

# Compiler principles

---

Compiler

Jakub Yaghob





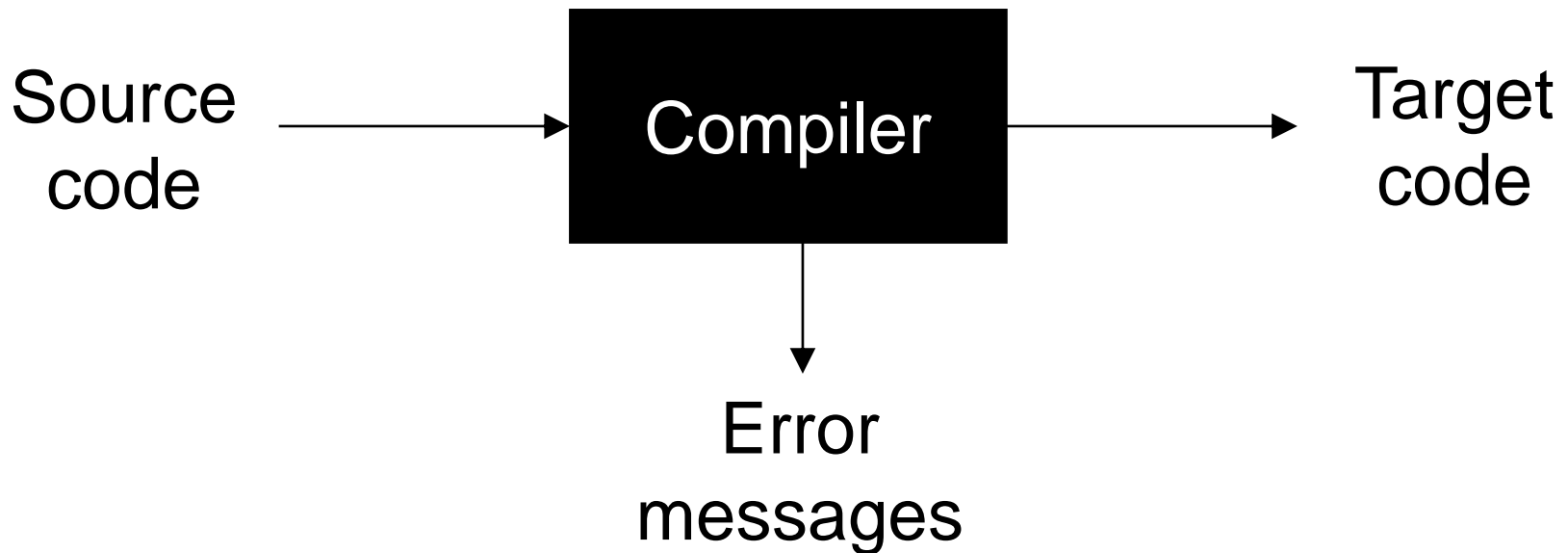
# Literature and slides

- The Dragon Book
  - Aho, Sethi, Ullman: Compilers - Principles, Techniques and Tools, Addison-Wesley 1986
  - Aho, Lam, Sethi, Ullman: Compilers - Principles, Techniques and Tools (2<sup>nd</sup> edition), Addison-Wesley 2006
- Advanced compiler techniques
  - Muchnick S.S.: Advanced compiler design and implementation, Morgan Kaufman Publishers 1997
- Slides
  - <http://www.ksi.mff.cuni.cz/lectures/NSWI098/html/index.html>



# What is a compiler?

- Naïve concept
  - A black-box compiling a source code to a target code



# What is a compiler? More formally



- Let's have an input language  $L_{in}$  generated by a grammar  $G_{in}$
- Let's have an output language  $L_{out}$  generated by a grammar  $G_{out}$  or accepted by an automaton  $A_{out}$
- The compiler is a mapping  $L_{in} \rightarrow L_{out}$ , where  $\forall w_{in} \in L_{in} \exists w_{out} \in L_{out}$ . The mapping does not exist for  $w_{in} \notin L_{in}$

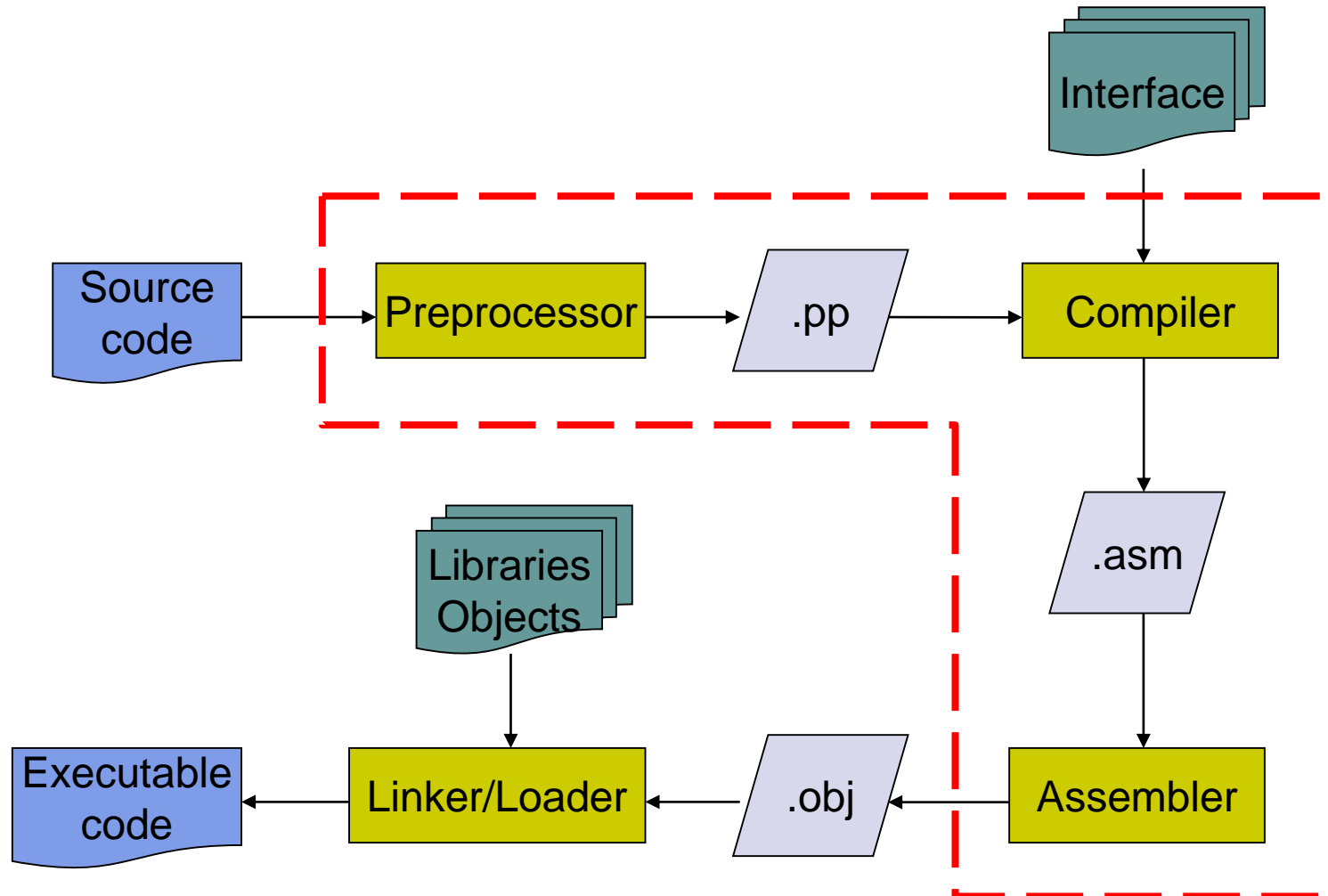


# Use cases

- Structured (e.g. MPS from JetBrains) or syntax-highlighting editor
- Pretty-printer
- Static program checker
  - LINT
- Interpreters
- Modelling languages compiler
  - Verilog, VHDL
- Query languages
  - SQL



# Program translation

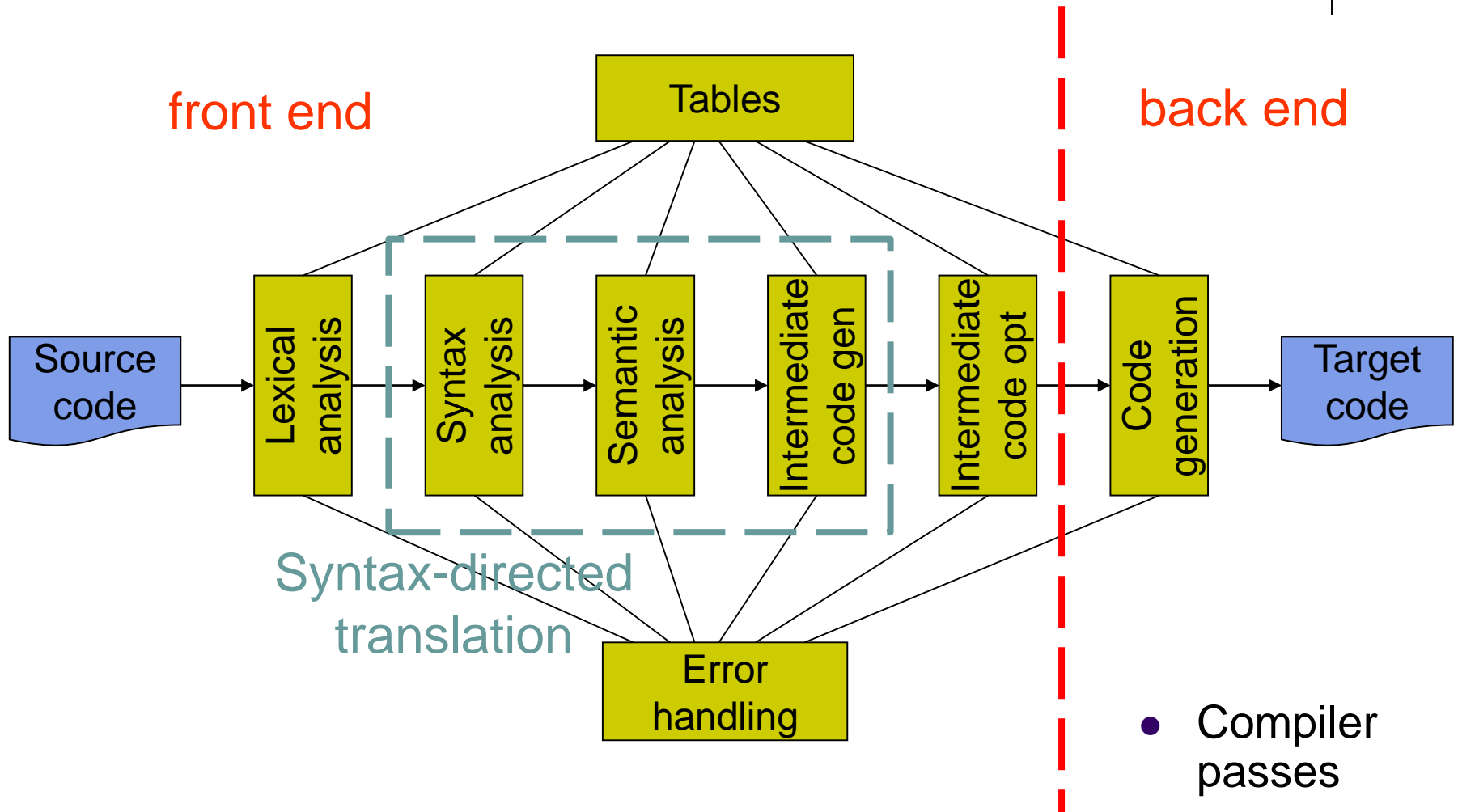




# Phases of a compiler

front end

back end





# Compiler-construction tools

- Parser generators
  - Produce syntax analyzers
  - Usually description based on a context-free grammar
  - Bison, Coco/R, ANTLR
- Scanner generators
  - Produce lexical analyzers
  - Usually description based on regular expressions
  - Flex
- Automatic code generators
  - Produce translations for each intermediate code instructions to the target code
  - A processor model and description
  - Mono JIT



# Compiler principles

---

Lexical analysis

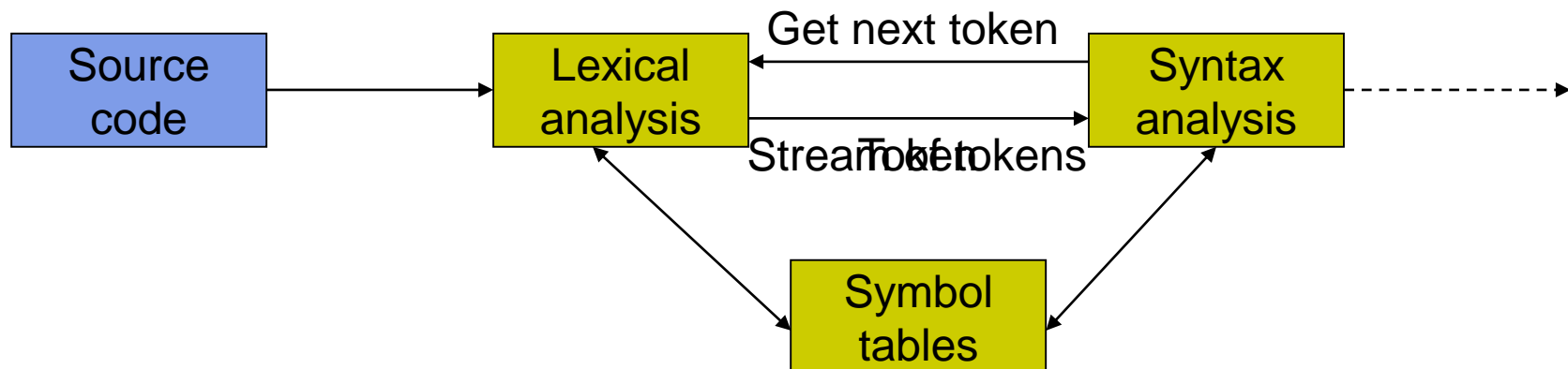
Jakub Yaghob





# Lexical analysis

- Reads input characters and produces a sequence of tokens
- Simpler design
  - Separation of lexical analysis from syntax analysis
- Specialization
  - Speedup
- Enhanced portability
  - Moving a compiler to other platform requires changes only in lexical analysis (EBCDIC)
- „Increasing“ grammar look-ahead
  - LR(1) does not mean look-ahead for 1 character
- Support for macros usually in lexical analysis





# Lexical analysis terms

- Token
  - Lexical analysis output and syntax analysis input
  - Called terminal on the syntax analysis/grammar side
  - A set of strings producing the same token
- Pattern
  - Rules describing a set of strings for given token
  - Usually described by regular expressions
- Lexical element
  - A sequence of characters in a source code corresponding with a pattern of a token
  - Some lexical elements do not have corresponding token
    - Comment
- Literal
  - A constant, has a value



