new E(t, ehat);  
}
```

Example (2/3): LL(k) analyzer written manually

Production rules:

$E \rightarrow T\check{E}$
 $\check{E} \rightarrow +T\check{E} \mid -T\check{E} \mid \epsilon$
 $T \rightarrow F\check{T}$
 $\check{T} \rightarrow *F\check{T} \mid /F\check{T} \mid \epsilon$
 $F \rightarrow x \mid y \mid z \mid (E)$

C# code:

```
EHat? ParseEHat() {  
    var la1 = LA(1);  Look-ahead  
    switch (la1)  
    {  
        case "+":  
            Match("+");  Step forward  
            var t1 = ParseT();  
            var ehat1 = ParseEHat();  
            return new EHat(t1, ehat1);  
        case "-":  
            Match("-");  
            var t2 = ParseT();  
            var ehat2 = ParseEHat();  
            return new EHat(t2, ehat2);  
        default:  
            return null;  
    }  
}
```

Example (3/3): LL(k) analyzer written manually

C# code:

Production rules:

$E \rightarrow T\check{E}$

$\check{E} \rightarrow +T\check{E} \mid -T\check{E} \mid \epsilon$

$T \rightarrow F\check{T}$

$\check{T} \rightarrow *F\check{T} \mid /F\check{T} \mid \epsilon$

$F \rightarrow x \mid y \mid z \mid (E)$

```
F? ParseF()
{
    var la1 = LA(1);
    switch (la1)
    {
        case "x": Match("x"); return new F("x");
        case "y": Match("y"); return new F("y");
        case "z": Match("z"); return new F("z");
        case "(":
            Match("(");
            var e = ParseE();
            Match(")"); ← Error, if something else comes!
            return new F(e);
        default:
            Unexpected(la1); ← Error if an unexpected token comes!
            return null;
    }
}
```

LL(*) analysis

- Idea:
 - > instead of a fixed **k** lookahead, predict the winning alternative using a state machine
 - > the state machine can read ahead any number of tokens
- Advantage: stronger than LL(k)
- In practice:
 - > ANTLR3 - LL(*): <http://wwwantlr.org/papers/LL-star-PLDI11.pdf>
 - > ANTLR4 - ALL(*) (Adaptive LL): <http://wwwantlr.org/papers/allstar-techreport.pdf>
 - even allows direct left recursion with precedence between the alternatives

This lecture: Syntax analysis

I. Context free (CF) grammars

II. Naïve methods (BFS, DFS)

III. Left analysis: $LL(k)$, $LL(*)$

IV. Right analysis: $LR(k)$, LALR

V. Error handling



Right analysis: LR(k)

- Prediction instead of backtracking
- Idea:
 - > L: read from left to right
 - > R: always the rightmost part of the derivation is changed (top of the stack)
 - > decision between shift-reduce: based on the lookahead of **k** terminals (token)
 - shift: put the next terminal onto the top of the stack
 - reduce: replace some symbols on the top of the stack, a potential right side of a rule, with the left side of the rule
- Automatic error recognition: the first error can be recognized precisely
 - > backtracking is necessary to continue the analysis
- Left recursion is not a problem

LR(1) analysis

Production rules:

$E \rightarrow E+T \mid E-T \mid T$
 $T \rightarrow T*F \mid T/F \mid F$
 $F \rightarrow x \mid y \mid z \mid (E)$

Program code:

$x+y*z$

Lexical analysis (tokens):

$x \ + \ y \ * \ z$

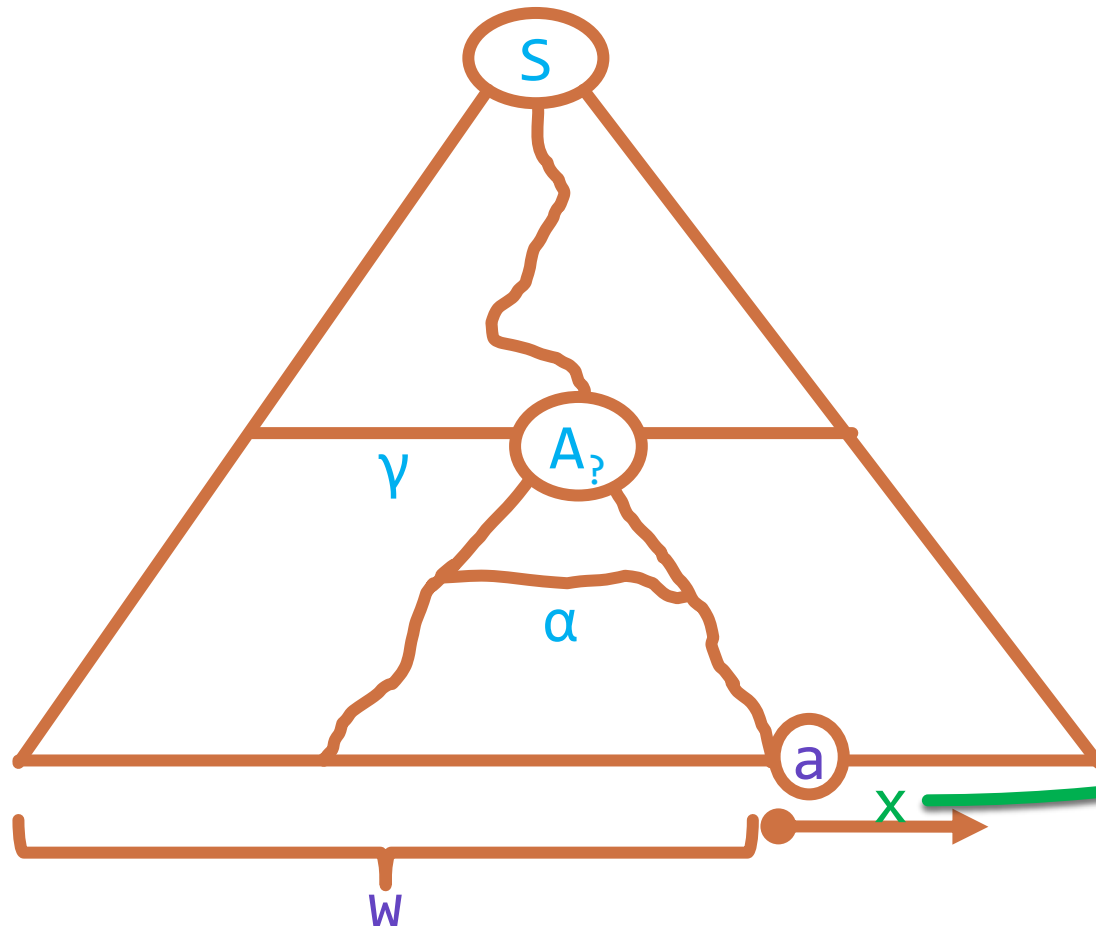
LR(1) analysis:

Stack: ϵ , Read position: $x+y*z$, Look-ahead: S

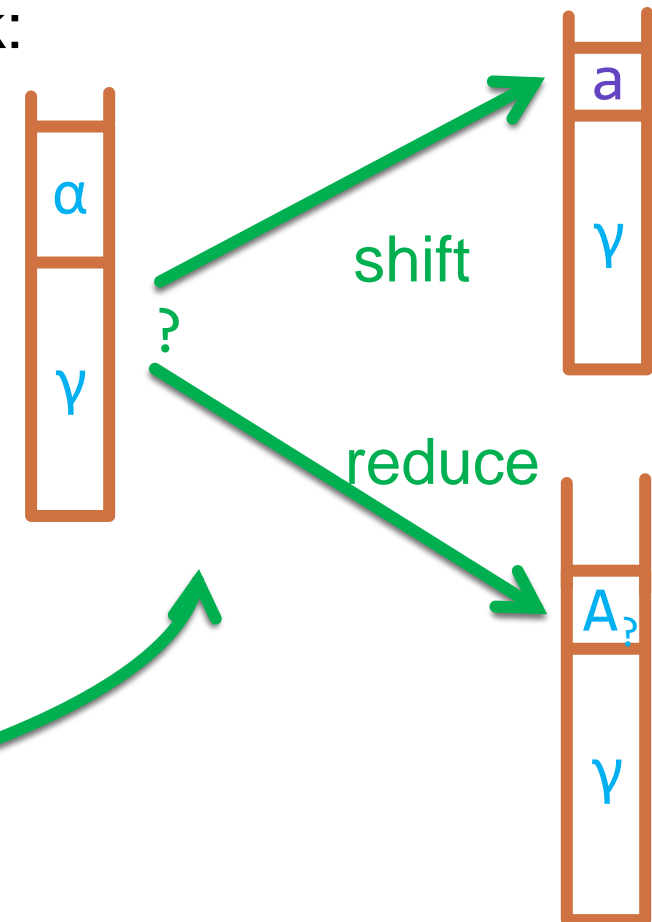
$\epsilon, \mid x+y*z \xrightarrow{S} x, x \mid +y*z \xrightarrow{7} F, x \mid +y*z \rightarrow$
 $\xrightarrow{6} T, x \mid +y*z \xrightarrow{3} E, x \mid +y*z \xrightarrow{S} E+, x+ \mid y*z \rightarrow$
 $\xrightarrow{S} E+y, x+y \mid *z \xrightarrow{8} E+F, x+y \mid *z \rightarrow$