# Examples

Token	Lexical element	Regular expression
while	while	while
relop	<, <=, =, <>, >, >=	\< \<= = \<> \> \>=
uint	0, 123	[0-9]+
	/* comment */	\* → cmt, <cmt>., <cmt>\*V

# Interesting problems in lexical analysis

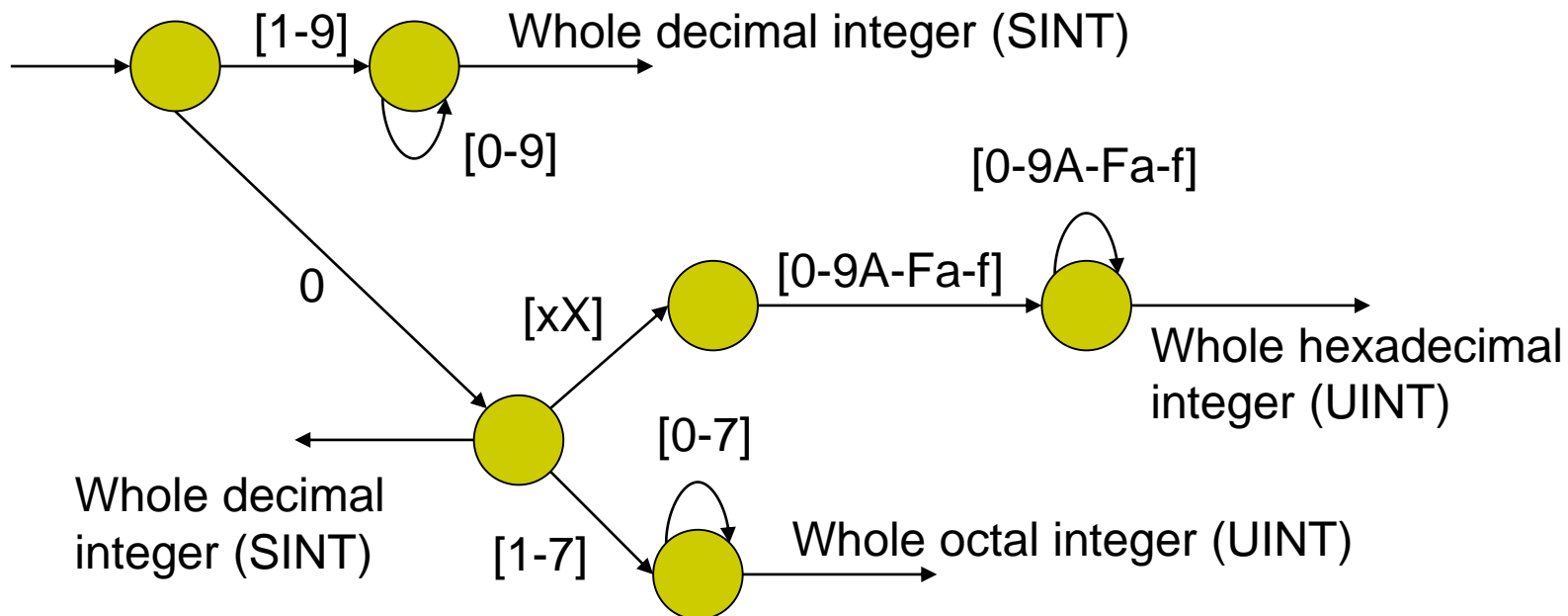


- Line indentation
  - Some languages have the indentation as a syntax construction
    - Python, Flex
- Identifiers
  - Identifiers with spaces
    - DO 5 l = 1.25
    - DO 5 l = 1,25
  - Keywords as identifiers
- Context dependent tokens
  - The value of a token depends on some other information
    - $a*b$ ;



# Lexical analysis background

- Patterns using regular expressions → regular languages → accepted by finite automata
- Restarting the automaton after each recognized/accepted token
- Finite automaton for an integer in C:





# Token attributes

- More patterns recognized as a token or the token is a literal
- Usually one attribute specifying more precisely a token or a literal value
  - Token=relop, specification='<='
  - Token=uint, specification='123'



# Lexical errors

- The finite automaton cannot continue and it is not in a final state
  - Unknown character
  - Unfinished string at the end of line
- Recovery
  - Ignore it
  - Deduce missing character(s)
- Typo in a keyword is not lexical error
  - `whle(f());`
- It can significantly influence syntax analysis





# Input buffering

- Lexical analysis takes 60-80% from the compile time
- One possible speedup: read the input file in blocks (buffers), the automaton works in the buffer memory
- Problems
  - Including a file means „including“ a buffer
    - #include

# Compiler principles

---

Syntax analysis

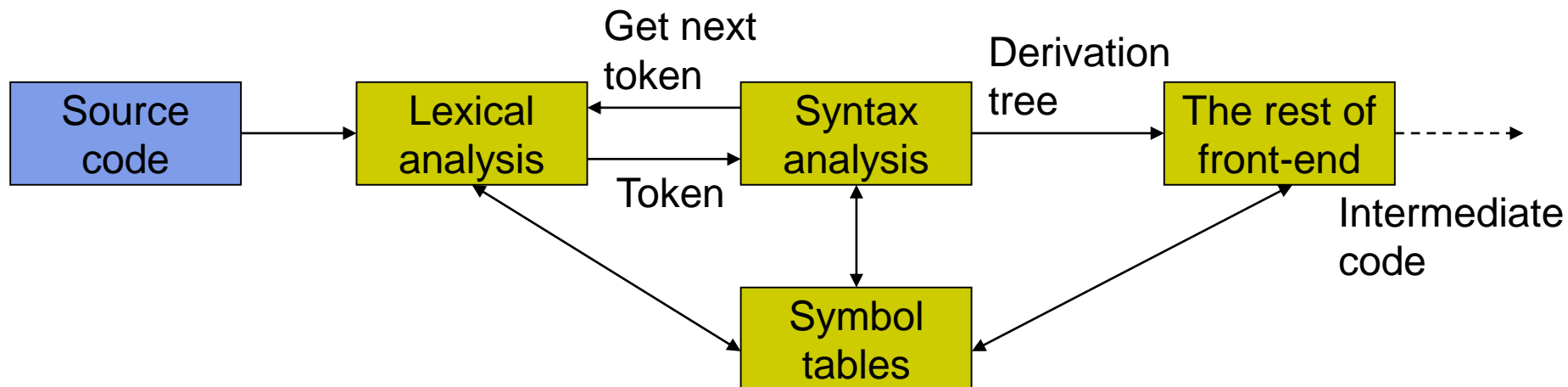
Jakub Yaghob





# Syntax analysis

- The main task
  - Decide, whether an input word is a word from an input language
- Other important tasks
  - Syntax-directed translation is the main loop of the compiler
  - Build the derivation tree
- Automaton type
  - We are talking about (deterministic) context-free grammars, therefore we are using (deterministic) pushdown automata





# Our grammar

1.  $E \rightarrow E + T$
2.  $E \rightarrow T$
3.  $T \rightarrow T * F$
4.  $T \rightarrow F$
5.  $F \rightarrow ( E )$
6.  $F \rightarrow \text{id}$

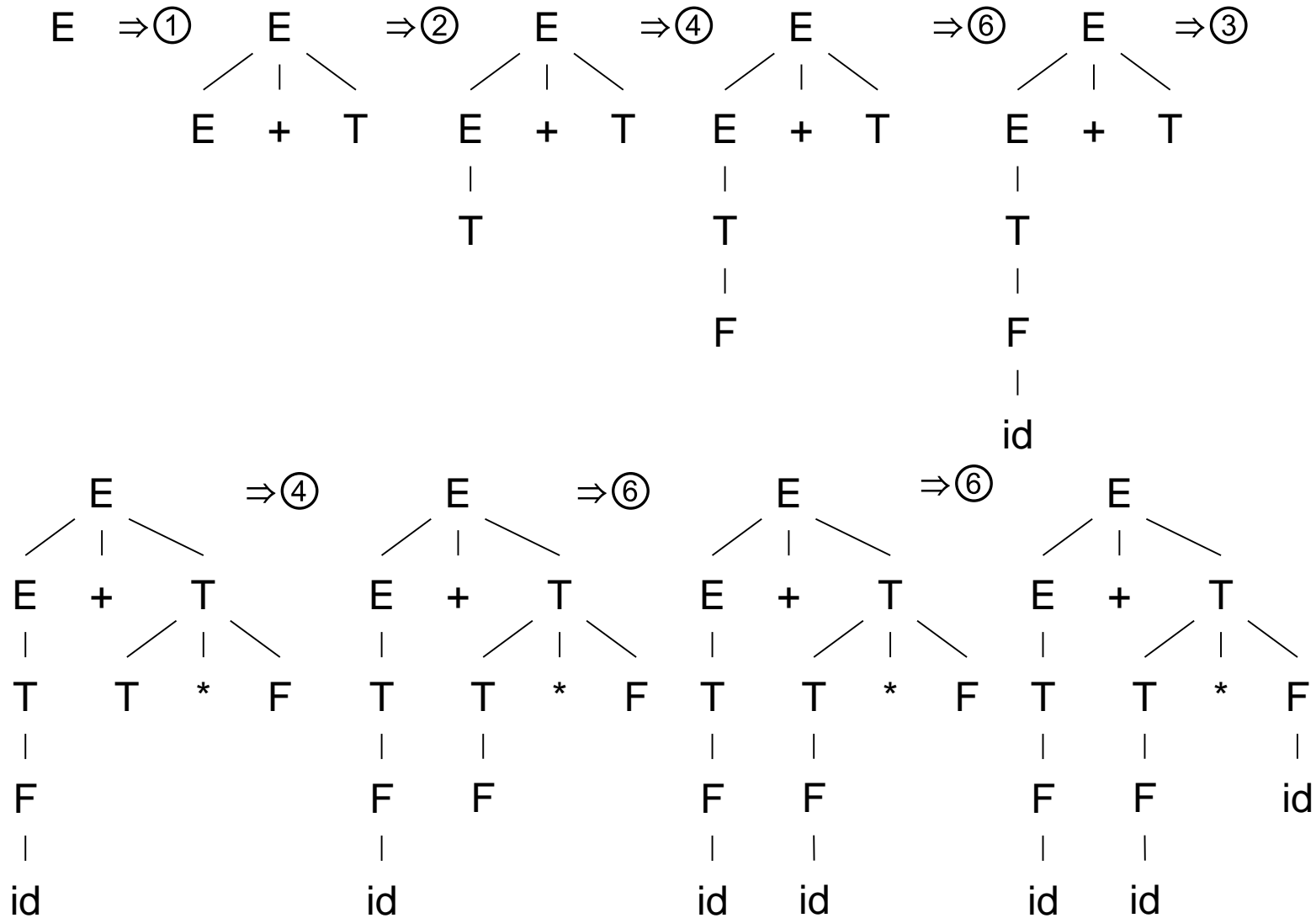


# Derivation (parse, syntax) tree

- Graphical representation of derivations using trees
  - Vertices are both non-terminals and terminals
  - Edges from inner vertex representing a non-terminal on the left side of a production rule to all symbols from the right side of a production rule
- $E \Rightarrow \textcircled{1} E+T \Rightarrow \textcircled{2} T+T \Rightarrow \textcircled{4} F+T \Rightarrow \textcircled{6} \text{id}+T \Rightarrow \textcircled{3} \text{id}+T * F \Rightarrow \textcircled{4} \text{id}+F * F \Rightarrow \textcircled{6} \text{id}+\text{id} * F \Rightarrow \textcircled{6} \text{id}+\text{id} * \text{id}$



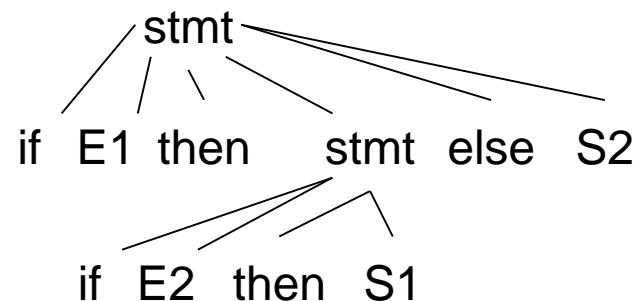
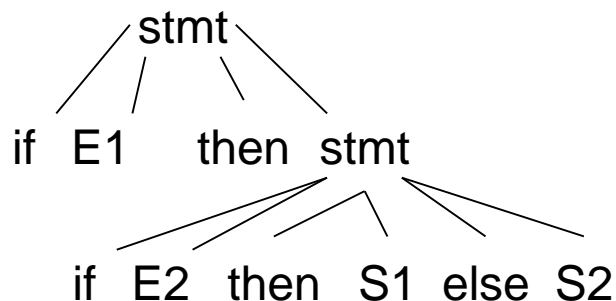
# Example





# Ambiguous grammar

- We can construct distinct derivation trees for the same input word
- Real-life example (dangling else):
  - $\text{stmt} \rightarrow$  **if** expr **then** stmt  
                  | **if** expr **then** stmt **else** stmt  
                  | **while** expr **do** stmt  
                  | **goto** num
  - Input word: **if**  $E_1$  **then** **if**  $E_2$  **then**  $S_1$  **else**  $S_2$





# Disambiguation

- Clarify, which tree is the right one
- In our case: **else** pairs with nearest “free” **if** (without **else**)
- Idea: “paired” statement is always between **if** and **else**
  - $\text{stmt} \rightarrow \text{m\_stmt}$
  - $\quad \quad \quad | \text{u\_stmt}$
  - $\text{m\_stmt} \rightarrow \text{if expr then m\_stmt else m\_stmt}$ 
    - $\quad \quad \quad | \text{while expr do m\_stmt}$
    - $\quad \quad \quad | \text{goto num}$
  - $\text{u\_stmt} \rightarrow \text{if expr then stmt}$ 
    - $\quad \quad \quad | \text{if expr then m\_stmt else u\_stmt}$
    - $\quad \quad \quad | \text{while expr do u\_stmt}$





# Left recursion elimination

- A grammar is a left-recursive grammar, when there is a non-terminal  $A$  for which it is true that  $A \Rightarrow^+ A\alpha$  for a string  $\alpha$
- It is a problem for top-down parsing
- A simple solution for  $\beta\alpha^m$ :
  - $A \rightarrow A\alpha$
  - $A \rightarrow \beta$
  - $A \rightarrow \beta A'$
  - $A' \rightarrow \alpha A'$
  - $A' \rightarrow \Lambda$

# Removing left recursion from our grammar



1.  $E \rightarrow E + T$

2.  $E \rightarrow T$

3.  $T \rightarrow T * F$

4.  $T \rightarrow F$

5.  $F \rightarrow ( E )$

6.  $F \rightarrow id$

1.  $E \rightarrow TE'$

2.  $E' \rightarrow + TE'$

3.  $E' \rightarrow \Lambda$

4.  $T \rightarrow FT'$

5.  $T' \rightarrow * FT'$

6.  $T' \rightarrow \Lambda$

7.  $F \rightarrow ( E )$

8.  $F \rightarrow id$



# Left factoring

- It is not clear, which option we should choose
  - $A \rightarrow \alpha\beta_1$
  - $A \rightarrow \alpha\beta_2$
  - $A \rightarrow \alpha A'$
  - $A' \rightarrow \beta_1$
  - $A' \rightarrow \beta_2$

# Non-context-free language constructions



- $L_1 = \{ w c w \mid w = (a|b)^* \}$ 
  - Check, whether an identifier **w** is declared before using
- $L_2 = \{ a^n b^m c^n d^m \mid n \geq 1, m \geq 1 \}$ 
  - Check, whether number of parameters in function call confirms to the function declaration
- $L_3 = \{ a^n b^n c^n \mid n \geq 0 \}$ 
  - The problem of “underscoring” a word
    - **a** is a char, **b** is BS, **c** is underscore
  - $(abc)^*$  is a regular expression

# Operators FIRST and FOLLOW

## – definitions



- If  $\alpha$  is any string of grammar symbols, let  $\text{FIRST}(\alpha)$  be the set of terminals that begin the strings derived from  $\alpha$ . If  $\alpha$  can be derived to  $\Lambda$ , then  $\Lambda$  is also in  $\text{FIRST}(\alpha)$
- Define  $\text{FOLLOW}(A)$ , for nonterminal  $A$ , to be the set of terminals that can appear immediately to the right of  $A$  in some string, where exists a derivation of the form  $S \Rightarrow^* \alpha A a \beta$  for some  $\alpha$  and  $\beta$ . If  $A$  can be the rightmost symbol in some sentential form, then  $\$$  is in  $\text{FOLLOW}(A)$ .

# Construction of the FIRST operator



- Construction for a grammar symbol  $X$ 
  - If  $X$  is terminal, then  $\text{FIRST}(X) = \{X\}$
  - If  $X \rightarrow \Lambda$  is a production, then add  $\Lambda$  to  $\text{FIRST}(X)$
  - If  $X$  is nonterminal and  $X \rightarrow Y_1 Y_2 \dots Y_k$  is a production, then place  $a$  in  $\text{FIRST}(X)$ , if for some  $i$ ,  $a$  is in  $\text{FIRST}(Y_i)$  and  $\Lambda \in \text{FIRST}(Y_j) \ \forall j < i$ . If  $\Lambda \in \text{FIRST}(Y_j) \ \forall j$ , then add  $\Lambda$  to  $\text{FIRST}(X)$
- Construction for any string
  - The construction of the FIRST operator for a string  $X_1 X_2 \dots X_n$  is similar as for nonterminal.

# Construction of the FOLLOW operator



- Construction for a nonterminal  $A$ 
  - Place  $\$$  in  $\text{FOLLOW}(S)$ , where  $S$  is the start symbol of a grammar and  $\$$  is EOS
  - If there is a production  $A \rightarrow \alpha B \beta$ , then everything in  $\text{FIRST}(\beta)$  except for  $\Lambda$  is placed in  $\text{FOLLOW}(B)$
  - If there is a production  $A \rightarrow \alpha B$  or  $A \rightarrow \alpha B \beta$  where  $\Lambda \in \text{FIRST}(\beta)$ , then everything in  $\text{FOLLOW}(A)$  is in  $\text{FOLLOW}(B)$

# FIRST and FOLLOW – an example for our grammar



- $\text{FIRST}(E) = \{ (, \text{id} \}$
- $\text{FIRST}(T) = \{ (, \text{id} \}$
- $\text{FIRST}(F) = \{ (, \text{id} \}$
- $\text{FIRST}(E') = \{ +, \Lambda \}$
- $\text{FIRST}(T') = \{ *, \Lambda \}$
- $\text{FOLLOW}(E) = \{ ), \$ \}$
- $\text{FOLLOW}(E') = \{ ), \$ \}$
- $\text{FOLLOW}(T) = \{ +, ), \$ \}$
- $\text{FOLLOW}(T') = \{ +, ), \$ \}$
- $\text{FOLLOW}(F) = \{ +, *, ), \$ \}$





# Top-down parsing

- An attempt to find a leftmost derivation for an input string
- An attempt to construct a parse tree for the input starting from the root and creating the nodes of the tree in preorder
- Recursive-descent parsing
  - Recursive descent using procedures
- Nonrecursive predictive parsing
  - An automaton with an explicit stack
- Both solutions have a problem with left recursion in a grammar
- Many current parser generators use top-down parsing
  - ANTLR, CocoR – LL(1) grammars with conflict resolution using dynamic look-ahead expansion to LL(k)



# Recursive-descent parsing

- One procedure/function for each nonterminal of a grammar
- Each procedure does two things
  - It decides, which grammar production with given nonterminal on the left side will be used using look-ahead. A production with right side  $\alpha$  will be used, when the look-ahead is in  $FIRST(\alpha)$ . If there is a conflict for the look-ahead among some production right sides, the grammar is not suitable for recursive-descent parsing. A production with  $\Lambda$  on the right side will be used, if the look-ahead is not in  $FIRST$  of any right side.
  - Procedure code copies the right side of a production. Nonterminal means calling a procedure for this nonterminal. Terminal is compared with the look-ahead. If they are equal, a next terminal is read. If they are not equal, it is an error.

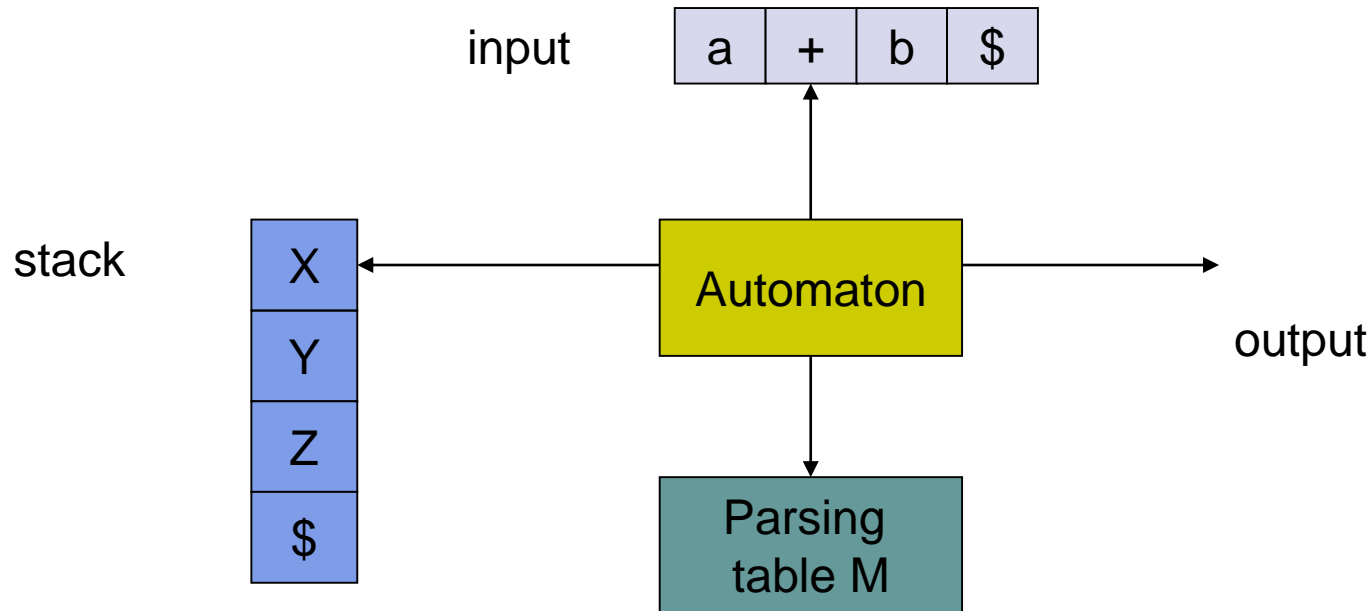
# Recursive-descent parsing – example for our grammar



```
void match(token t) {
    if(lookahead==t)
        lookahead = nexttoken();
    else error();
}
void E(void) {
    T(); Eap();
}
void Eap(void) {
    if(lookahead=='+') {
        match('+'); T(); Eap();
    }
}
void T(void) {
    F(); Tap();
}
```

```
void Tap(void) {
    if(lookahead=='*') {
        match('*'); F(); Tap();
    }
}
void F(void) {
    switch(lookahead) {
        case '(': match('('); E();
                    match(')'); break;
        case 'id':
                    match('id'); break;
        default:
                    error();
    }
}
```

# Nonrecursive predictive parsing



- Parsing table  $M[A, a]$ , where  $A$  is nonterminal and  $a$  is terminal
- The stack contains grammar symbols



# LL(1) automaton behavior

- Initial configuration
  - Input pointer points to the first terminal in the input string
  - The stack contains the start symbol of the grammar on top of \$
- In each step, the automaton decides, what to do, using a symbol  $X$  on top of the stack and a terminal  $a$ , pointed by input pointer
  - If  $X=a=\$$ , the parser halts, parsing finished successfully
  - If  $X=a\neq \$$ , the parser pops  $X$  from the stack and advances the input pointer to the next input symbol
  - If  $X\neq a$  and  $X\in T$ , the parser reports error
  - If  $X$  is a nonterminal, the parser uses entry  $M[X, a]$ . If this entry is a production, the parser replaces  $X$  on top of the stack by the right side (leftmost symbol on top of the stack). At the same time, the parser generates an output about using the production. If the entry is **error**, the parser informs about a syntax error.

# Construction of predictive parsing tables



- For each production  $A \rightarrow \alpha$  do following steps
  - For  $\forall a \in \text{FIRST}(\alpha)$  add  $A \rightarrow \alpha$  to  $M[A, a]$
  - If  $\Lambda \in \text{FIRST}(\alpha)$ , add  $A \rightarrow \alpha$  to  $M[A, b]$   
 $\forall b \in \text{FOLLOW}(A)$ . Moreover, if  $\$ \in \text{FOLLOW}(A)$ ,  
add  $A \rightarrow \alpha$  to  $M[A, \$]$
- Mark each empty entry in  $M$  as **error**

# Example of table construction for our grammar



	id	+	*	(	)	\$
E	$E \rightarrow TE'$			$E \rightarrow TE'$		
E'		$E' \rightarrow +TE'$			$E' \rightarrow \Lambda$	$E' \rightarrow \Lambda$
T	$T \rightarrow FT'$			$T \rightarrow FT'$		
T'		$T' \rightarrow \Lambda$	$T' \rightarrow *FT'$		$T' \rightarrow \Lambda$	$T' \rightarrow \Lambda$
F	$F \rightarrow \text{id}$			$F \rightarrow (E)$		

# Example of parser behavior for our grammar



Stack	Input	Output
\$E	id+id*id\$	
\$E'T	id+id*id\$	$E \rightarrow TE'$
\$E'T'F	id+id*id\$	$T \rightarrow FT'$
\$E'T'id	id+id*id\$	$F \rightarrow id$
\$E'T'	+id*id\$	
\$E'	+id*id\$	$T' \rightarrow \Lambda$
\$E'T+	+id*id\$	$E' \rightarrow +TE'$
\$E'T	id*id\$	
\$E'T'F	id*id\$	$T \rightarrow FT'$

Stack	Input	Output
\$E'T'id	id*id\$	$F \rightarrow id$
\$E'T'	*id\$	
\$E'T'F*	*id\$	$T' \rightarrow *FT'$
\$E'T'F	id\$	
\$E'T'id	id\$	$F \rightarrow id$
\$E'T'	\$	
\$E'	\$	$T' \rightarrow \Lambda$
\$	\$	$E' \rightarrow \Lambda$





# LL(1) grammar

- Context-free grammar  $G=(T,N,S,P)$  is a LL(1) grammar, if and only if whenever  $A \rightarrow \alpha$ ,  $A \rightarrow \beta \in P$  are two distinct ( $\alpha \neq \beta$ ) productions of  $G$  and we have any left sentential forms  $uA\gamma$ ,  $vA\delta$ , where  $u,v \in T^*$  and  $\gamma, \delta \in (T \cup N)^*$ , the following condition holds:
  - $\text{FIRST}(\alpha\gamma) \cap \text{FIRST}(\beta\delta) = \emptyset$
- Simplified detection: no ambiguous or left-recursive grammar can be LL(1)



# Grammar terminology

- $PXY(k)$
- $X$  – direction of the input reading
  - In our case always  $L$ , i.e. from left to right
- $Y$  – kind of derivation
  - $L$  – left derivation
  - $R$  – right derivation
- $P$  – prefix
  - Subtle division of some grammar classes
- $k$  – look-ahead
  - An integer, usually 1, can be 0 or more generally  $k$
- Examples
  - $LL(1)$ ,  $LR(0)$ ,  $LR(1)$ ,  $LL(k)$ ,  $SLR(1)$ ,  $LALR(1)$

# Expanding definition of FIRST and FOLLOW on k



- If  $\alpha$  is a string from grammar symbols, then  $\text{FIRST}_k(\alpha)$  is a set of terminal words with maximal length  $k$ , which are on the beginning of at least one string derived from  $\alpha$ . If  $\alpha$  can be derived on  $\Lambda$ , then  $\Lambda$  is in  $\text{FIRST}_k(\alpha)$ .
- $\text{FOLLOW}_k(A)$  for nonterminal  $A$  is a set of terminal words with maximal length  $k$ , which can be on the right side of  $A$  in any string derived from the start nonterminal ( $S \Rightarrow^* \alpha A u \beta$  for some  $\alpha$  and  $\beta$ ). If  $A$  is the right-most symbol in any sentential form, then  $\$$  is in  $\text{FOLLOW}_k(A)$ .