 $\xrightarrow{6} E+T, x+y \mid *z \xrightarrow{S} E+T*, x+y* \mid z \rightarrow$
 $\xrightarrow{S} E+T*z, x+y*z \mid \epsilon \xrightarrow{9} E+T*F, x+y*z \mid \epsilon \rightarrow$
 $\xrightarrow{4} E+T, x+y*z \mid \epsilon \xrightarrow{1} E, x+y*z \mid \epsilon$

LR(k) overview

$S \Rightarrow \dots \Rightarrow \gamma A_? x \Rightarrow \gamma \alpha x \Rightarrow \dots \Rightarrow w x$



Stack:



Shift-Reduce conflict: Dangling Else

What is the problem with the following grammar?

Statement \rightarrow IfStatement | Expression

Expression \rightarrow IDENTIFIER

IfStatement \rightarrow

① IF Expression THEN Statement |

② IF Expression THEN Statement ELSE Statement

Possible solutions:

1. Changing the grammar
2. Prefer either shift or reduce
3. Precedence between rules

How does it analyze the following?

IF Cond1 THEN

IF Cond2 THEN DoSomething1

② ELSE DoSomething2

Shift?

How does it analyze the following?

IF Cond1 THEN

IF Cond2 THEN DoSomething1

ELSE DoSomething2

Reduce?

①

(Problems like „dangling else” are an issue in general, not just for LR analysis.)

Resolution of the Shift-Reduce conflict: changing the grammar

Statement \rightarrow IfStatement | Expression

Expression \rightarrow IDENTIFIER

IfStatement \rightarrow

IF Expression THEN Statement |

IF Expression THEN Statement ELSE Statement



Statement \rightarrow IfStatement | Expression

Expression \rightarrow IDENTIFIER

IfStatement \rightarrow

IF Expression THEN Statement |

IF Expression THEN IfThenElseStatement ELSE Statement

IfThenElseStatement \rightarrow

IF Expression THEN IfThenElseStatement ELSE IfThenElseStatement
Expression

ELSE is always bound to
the nearest IF

LR(k) analysis

- Advantages:
 - > Stronger than LL(k)