# LL(k) grammar

- Context-free grammar  $G=(T,N,S,P)$  is a strong LL(k) grammar for  $k \geq 1$ , if and only if whenever  $A \rightarrow \alpha, A \rightarrow \beta \in P$  are two distinct ( $\alpha \neq \beta$ ) productions and we have any left sentential forms  $uA\gamma, vA\delta$ , where  $u, v \in T^*$  and  $\gamma, \delta \in (T \cup N)^*$ , the following condition holds:
  - $FIRST_k(\alpha\gamma) \cap FIRST_k(\beta\delta) = \emptyset$ .
- LL(k) (not strong)
  - $u=v, \gamma=\delta$

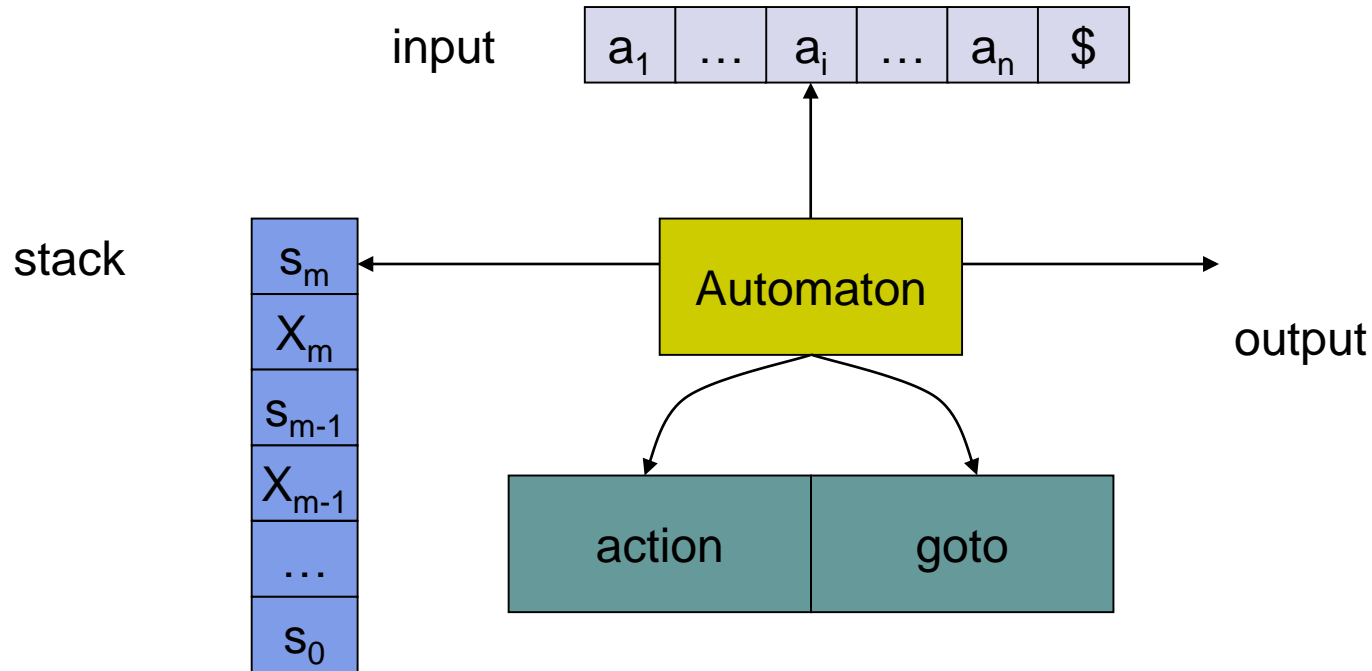


# Bottom-up analysis

- Attempts to find in reverse the rightmost derivation for an input string
- Attempts to construct a parse tree for an input string beginning at the leaves and working up towards the root
- Replace a substring corresponding to a right side of a production by a nonterminal from the left side of the production in each reduce step
- Used in parser generators
  - Bison – LALR(1), GLR(1)
- Advantages against LL(1) parsers
  - It can be implemented with the same efficiency as top-down parsing
  - The class of decidable languages LR(1) is a proper superset of LL(1)
- SLR(1), LR(1), LALR(1)



# LR parser automaton



- $s_i$  are states
  - A state on the top of the stack is the current state of the automaton
- $x_i$  are grammar symbols



# LR(1) automaton behavior

- Initial configuration
  - Input pointer points to the first terminal in the input string
  - Initial state  $s_0$  is on the stack
- In each step address table action[ $s_m, a_i$ ] using  $s_m$  and  $a_i$ 
  - Shift  $s$ , where  $s$  is a new state
    - It shifts the input tape by 1 terminal and add  $a_i$  and  $s$  on the top of the stack
  - Reduce using production  $A \rightarrow \alpha$ 
    - Remove  $r=|\alpha|$  pairs  $(s_k, X_k)$  from the top of the stack, add  $A$  on the top of the stack and then goto[ $s_{m-r}, A$ ] ( $s_{m-r}$  is a state on the top of the stack after erasing pairs)
    - Generate an output
  - Accept
    - The input string is accepted
    - Generate an output
  - Error
    - The input string is not in the input language

# LR automaton tables for our grammar



state	action						goto		
	id	+	*	(	)	\$	E	T	F
0	s5			s4			1	2	3
1		s6				acc			
2		r2	s7		r2	r2			
3		r4	r4		r4	r4			
4	s5			s4			8	2	3
5		r6	r6		r6	r6			
6	s5			s4				9	3
7	s5			s4					10
8		s6			s11				
9		r1	s7		r1	r1			
10		r3	r3		r3	r3			
11		r5	r5		r5	r5			





# Example of LR parser behavior

Stack	Input	Action
0	<b>id+id*id\$</b>	s5
0 <b>id</b> 5	<b>+id*id\$</b>	r6: F→ <b>id</b>
0 F 3	<b>+id*id\$</b>	r4: T→F
0 T 2	<b>+id*id\$</b>	r2: E→T
0 E 1	<b>+id*id\$</b>	s6
0 E 1 <b>+</b> 6	<b>id*id\$</b>	s5
0 E 1 <b>+</b> 6 <b>id</b> 5	<b>*id\$</b>	r6: F→ <b>id</b>
0 E 1 <b>+</b> 6 F 3	<b>*id\$</b>	r4: T→F
0 E 1 <b>+</b> 6 T 9	<b>*id\$</b>	s7
0 E 1 <b>+</b> 6 T 9 <b>*</b> 7	<b>id\$</b>	s5
0 E 1 <b>+</b> 6 T 9 <b>*</b> 7 <b>id</b> 5	<b>\$</b>	r6: F→ <b>id</b>
0 E 1 <b>+</b> 6 T 9 <b>*</b> 7 F 10	<b>\$</b>	r3: T→T * F
0 E 1 <b>+</b> 6 T 9	<b>\$</b>	r1: E→E + T
0 E 1	<b>\$</b>	acc



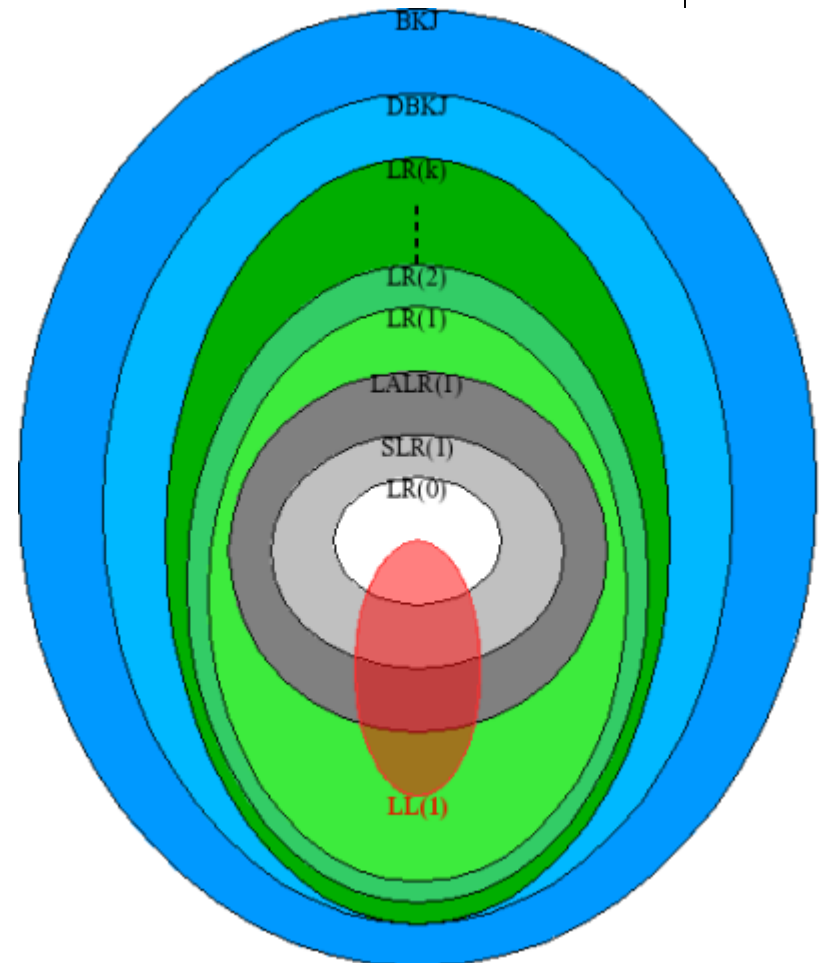
# LR(k) grammar

- Context-free grammar  $G=(T,N,S,P)$  is LR(k) grammar for  $k \geq 1$ , if and only if whenever  $A \rightarrow \alpha, A \rightarrow \beta \in P$  are two distinct ( $\alpha \neq \beta$ ) productions of  $G$  and we have any two right sentential forms  $\gamma Au, \delta Av$ , where  $u, v \in T^*$  and  $\gamma, \delta \in (T \cup N)^*$ , the following condition holds:
  - $\text{FIRST}_k(u) \cap \text{FIRST}_k(v) = \emptyset$

# Grammars (languages) strength



- Union of all  $LR(k)$  are deterministic context-free languages (DBKJ) and it is a proper subset of all context-free languages (BKJ)





# Grammar augmentation

- Augmentation of a grammar  $G=(T,N,S,P)$  is a grammar  $G'=(T,N',S',P')$ , where  $N'=N\cup\{S'\}$  and  $P'=P\cup\{S'\rightarrow S\}$
- The augmentation is not necessary whenever  $S$  is on the left side of one production and it isn't on any right side of grammar productions
- It helps recognize the end of parsing
- For our grammar:
  - $S'\rightarrow E$



# LR(0) items

- LR(0) item of a grammar  $G$  is a production with a special symbol dot on the right side
  - Special symbol is a valid symbol for comparison of two LR(0) items of a same production. LR(0) items of the same production are different, whenever the dot is on different position. Moreover, the dot is not a grammar symbol
- An example for production  $E \rightarrow E + T$ :

$E \rightarrow \blacklozenge E + T$	$E \rightarrow E + \blacklozenge T$
$E \rightarrow E \blacklozenge + T$	$E \rightarrow E + T \blacklozenge$



# The closure operation

- If  $I$  is a set of LR(0) items for a grammar  $G$ , then  $\text{CLOSURE}(I)$  is a set of LR(0) items constructed from  $I$  by following rules:
  - Add  $I$  to the  $\text{CLOSURE}(I)$
  - $\forall A \rightarrow \alpha \blacklozenge B \beta \in \text{CLOSURE}(I)$ , where  $B \in N$ , add  $\forall B \rightarrow \gamma \in P$  to  $\text{CLOSURE}(I)$  LR(0) item  $B \rightarrow \blacklozenge \gamma$ , if it is not already there. Apply this rule until no more new LR(0) items can be added to  $\text{CLOSURE}(I)$

# Example of closure for our grammar



- $I = \{S' \rightarrow \blacklozenge E\}$
- $\text{CLOSURE}(I) =$ 
  - $S' \rightarrow \blacklozenge E$
  - $E \rightarrow \blacklozenge E + T$
  - $E \rightarrow \blacklozenge T$
  - $T \rightarrow \blacklozenge T * F$
  - $T \rightarrow \blacklozenge F$
  - $F \rightarrow \blacklozenge ( E )$
  - $F \rightarrow \blacklozenge \text{id}$



# GOTO operation

- GOTO( $I, X$ ) operation for a set  $I$  of LR(0) items and a grammar symbol  $X$  is defined to be the closure of the set of all LR(0) items  $A \rightarrow \alpha X \diamond \beta$  such that  $A \rightarrow \alpha \diamond X \beta \in I$

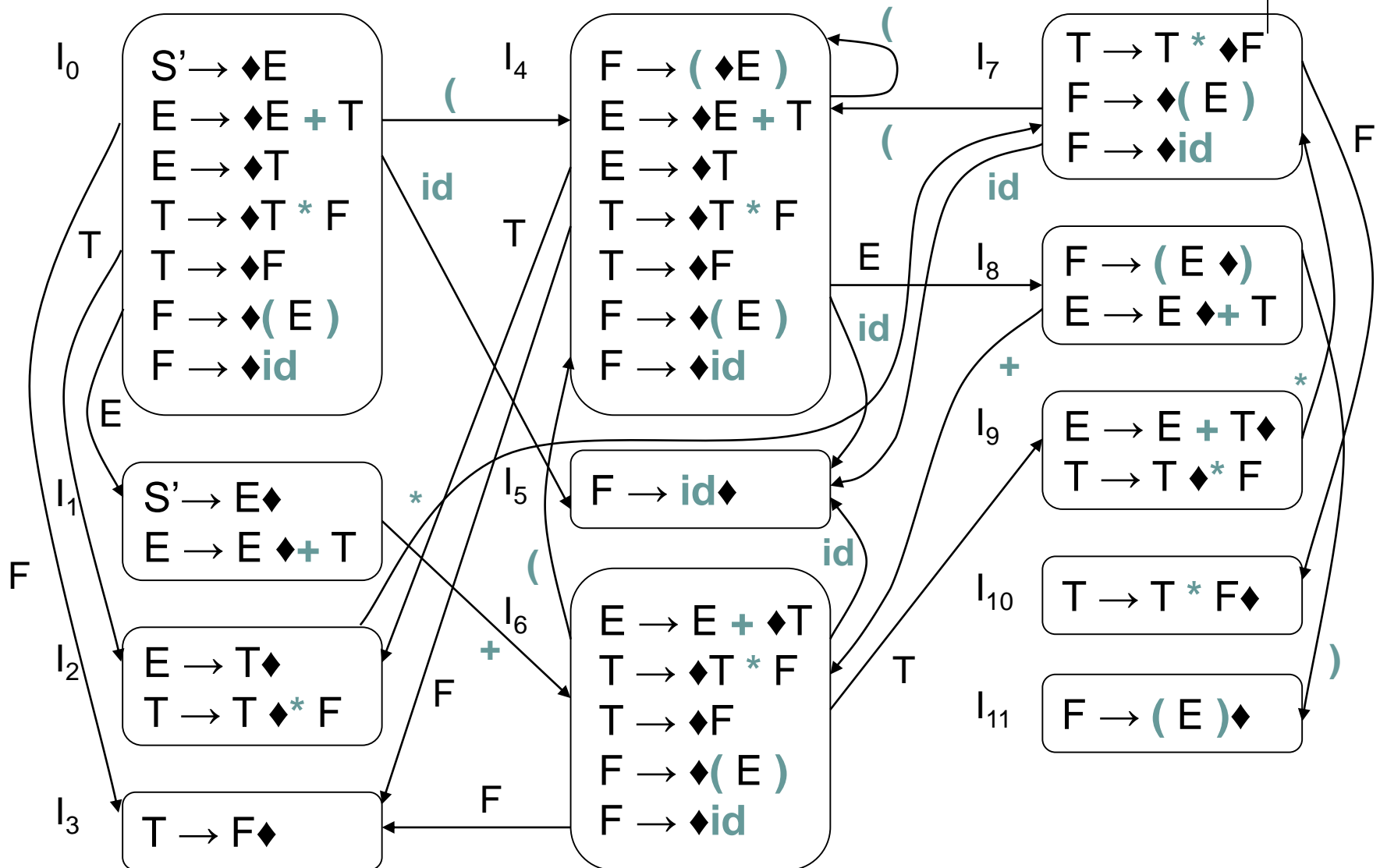


# Construction of canonical collection of sets of LR(0) items



- We have an augmented grammar  $G'=(T,N',S',P')$
- Construction of canonical collection  $C$  of sets of LR(0) items:
  - We start with  $C=\{ \text{CLOSURE}(\{S' \rightarrow \blacklozenge S\}) \}$
  - $\forall I \in C$  and  $\forall X \in T \cup N'$  such as  $\text{GOTO}(I, X) \notin C \wedge \text{GOTO}(I, X) \neq \emptyset$ , add  $\text{GOTO}(I, X)$  to  $C$ . Repeat this step, until something new is added to  $C$ .

# Construction of canonical collection for our grammar





# Valid items

- LR(0) item  $A \rightarrow \beta_1 \blacklozenge \beta_2$  is a valid item for a viable prefix  $\alpha\beta_1$ , if there is a rightmost derivation  $S' \Rightarrow^+ \alpha A w \Rightarrow \alpha\beta_1\beta_2 w$
- It is a great hint for a parser. It helps to decide, if the parser should make a shift or a reduction, if  $\alpha\beta_1$  is on top of the stack
- Basic LR parsing theorem: A set of valid items for a viable prefix  $\gamma$  is exactly a set of items reachable from the initial state through the prefix  $\gamma$  by deterministic finite automaton constructed from canonical collection with GOTO transitions.

# SLR(1) automaton construction



- We have an augmented grammar  $G'$ . Tables of SLR(1) automaton are constructed by following algorithm
  - Construct a canonical collection  $C$  of sets of LR(0) items
  - State  $i$  is constructed from  $I_i$ . The parsing actions for state  $i$  are determined as follows
    - $A \rightarrow \alpha \diamond a \beta \in I_i, a \in T \wedge \text{GOTO}(I_i, a) = I_j$ , then  $\text{action}[i, a] = \text{shift } j$
    - $A \rightarrow \alpha \diamond \in I_i$ , then  $\forall a \in \text{FOLLOW}(A) \wedge A \neq S'$  is  $\text{action}[i, a] = \text{reduce } A \rightarrow \alpha$
    - $S' \rightarrow S \diamond \in I_i$ , then  $\text{action}[i, \$] = \text{accept}$
  - If there is a conflict in the previous step, the grammar is not a SLR(1) grammar and the automaton cannot be constructed
  - Table goto is indexed by state  $i$  and  $A \in N'$ : whenever  $\text{GOTO}(I_i, A) = I_j$ , then  $\text{goto}[i, A] = j$
  - All empty cells are filled by error instruction
  - The initial state of the parser is the state, which contains LR(0) item  $S' \rightarrow \diamond S$



# Full LR(1) automata

- $\text{action}[i, a]$  is set to reduction  $A \rightarrow \alpha$ , when  $A \rightarrow \alpha \diamond \in I_i$ ,  $\forall a \in \text{FOLLOW}(A)$  for a state  $i$  during SLR(1) construction
- In some situation, when  $i$  is on top of the stack, the viable prefix  $\beta\alpha$  is in form, where  $\beta A$  cannot be followed by a terminal  $a$  in any right sentential form. Therefore reduction  $A \rightarrow \alpha$  is for lookahead  $a$  invalid.
- Solution: add more information to states, so we can avoid invalid reductions.



# LR(1) items

- The added information is stored as an additional terminal for each LR(0) item. Such item has a form  $[A \rightarrow \alpha \blacklozenge \beta, a]$ , where  $A \rightarrow \alpha \beta \in P$ ,  $a \in T$ , and we call it LR(1) item. The terminal **a** is called lookahead.
  - The lookahead has no meaning for  $A \rightarrow \alpha \blacklozenge \beta$ , where  $\beta \neq \Lambda$
  - Reduction  $A \rightarrow \alpha$  is set only when  $[A \rightarrow \alpha \blacklozenge, a] \in I_i$  for current state  $i$  and a terminal **a** on the input
  - A set of terminals created from lookaheads of LR(1) items  $\subseteq \text{FOLLOW}(A)$
- LR(1) item  $[A \rightarrow \alpha \blacklozenge \beta, a]$  is valid for viable prefix  $\gamma$ , whenever  $\exists$  right derivation  $S \Rightarrow^+ \delta A w \Rightarrow \delta \alpha \beta w$ , where
  - $\gamma = \delta \alpha$
  - Either **a** is the first symbol of  $w$  or  $w = \Lambda$  and **a** is \$



# Closure for LR(1) items

- We have a set of LR(1) items  $I$  for a grammar  $G$ . We define  $CLOSURE1(I)$  as a set of LR(1) items constructed from  $I$  by following procedure:
  - Add set  $I$  to  $CLOSURE1(I)$
  - $\forall [A \rightarrow \alpha \blacklozenge B \beta, a] \in CLOSURE1(I)$ , where  $B \in N$ , add LR(1) item  $[B \rightarrow \blacklozenge \gamma, b] \forall B \rightarrow \gamma \in P$  and  $\forall b \in FIRST(\beta a)$  to  $CLOSURE1(I)$ , if it isn't there already. Repeat this step, until something is added to  $CLOSURE1(I)$ .

# GOTO operation for LR(1) items



- We define  $GOTO_1(I, X)$  operation for a set  $I$  of LR(1) items and a grammar symbol  $X$  as a  $CLOSURE_1$  of a set of all items  $[A \rightarrow \alpha X \blacklozenge \beta, a]$  where  $[A \rightarrow \alpha \blacklozenge X \beta, a] \in I$



# Construction of canonical collection of sets of LR(1) items



- We have an augmented grammar  $G'=(T,N',S',P')$
- Construction of canonical collection  $C$  of LR(1) items:
  - We start with  $C=\{ \text{CLOSURE}_1(\{[S' \rightarrow \diamond S, \$]\}) \}$
  - Add  $\text{GOTO}_1(I, X)$  to  $C \forall I \in C$  and  $\forall X \in T \cup N'$ , where  $\text{GOTO}_1(I, X) \notin C \wedge \text{GOTO}_1(I, X) \neq \emptyset$ . Repeat this step, until something new is added to  $C$ .

# Example of LR(1) grammar, which is not SLR(1)



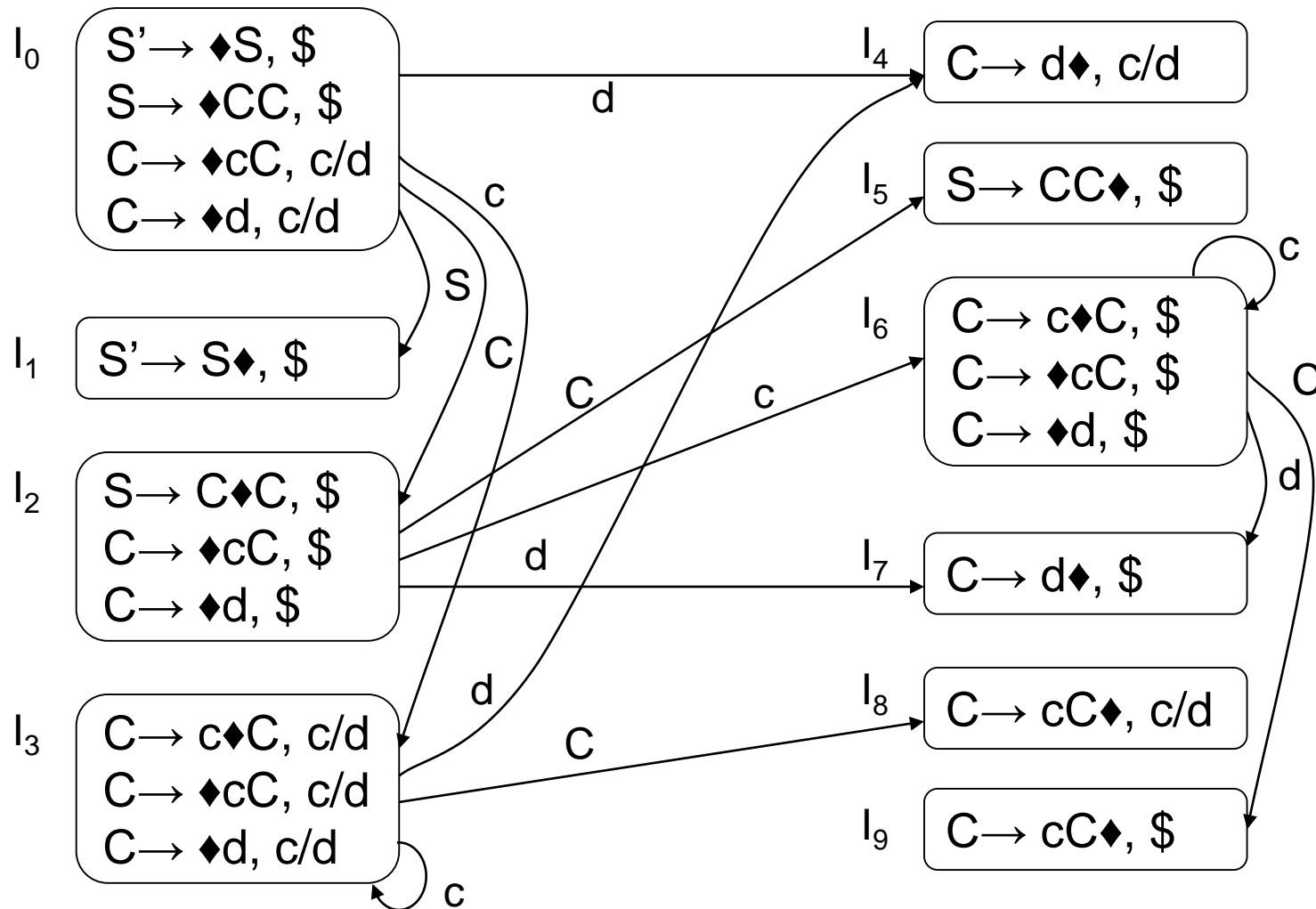
- $S' \rightarrow S$
- $S \rightarrow CC$
- $C \rightarrow cC$
- $C \rightarrow d$

# Example of closure construction for LR(1) items



- $I = \{[S' \rightarrow \blacklozenge S, \$]\}$
- $\text{CLOSURE}_1(I) =$ 
  - $S' \rightarrow \blacklozenge S, \$ \quad \beta = \Lambda, \text{FIRST}(\beta \$) = \text{FIRST}(\$) = \{\$\}$
  - $S \rightarrow \blacklozenge CC, \$ \quad \beta = C, \text{FIRST}(C \$) = \{c, d\}$
  - $C \rightarrow \blacklozenge cC, c/d$
  - $C \rightarrow \blacklozenge d, c/d$

# Example of construction of canonical collection of LR(1) items





# LR(1) parser construction

- We have an augmented grammar  $G'$ . LR(1) automaton tables are constructed by following algorithm
  - Construct a canonical collection  $C$  of sets of LR(1) items
  - State  $i$  is constructed from  $I_i$ . The parsing actions for state  $i$  are determined as follows
    - $[A \rightarrow \alpha \blacklozenge a \beta, b] \in I_i, a \in T \wedge \text{GOTO1}(I_i, a) = I_j$ , then  $\text{action}[i, a] = \text{shift } j$
    - $[A \rightarrow \alpha \blacklozenge, a] \in I_i \wedge A \neq S'$ , then  $\text{action}[i, a] = \text{reduce } A \rightarrow \alpha$
    - $[S' \rightarrow S \blacklozenge, \$] \in I_i$ , then  $\text{action}[i, \$] = \text{accept}$
  - If there is a conflict in the previous step, the grammar is not a LR(1) grammar and the automaton cannot be constructed
  - Table goto is indexed by state  $i$  and  $A \in N'$ : whenever  $\text{GOTO1}(I_i, A) = I_j$ , then  $\text{goto}[i, A] = j$
  - All empty cells are filled by error instruction
  - The initial state of the parser is the state, which contains LR(1) item  $[S' \rightarrow \blacklozenge S, \$]$



# LALR

- LALR=LookAhead-LR
- Often used in practice
  - Bison
  - Most common programming languages can be expressed by an LALR grammar
  - Parser tables are considerably smaller than LR(1) tables
- SLR and LALR parsers have the same number of states, LR parsers have greater number of states
  - Common languages have hundreds of states
  - LR(1) parsers have thousands of states for the same grammar



# How to make smaller tables?

- Idea: merge sets with the same core into one set including GOTO1 merge
  - Core: a set of LR(0) items (no lookahead)
  - Merge cannot produce shift/reduce conflict
    - Suppose in the union there is a conflict on lookahead **a** for LR(1) items  $[A \rightarrow \alpha \blacklozenge, a]$  and  $[B \rightarrow \beta \blacklozenge a \gamma, b]$
    - Cores are same, therefore in the set with  $[A \rightarrow \alpha \blacklozenge, a]$  must be  $[B \rightarrow \beta \blacklozenge a \gamma, c]$  as well for some **c**. There was already a shift/reduce conflict before merge
  - Merge can produce reduce/reduce conflict

# Easy LALR(1) table construction



- We have an augmented grammar  $G'$ . LALR(1) automaton tables are constructed by following algorithm
  - Construct a canonical collection  $C$  of sets of LR(1) items
  - For each core in collection  $C$ , find all sets having that core, and replace these sets by their union
  - Let  $C' = \{ J_0, J_1, \dots, J_m \}$  be the resulting collection of LR(1) items
  - Table action is constructed for  $C'$  in the same manner as for full LR(1) parser
  - If there is a conflict, the grammar is not LALR(1) grammar
  - If  $J \in C'$  is the union of sets of LR(1) items  $I_i$  ( $J = I_1 \cup I_2 \cup \dots \cup I_k$ ), then cores  $GOTO1(I_1, X), \dots, GOTO1(I_k, X)$  are the same, since  $I_1, \dots, I_k$  all have the same core. Let  $K$  be the union of all sets of items having the same core as  $goto(I_1, X)$ . Then  $GOTO1(J, X) = K$
- Important disadvantage – we need to construct full LR(1)



# Compiler principles

---

Syntax analysis

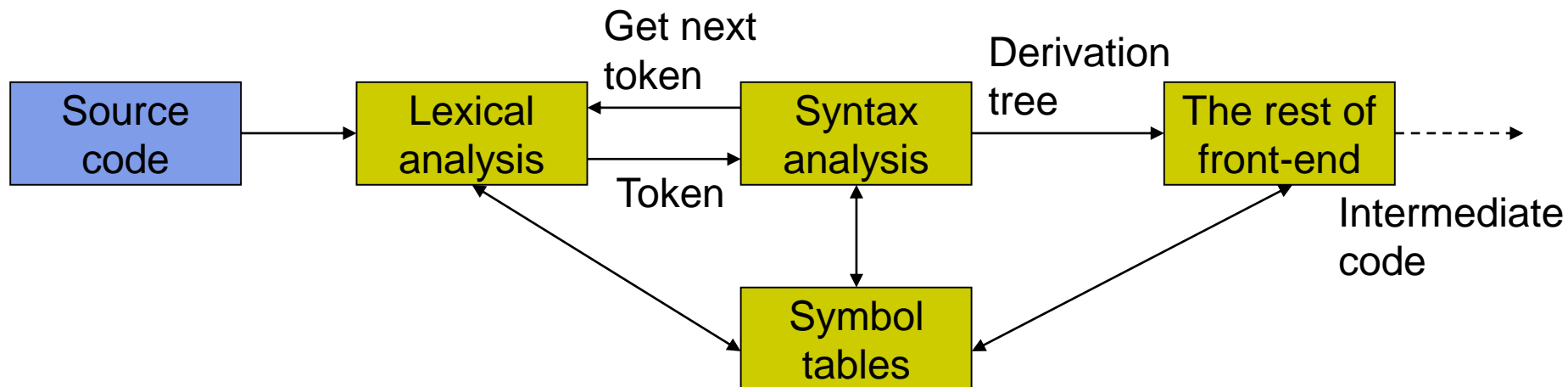
Jakub Yaghob





# Syntax analysis

- The main task
  - Decide, whether an input word is a word from an input language
- Other important tasks
  - Syntax-directed translation is the main loop of the compiler
  - Build the derivation tree
- Automaton type
  - We are talking about (deterministic) context-free grammars, therefore we are using (deterministic) pushdown automata