 - > Deterministic
 - > Fast: linear in time
 - > Left recursion is not a problem
 - > Can be generated automatically: based on a table
- Disadvantages:
 - > The table is very large, even for LR(1)
 - > Shift-reduce conflicts must be resolved

LALR analysis

- LALR = Look-Ahead LR parser
- Simplified right analyzer
 - > LR(1) analyzer: stronger than LR(0)
 - > LR(0) state space: smaller state space than LR(1)
- Can be generated automatically
 - > E.g., Yacc, Bison

This lecture: Syntax analysis

I. Context free (CF) grammars

II. Naïve methods (BFS, DFS)

III. Left analysis: $LL(k)$, $LL(*)$

IV. Right analysis: $LR(k)$, LALR

V. Error handling



Error handling

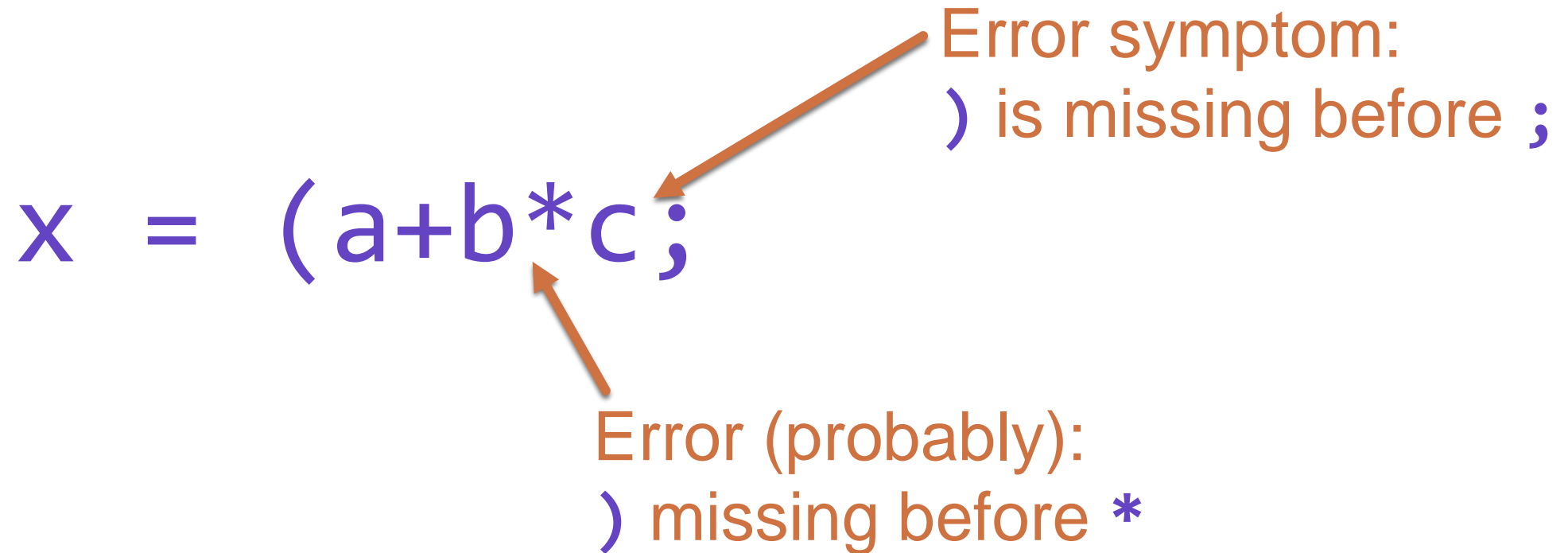
- Goal of syntax analysis:
 - > production of the syntax tree or finding as many errors as possible
- On error:
 - > error message, return to consistent state, continue analysis
- Report:
 - > position (file name, line, character)
 - > severity (error, warning, info, etc.)
 - > error message

Error position vs. error symptom

$x = (a+b*c;$

Error symptom:
) is missing before ;

Error (probably):
) missing before *



Position of the symptom is not necessarily the same as the position of the error!

Errors in various phases of compilation

- Lexical error:
 - > e.g., invalid character, missing end of a string or comment, premature end of file
- Syntax error:
 - > code is not valid according to the grammar
 - > correction: with the fewest operations (inserting, deleting or replacing a token)
- Semantic error:
 - > error in the static semantics (e.g., type error)
 - > some of these can only be detected during optimization (pl. invalid indexing)
- Lack of resources:
 - > e.g., insufficient memory
 - > can occur in any phase



Thank you!