# Our grammar

1.  $E \rightarrow E + T$
2.  $E \rightarrow T$
3.  $T \rightarrow T * F$
4.  $T \rightarrow F$
5.  $F \rightarrow ( E )$
6.  $F \rightarrow \text{id}$

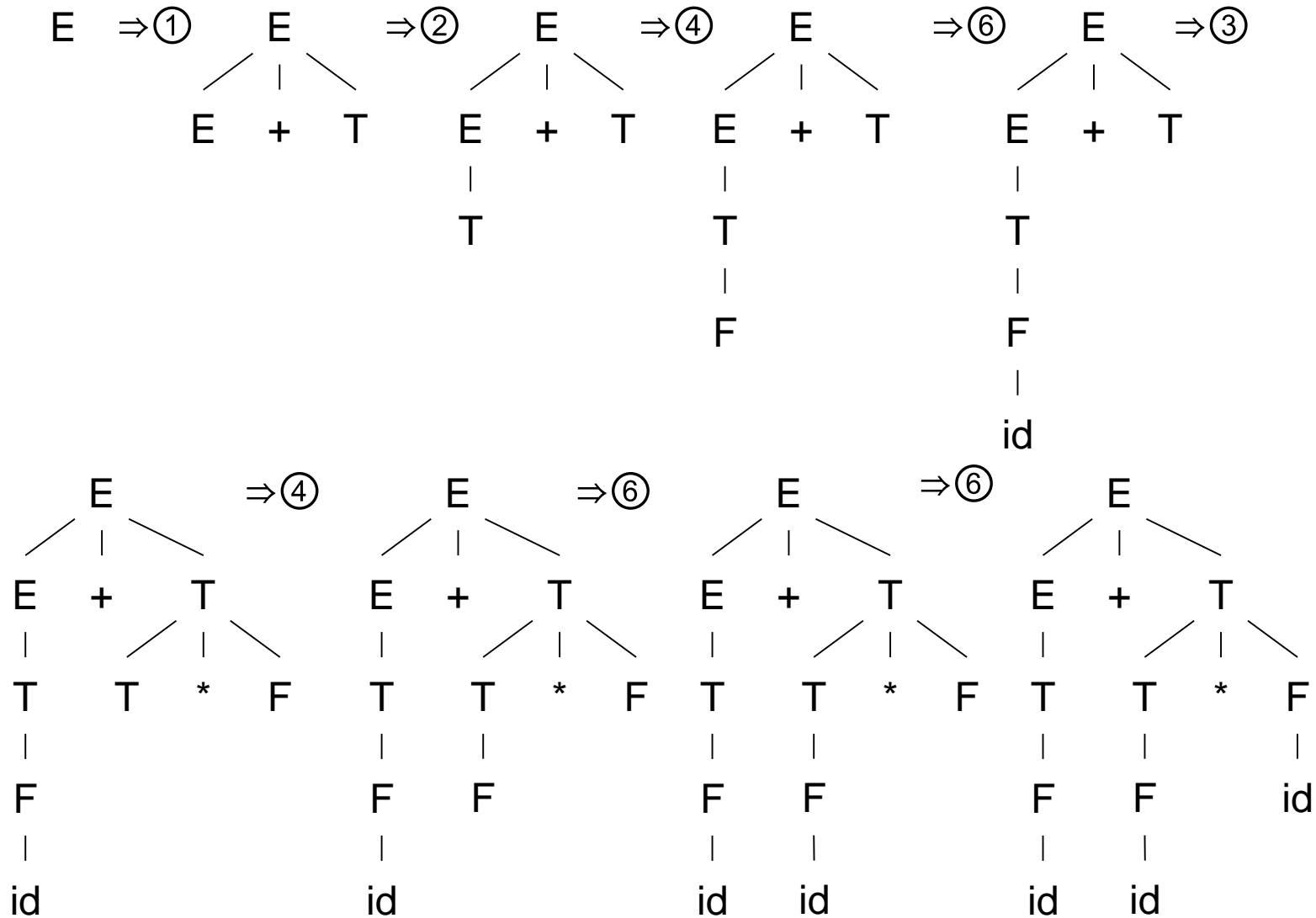


# Derivation (parse, syntax) tree

- Graphical representation of derivations using trees
  - Vertices are both non-terminals and terminals
  - Edges from inner vertex representing a non-terminal on the left side of a production rule to all symbols from the right side of a production rule
- $E \Rightarrow \textcircled{1} E+T \Rightarrow \textcircled{2} T+T \Rightarrow \textcircled{4} F+T \Rightarrow \textcircled{6} \text{id}+T \Rightarrow \textcircled{3} \text{id}+T^*F \Rightarrow \textcircled{4} \text{id}+F^*F \Rightarrow \textcircled{6} \text{id}+\text{id}^*F \Rightarrow \textcircled{6} \text{id}+\text{id}^*\text{id}$



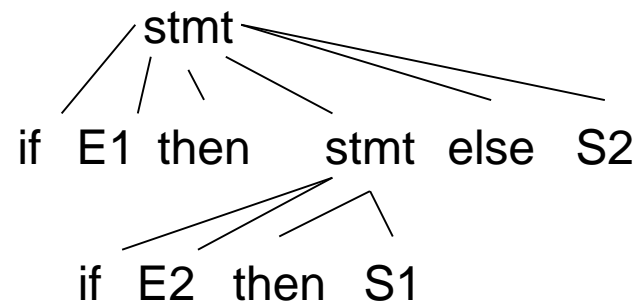
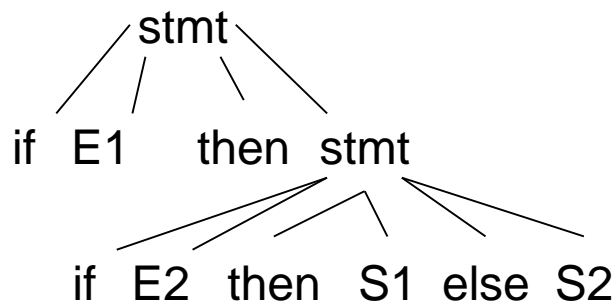
# Example





# Ambiguous grammar

- We can construct distinct derivation trees for the same input word
- Real-life example (dangling else):
  - $\text{stmt} \rightarrow$  **if** expr **then** stmt  
| **if** expr **then** stmt **else** stmt  
| **while** expr **do** stmt  
| **goto** num
  - Input word: **if**  $E_1$  **then** **if**  $E_2$  **then**  $S_1$  **else**  $S_2$





# Disambiguation

- Clarify, which tree is the right one
- In our case: **else** pairs with nearest “free” **if** (without **else**)
- Idea: “paired” statement is always between **if** and **else**
  - $\text{stmt} \rightarrow \text{m\_stmt}$
  - $\quad \quad \quad | \text{u\_stmt}$
  - $\text{m\_stmt} \rightarrow \text{if expr then m\_stmt else m\_stmt}$ 
    - $\quad \quad \quad | \text{while expr do m\_stmt}$
    - $\quad \quad \quad | \text{goto num}$
  - $\text{u\_stmt} \rightarrow \text{if expr then stmt}$ 
    - $\quad \quad \quad | \text{if expr then m\_stmt else u\_stmt}$
    - $\quad \quad \quad | \text{while expr do u\_stmt}$



# Left recursion elimination

- A grammar is a left-recursive grammar, when there is a non-terminal  $A$  for which it is true that  $A \Rightarrow^+ A\alpha$  for a string  $\alpha$
- It is a problem for top-down parsing
- A simple solution for  $\beta\alpha^m$ :
  - $A \rightarrow A\alpha$
  - $A \rightarrow \beta$
  - $A \rightarrow \beta A'$
  - $A' \rightarrow \alpha A'$
  - $A' \rightarrow \Lambda$



# Removing left recursion from our grammar



1.  $E \rightarrow E + T$

2.  $E \rightarrow T$

3.  $T \rightarrow T * F$

4.  $T \rightarrow F$

5.  $F \rightarrow ( E )$

6.  $F \rightarrow \text{id}$

1.  $E \rightarrow TE'$

2.  $E' \rightarrow + TE'$

3.  $E' \rightarrow \Lambda$

4.  $T \rightarrow FT'$

5.  $T' \rightarrow * FT'$

6.  $T' \rightarrow \Lambda$

7.  $F \rightarrow ( E )$

8.  $F \rightarrow \text{id}$



# Left factoring

- It is not clear, which option we should choose
  - $A \rightarrow \alpha\beta_1$
  - $A \rightarrow \alpha\beta_2$
  - $A \rightarrow \alpha A'$
  - $A' \rightarrow \beta_1$
  - $A' \rightarrow \beta_2$

# Non-context-free language constructions



- $L_1 = \{ w c w \mid w = (a|b)^* \}$ 
  - Check, whether an identifier **w** is declared before using
- $L_2 = \{ a^n b^m c^n d^m \mid n \geq 1, m \geq 1 \}$ 
  - Check, whether number of parameters in function call confirms to the function declaration
- $L_3 = \{ a^n b^n c^n \mid n \geq 0 \}$ 
  - The problem of “underscoring” a word
    - **a** is a char, **b** is BS, **c** is underscore
  - $(abc)^*$  is a regular expression

# Operators FIRST and FOLLOW

## – definitions



- If  $\alpha$  is any string of grammar symbols, let  $\text{FIRST}(\alpha)$  be the set of terminals that begin the strings derived from  $\alpha$ . If  $\alpha$  can be derived to  $\Lambda$ , then  $\Lambda$  is also in  $\text{FIRST}(\alpha)$
- Define  $\text{FOLLOW}(A)$ , for nonterminal  $A$ , to be the set of terminals that can appear immediately to the right of  $A$  in some string, where exists a derivation of the form  $S \Rightarrow^* \alpha A a \beta$  for some  $\alpha$  and  $\beta$ . If  $A$  can be the rightmost symbol in some sentential form, then  $\$$  is in  $\text{FOLLOW}(A)$ .

# Construction of the FIRST operator



- Construction for a grammar symbol  $X$ 
  - If  $X$  is terminal, then  $\text{FIRST}(X) = \{X\}$
  - If  $X \rightarrow \Lambda$  is a production, then add  $\Lambda$  to  $\text{FIRST}(X)$
  - If  $X$  is nonterminal and  $X \rightarrow Y_1 Y_2 \dots Y_k$  is a production, then place  $a$  in  $\text{FIRST}(X)$ , if for some  $i$ ,  $a$  is in  $\text{FIRST}(Y_i)$  and  $\Lambda \in \text{FIRST}(Y_j) \ \forall j < i$ . If  $\Lambda \in \text{FIRST}(Y_j) \ \forall j$ , then add  $\Lambda$  to  $\text{FIRST}(X)$
- Construction for any string
  - The construction of the FIRST operator for a string  $X_1 X_2 \dots X_n$  is similar as for nonterminal.

# Construction of the FOLLOW operator



- Construction for a nonterminal  $A$ 
  - Place  $\$$  in  $\text{FOLLOW}(S)$ , where  $S$  is the start symbol of a grammar and  $\$$  is EOS
  - If there is a production  $A \rightarrow \alpha B \beta$ , then everything in  $\text{FIRST}(\beta)$  except for  $\Lambda$  is placed in  $\text{FOLLOW}(B)$
  - If there is a production  $A \rightarrow \alpha B$  or  $A \rightarrow \alpha B \beta$  where  $\Lambda \in \text{FIRST}(\beta)$ , then everything in  $\text{FOLLOW}(A)$  is in  $\text{FOLLOW}(B)$

# FIRST and FOLLOW – an example for our grammar



- $\text{FIRST}(E) = \{ (, \text{id} \}$
- $\text{FIRST}(T) = \{ (, \text{id} \}$
- $\text{FIRST}(F) = \{ (, \text{id} \}$
- $\text{FIRST}(E') = \{ +, \Lambda \}$
- $\text{FIRST}(T') = \{ *, \Lambda \}$
- $\text{FOLLOW}(E) = \{ ), \$ \}$
- $\text{FOLLOW}(E') = \{ ), \$ \}$
- $\text{FOLLOW}(T) = \{ +, ), \$ \}$
- $\text{FOLLOW}(T') = \{ +, ), \$ \}$
- $\text{FOLLOW}(F) = \{ +, *, ), \$ \}$



# Top-down parsing

- An attempt to find a leftmost derivation for an input string
- An attempt to construct a parse tree for the input starting from the root and creating the nodes of the tree in preorder
- Recursive-descent parsing
  - Recursive descent using procedures
- Nonrecursive predictive parsing
  - An automaton with an explicit stack
- Both solutions have a problem with left recursion in a grammar
- Many current parser generators use top-down parsing
  - ANTLR, CocoR – LL(1) grammars with conflict resolution using dynamic look-ahead expansion to LL(k)





# Recursive-descent parsing

- One procedure/function for each nonterminal of a grammar
- Each procedure does two things
  - It decides, which grammar production with given nonterminal on the left side will be used using look-ahead. A production with right side  $\alpha$  will be used, when the look-ahead is in  $FIRST(\alpha)$ . If there is a conflict for the look-ahead among some production right sides, the grammar is not suitable for recursive-descent parsing. A production with  $\Lambda$  on the right side will be used, if the look-ahead is not in  $FIRST$  of any right side.
  - Procedure code copies the right side of a production. Nonterminal means calling a procedure for this nonterminal. Terminal is compared with the look-ahead. If they are equal, a next terminal is read. If they are not equal, it is an error.

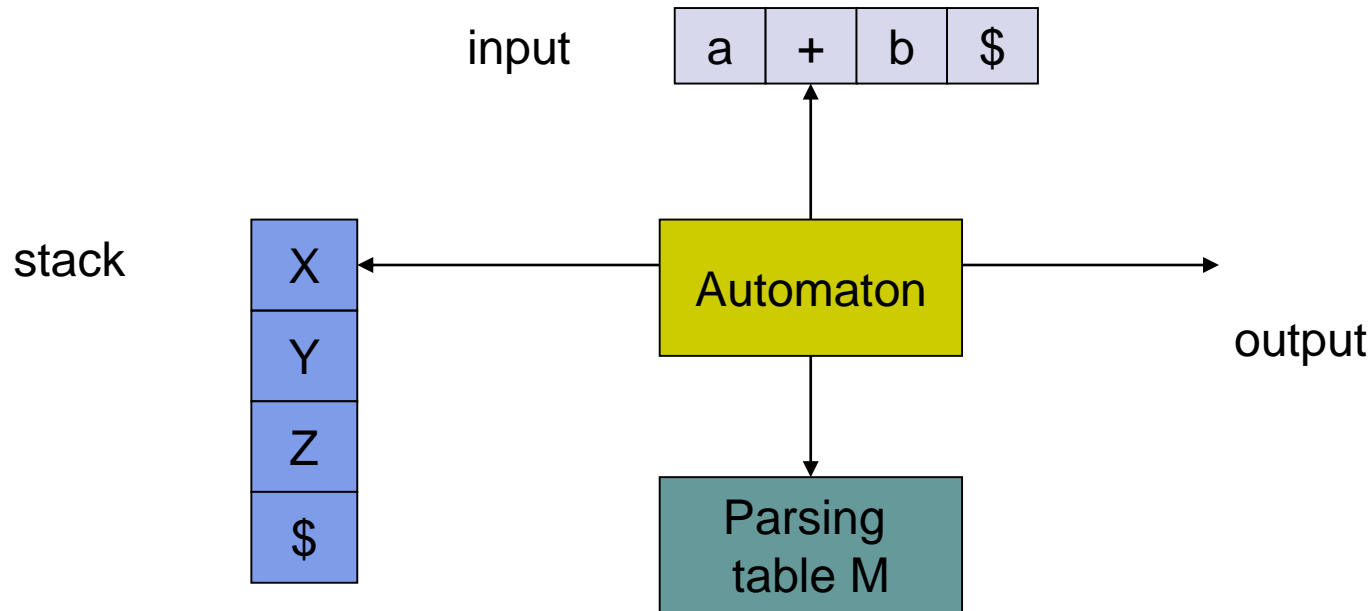
# Recursive-descent parsing – example for our grammar



```
void match(token t) {
    if(lookahead==t)
        lookahead = nexttoken();
    else error();
}
void E(void) {
    T(); Eap();
}
void Eap(void) {
    if(lookahead=='+') {
        match('+'); T(); Eap();
    }
}
void T(void) {
    F(); Tap();
}
```

```
void Tap(void) {
    if(lookahead=='*') {
        match('*'); F(); Tap();
    }
}
void F(void) {
    switch(lookahead) {
        case '(': match('('); E();
                    match(')'); break;
        case 'id':
                    match('id'); break;
        default:
                    error();
    }
}
```

# Nonrecursive predictive parsing



- Parsing table  $M[A, a]$ , where  $A$  is nonterminal and  $a$  is terminal
- The stack contains grammar symbols



# LL(1) automaton behavior

- Initial configuration
  - Input pointer points to the first terminal in the input string
  - The stack contains the start symbol of the grammar on top of \$
- In each step, the automaton decides, what to do, using a symbol  $X$  on top of the stack and a terminal  $a$ , pointed by input pointer
  - If  $X=a=\$$ , the parser halts, parsing finished successfully
  - If  $X=a\neq \$$ , the parser pops  $X$  from the stack and advances the input pointer to the next input symbol
  - If  $X\neq a$  and  $X\in T$ , the parser reports error
  - If  $X$  is a nonterminal, the parser uses entry  $M[X, a]$ . If this entry is a production, the parser replaces  $X$  on top of the stack by the right side (leftmost symbol on top of the stack). At the same time, the parser generates an output about using the production. If the entry is **error**, the parser informs about a syntax error.

# Construction of predictive parsing tables



- For each production  $A \rightarrow \alpha$  do following steps
  - For  $\forall a \in \text{FIRST}(\alpha)$  add  $A \rightarrow \alpha$  to  $M[A, a]$
  - If  $\Lambda \in \text{FIRST}(\alpha)$ , add  $A \rightarrow \alpha$  to  $M[A, b]$   
 $\forall b \in \text{FOLLOW}(A)$ . Moreover, if  $\$ \in \text{FOLLOW}(A)$ ,  
add  $A \rightarrow \alpha$  to  $M[A, \$]$
- Mark each empty entry in  $M$  as **error**

# Example of table construction for our grammar



	id	+	*	(	)	\$
E	$E \rightarrow TE'$			$E \rightarrow TE'$		
E'		$E' \rightarrow +TE'$			$E' \rightarrow \Lambda$	$E' \rightarrow \Lambda$
T	$T \rightarrow FT'$			$T \rightarrow FT'$		
T'		$T' \rightarrow \Lambda$	$T' \rightarrow *FT'$		$T' \rightarrow \Lambda$	$T' \rightarrow \Lambda$
F	$F \rightarrow \text{id}$			$F \rightarrow (E)$		

# Example of parser behavior for our grammar



Stack	Input	Output
\$E	id+id*id\$	
\$E'T	id+id*id\$	$E \rightarrow TE'$
\$E'T'F	id+id*id\$	$T \rightarrow FT'$
\$E'T'id	id+id*id\$	$F \rightarrow id$
\$E'T'	+id*id\$	
\$E'	+id*id\$	$T' \rightarrow \Lambda$
\$E'T+	+id*id\$	$E' \rightarrow +TE'$
\$E'T	id*id\$	
\$E'T'F	id*id\$	$T \rightarrow FT'$

Stack	Input	Output
\$E'T'id	id*id\$	$F \rightarrow id$
\$E'T'	*id\$	
\$E'T'F*	*id\$	$T' \rightarrow *FT'$
\$E'T'F	id\$	
\$E'T'id	id\$	$F \rightarrow id$
\$E'T'	\$	
\$E'	\$	$T' \rightarrow \Lambda$
\$	\$	$E' \rightarrow \Lambda$



# LL(1) grammar

- Context-free grammar  $G=(T,N,S,P)$  is a LL(1) grammar, if and only if whenever  $A \rightarrow \alpha$ ,  $A \rightarrow \beta \in P$  are two distinct ( $\alpha \neq \beta$ ) productions of  $G$  and we have any left sentential forms  $uA\gamma$ ,  $vA\delta$ , where  $u,v \in T^*$  and  $\gamma, \delta \in (T \cup N)^*$ , the following condition holds:
  - $\text{FIRST}(\alpha\gamma) \cap \text{FIRST}(\beta\delta) = \emptyset$
- Simplified detection: no ambiguous or left-recursive grammar can be LL(1)





# Grammar terminology

- $PXY(k)$
- $X$  – direction of the input reading
  - In our case always  $L$ , i.e. from left to right
- $Y$  – kind of derivation
  - $L$  – left derivation
  - $R$  – right derivation
- $P$  – prefix
  - Subtle division of some grammar classes
- $k$  – look-ahead
  - An integer, usually 1, can be 0 or more generally  $k$
- Examples
  - $LL(1)$ ,  $LR(0)$ ,  $LR(1)$ ,  $LL(k)$ ,  $SLR(1)$ ,  $LALR(1)$

# Expanding definition of FIRST and FOLLOW on k



- If  $\alpha$  is a string from grammar symbols, then  $\text{FIRST}_k(\alpha)$  is a set of terminal words with maximal length  $k$ , which are on the beginning of at least one string derived from  $\alpha$ . If  $\alpha$  can be derived on  $\Lambda$ , then  $\Lambda$  is in  $\text{FIRST}_k(\alpha)$ .
- $\text{FOLLOW}_k(A)$  for nonterminal  $A$  is a set of terminal words with maximal length  $k$ , which can be on the right side of  $A$  in any string derived from the start nonterminal ( $S \Rightarrow^* \alpha A u \beta$  for some  $\alpha$  and  $\beta$ ). If  $A$  is the right-most symbol in any sentential form, then  $\$$  is in  $\text{FOLLOW}_k(A)$ .



# LL(k) grammar

- Context-free grammar  $G=(T,N,S,P)$  is a strong LL(k) grammar for  $k \geq 1$ , if and only if whenever  $A \rightarrow \alpha, A \rightarrow \beta \in P$  are two distinct ( $\alpha \neq \beta$ ) productions and we have any left sentential forms  $uA\gamma, vA\delta$ , where  $u, v \in T^*$  and  $\gamma, \delta \in (T \cup N)^*$ , the following condition holds:
  - $\text{FIRST}_k(\alpha\gamma) \cap \text{FIRST}_k(\beta\delta) = \emptyset$ .
- LL(k) (not strong)
  - $u=v, \gamma=\delta$

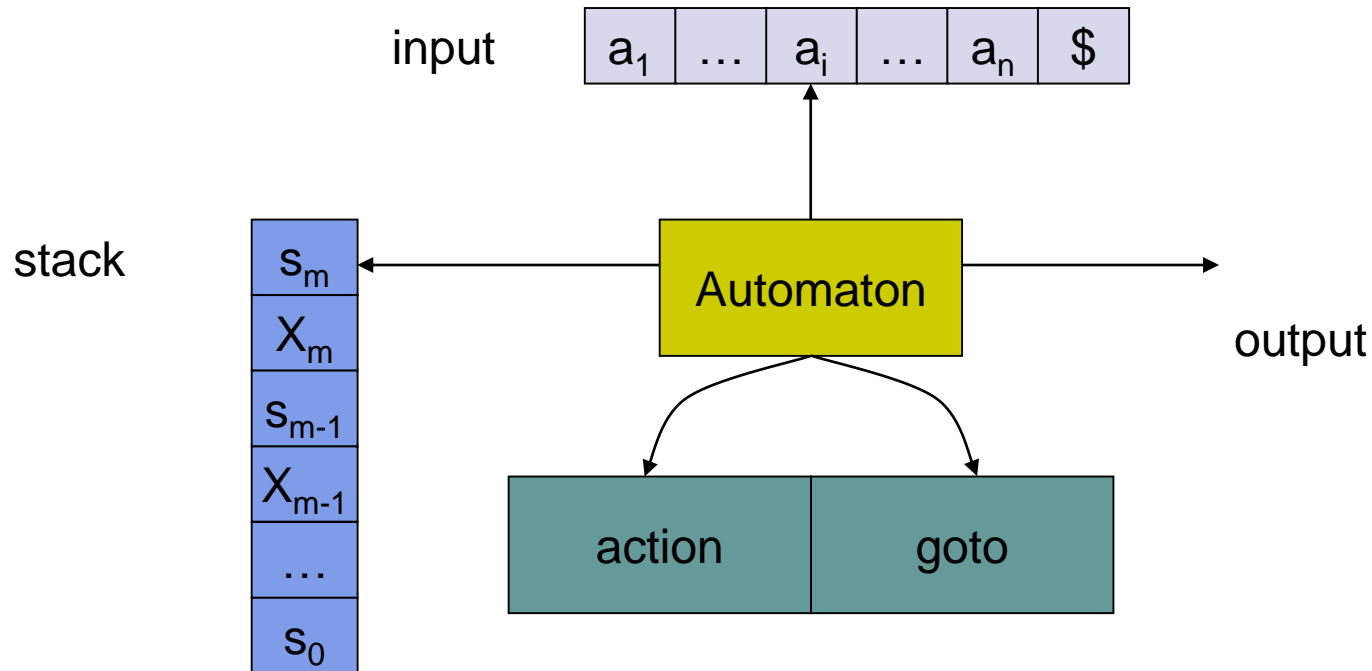


# Bottom-up analysis

- Attempts to find in reverse the rightmost derivation for an input string
- Attempts to construct a parse tree for an input string beginning at the leaves and working up towards the root
- Replace a substring corresponding to a right side of a production by a nonterminal from the left side of the production in each reduce step
- Used in parser generators
  - Bison – LALR(1), GLR(1)
- Advantages against LL(1) parsers
  - It can be implemented with the same efficiency as top-down parsing
  - The class of decidable languages LR(1) is a proper superset of LL(1)
- SLR(1), LR(1), LALR(1)



# LR parser automaton



- $s_i$  are states
  - A state on the top of the stack is the current state of the automaton
- $x_i$  are grammar symbols



# LR(1) automaton behavior

- Initial configuration
  - Input pointer points to the first terminal in the input string
  - Initial state  $s_0$  is on the stack
- In each step address table action[ $s_m, a_i$ ] using  $s_m$  and  $a_i$ 
  - Shift  $s$ , where  $s$  is a new state
    - It shifts the input tape by 1 terminal and add  $a_i$  and  $s$  on the top of the stack
  - Reduce using production  $A \rightarrow \alpha$ 
    - Remove  $r=|\alpha|$  pairs  $(s_k, X_k)$  from the top of the stack, add  $A$  on the top of the stack and then goto[ $s_{m-r}, A$ ] ( $s_{m-r}$  is a state on the top of the stack after erasing pairs)
    - Generate an output
  - Accept
    - The input string is accepted
    - Generate an output
  - Error
    - The input string is not in the input language

# LR automaton tables for our grammar



state	action						goto		
	id	+	*	(	)	\$	E	T	F
0	s5			s4			1	2	3
1		s6				acc			
2		r2	s7		r2	r2			
3		r4	r4		r4	r4			
4	s5			s4			8	2	3
5		r6	r6		r6	r6			
6	s5			s4				9	3
7	s5			s4					10
8		s6			s11				
9		r1	s7		r1	r1			
10		r3	r3		r3	r3			
11		r5	r5		r5	r5			



# Example of LR parser behavior

Stack	Input	Action
0	<b>id+id*id\$</b>	s5
0 <b>id</b> 5	<b>+id*id\$</b>	r6: F→ <b>id</b>
0 F 3	<b>+id*id\$</b>	r4: T→F
0 T 2	<b>+id*id\$</b>	r2: E→T
0 E 1	<b>+id*id\$</b>	s6
0 E 1 <b>+</b> 6	<b>id*id\$</b>	s5
0 E 1 <b>+</b> 6 <b>id</b> 5	<b>*id\$</b>	r6: F→ <b>id</b>
0 E 1 <b>+</b> 6 F 3	<b>*id\$</b>	r4: T→F
0 E 1 <b>+</b> 6 T 9	<b>*id\$</b>	s7
0 E 1 <b>+</b> 6 T 9 <b>*</b> 7	<b>id\$</b>	s5
0 E 1 <b>+</b> 6 T 9 <b>*</b> 7 <b>id</b> 5	<b>\$</b>	r6: F→ <b>id</b>
0 E 1 <b>+</b> 6 T 9 <b>*</b> 7 F 10	<b>\$</b>	r3: T→T * F
0 E 1 <b>+</b> 6 T 9	<b>\$</b>	r1: E→E + T
0 E 1	<b>\$</b>	acc





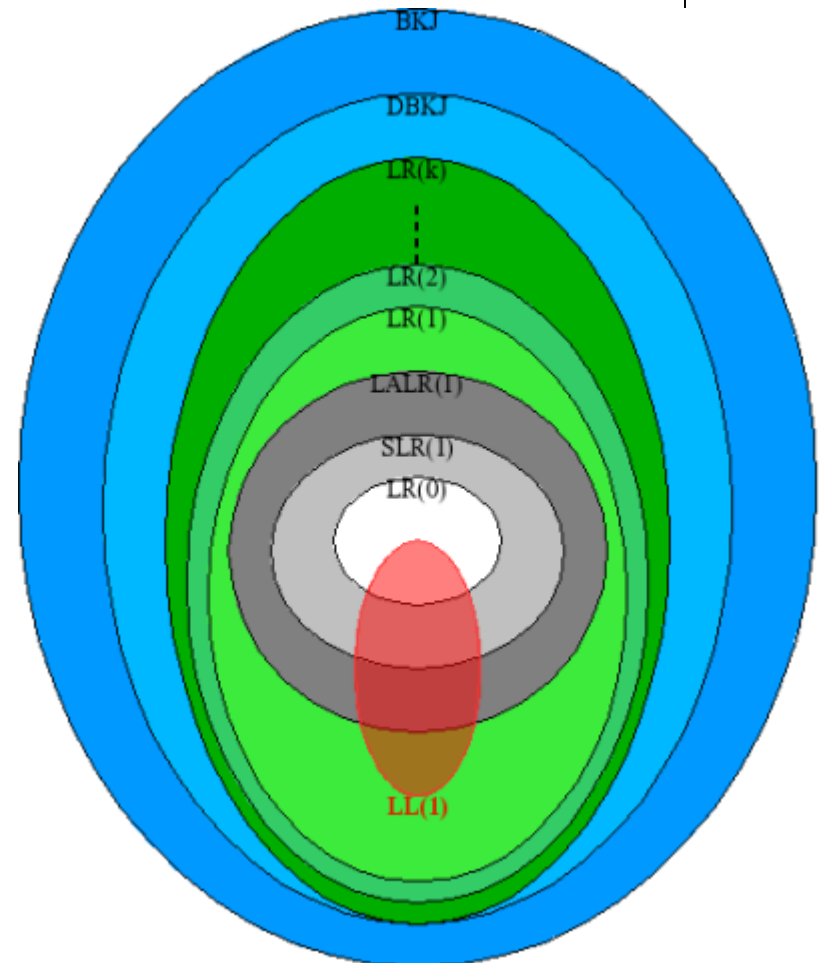
# LR(k) grammar

- Context-free grammar  $G=(T,N,S,P)$  is LR(k) grammar for  $k \geq 1$ , if and only if whenever  $A \rightarrow \alpha, A \rightarrow \beta \in P$  are two distinct ( $\alpha \neq \beta$ ) productions of  $G$  and we have any two right sentential forms  $\gamma Au, \delta Av$ , where  $u, v \in T^*$  and  $\gamma, \delta \in (T \cup N)^*$ , the following condition holds:
  - $\text{FIRST}_k(u) \cap \text{FIRST}_k(v) = \emptyset$

# Grammars (languages) strength



- Union of all  $LR(k)$  are deterministic context-free languages (DBKJ) and it is a proper subset of all context-free languages (BKJ)





# Grammar augmentation

- Augmentation of a grammar  $G=(T,N,S,P)$  is a grammar  $G'=(T,N',S',P')$ , where  $N'=N\cup\{S'\}$  and  $P'=P\cup\{S'\rightarrow S\}$
- The augmentation is not necessary whenever  $S$  is on the left side of one production and it isn't on any right side of grammar productions
- It helps recognize the end of parsing
- For our grammar:
  - $S'\rightarrow E$



# LR(0) items

- LR(0) item of a grammar  $G$  is a production with a special symbol dot on the right side
  - Special symbol is a valid symbol for comparison of two LR(0) items of a same production. LR(0) items of the same production are different, whenever the dot is on different position. Moreover, the dot is not a grammar symbol
- An example for production  $E \rightarrow E + T$ :

$E \rightarrow \blacklozenge E + T$	$E \rightarrow E + \blacklozenge T$
$E \rightarrow E \blacklozenge + T$	$E \rightarrow E + T \blacklozenge$



# The closure operation

- If  $I$  is a set of LR(0) items for a grammar  $G$ , then  $\text{CLOSURE}(I)$  is a set of LR(0) items constructed from  $I$  by following rules:
  - Add  $I$  to the  $\text{CLOSURE}(I)$
  - $\forall A \rightarrow \alpha \blacklozenge B \beta \in \text{CLOSURE}(I)$ , where  $B \in N$ , add  $\forall B \rightarrow \gamma \in P$  to  $\text{CLOSURE}(I)$  LR(0) item  $B \rightarrow \blacklozenge \gamma$ , if it is not already there. Apply this rule until no more new LR(0) items can be added to  $\text{CLOSURE}(I)$

# Example of closure for our grammar



- $I = \{S' \rightarrow \blacklozenge E\}$
- $\text{CLOSURE}(I) =$ 
  - $S' \rightarrow \blacklozenge E$
  - $E \rightarrow \blacklozenge E + T$
  - $E \rightarrow \blacklozenge T$
  - $T \rightarrow \blacklozenge T * F$
  - $T \rightarrow \blacklozenge F$
  - $F \rightarrow \blacklozenge ( E )$
  - $F \rightarrow \blacklozenge \text{id}$



# GOTO operation

- GOTO( $I$ ,  $X$ ) operation for a set  $I$  of LR(0) items and a grammar symbol  $X$  is defined to be the closure of the set of all LR(0) items  $A \rightarrow \alpha X \diamond \beta$  such that  $A \rightarrow \alpha \diamond X \beta \in I$

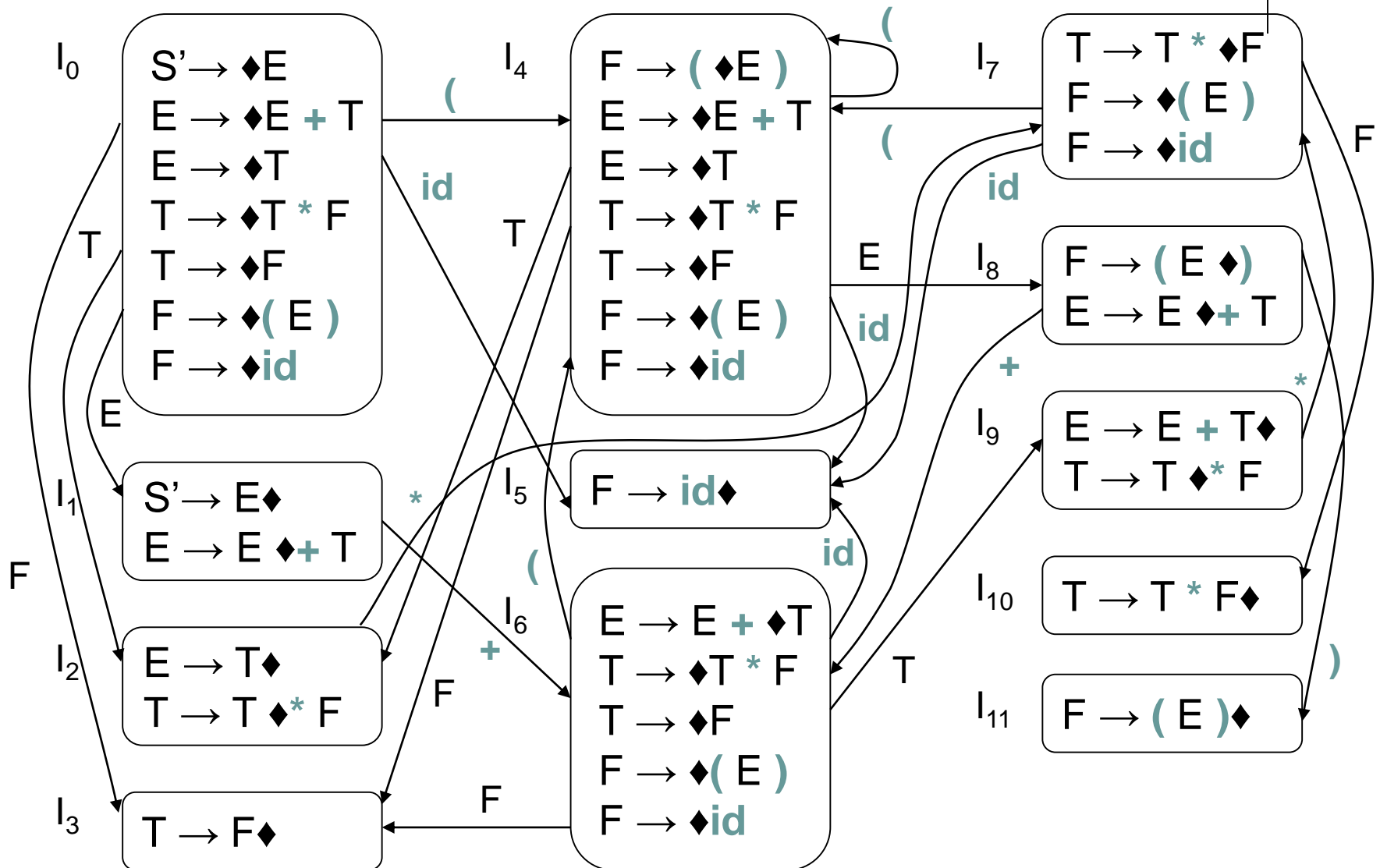
# Construction of canonical collection of sets of LR(0) items



- We have an augmented grammar  $G'=(T,N',S',P')$
- Construction of canonical collection  $C$  of sets of LR(0) items:
  - We start with  $C=\{ \text{CLOSURE}(\{S' \rightarrow \blacklozenge S\}) \}$
  - $\forall I \in C$  and  $\forall X \in T \cup N'$  such as  $\text{GOTO}(I, X) \notin C \wedge \text{GOTO}(I, X) \neq \emptyset$ , add  $\text{GOTO}(I, X)$  to  $C$ . Repeat this step, until something new is added to  $C$ .



# Construction of canonical collection for our grammar





# Valid items

- LR(0) item  $A \rightarrow \beta_1 \blacklozenge \beta_2$  is a valid item for a viable prefix  $\alpha\beta_1$ , if there is a rightmost derivation  $S' \Rightarrow^+ \alpha A w \Rightarrow \alpha\beta_1\beta_2 w$
- It is a great hint for a parser. It helps to decide, if the parser should make a shift or a reduction, if  $\alpha\beta_1$  is on top of the stack
- Basic LR parsing theorem: A set of valid items for a viable prefix  $\gamma$  is exactly a set of items reachable from the initial state through the prefix  $\gamma$  by deterministic finite automaton constructed from canonical collection with GOTO transitions.

# SLR(1) automaton construction



- We have an augmented grammar  $G'$ . Tables of SLR(1) automaton are constructed by following algorithm
  - Construct a canonical collection  $C$  of sets of LR(0) items
  - State  $i$  is constructed from  $I_i$ . The parsing actions for state  $i$  are determined as follows
    - $A \rightarrow \alpha \diamond a \beta \in I_i, a \in T \wedge \text{GOTO}(I_i, a) = I_j$ , then  $\text{action}[i, a] = \text{shift } j$
    - $A \rightarrow \alpha \diamond \in I_i$ , then  $\forall a \in \text{FOLLOW}(A) \wedge A \neq S'$  is  $\text{action}[i, a] = \text{reduce } A \rightarrow \alpha$
    - $S' \rightarrow S \diamond \in I_i$ , then  $\text{action}[i, \$] = \text{accept}$
  - If there is a conflict in the previous step, the grammar is not a SLR(1) grammar and the automaton cannot be constructed
  - Table goto is indexed by state  $i$  and  $A \in N'$ : whenever  $\text{GOTO}(I_i, A) = I_j$ , then  $\text{goto}[i, A] = j$
  - All empty cells are filled by error instruction
  - The initial state of the parser is the state, which contains LR(0) item  $S' \rightarrow \diamond S$



# Full LR(1) automata

- $\text{action}[i, a]$  is set to reduction  $A \rightarrow \alpha$ , when  $A \rightarrow \alpha \diamond \in I_i$ ,  $\forall a \in \text{FOLLOW}(A)$  for a state  $i$  during SLR(1) construction
- In some situation, when  $i$  is on top of the stack, the viable prefix  $\beta\alpha$  is in form, where  $\beta A$  cannot be followed by a terminal  $a$  in any right sentential form. Therefore reduction  $A \rightarrow \alpha$  is for lookahead  $a$  invalid.
- Solution: add more information to states, so we can avoid invalid reductions.



# LR(1) items

- The added information is stored as an additional terminal for each LR(0) item. Such item has a form  $[A \rightarrow \alpha \blacklozenge \beta, a]$ , where  $A \rightarrow \alpha \beta \in P$ ,  $a \in T$ , and we call it LR(1) item. The terminal **a** is called lookahead.
  - The lookahead has no meaning for  $A \rightarrow \alpha \blacklozenge \beta$ , where  $\beta \neq \Lambda$
  - Reduction  $A \rightarrow \alpha$  is set only when  $[A \rightarrow \alpha \blacklozenge, a] \in I_i$  for current state  $i$  and a terminal **a** on the input
  - A set of terminals created from lookaheads of LR(1) items  $\subseteq \text{FOLLOW}(A)$
- LR(1) item  $[A \rightarrow \alpha \blacklozenge \beta, a]$  is valid for viable prefix  $\gamma$ , whenever  $\exists$  right derivation  $S \Rightarrow^+ \delta A w \Rightarrow \delta \alpha \beta w$ , where
  - $\gamma = \delta \alpha$
  - Either **a** is the first symbol of  $w$  or  $w = \Lambda$  and **a** is \$



# Closure for LR(1) items

- We have a set of LR(1) items  $I$  for a grammar  $G$ . We define  $CLOSURE_1(I)$  as a set of LR(1) items constructed from  $I$  by following procedure:
  - Add set  $I$  to  $CLOSURE_1(I)$
  - $\forall [A \rightarrow \alpha \blacklozenge B \beta, a] \in CLOSURE_1(I)$ , where  $B \in N$ , add LR(1) item  $[B \rightarrow \blacklozenge \gamma, b] \forall B \rightarrow \gamma \in P$  and  $\forall b \in FIRST(\beta a)$  to  $CLOSURE_1(I)$ , if it isn't there already. Repeat this step, until something is added to  $CLOSURE_1(I)$ .

# GOTO operation for LR(1) items



- We define  $GOTO_1(I, X)$  operation for a set  $I$  of LR(1) items and a grammar symbol  $X$  as a  $CLOSURE_1$  of a set of all items  $[A \rightarrow \alpha X \blacklozenge \beta, a]$  where  $[A \rightarrow \alpha \blacklozenge X \beta, a] \in I$

# Construction of canonical collection of sets of LR(1) items



- We have an augmented grammar  $G'=(T,N',S',P')$
- Construction of canonical collection  $C$  of LR(1) items:
  - We start with  $C=\{ \text{CLOSURE}_1(\{[S' \rightarrow \diamond S, \$]\}) \}$
  - Add  $\text{GOTO}_1(I, X)$  to  $C \forall I \in C$  and  $\forall X \in T \cup N'$ , where  $\text{GOTO}_1(I, X) \notin C \wedge \text{GOTO}_1(I, X) \neq \emptyset$ . Repeat this step, until something new is added to  $C$ .



# Example of LR(1) grammar, which is not SLR(1)



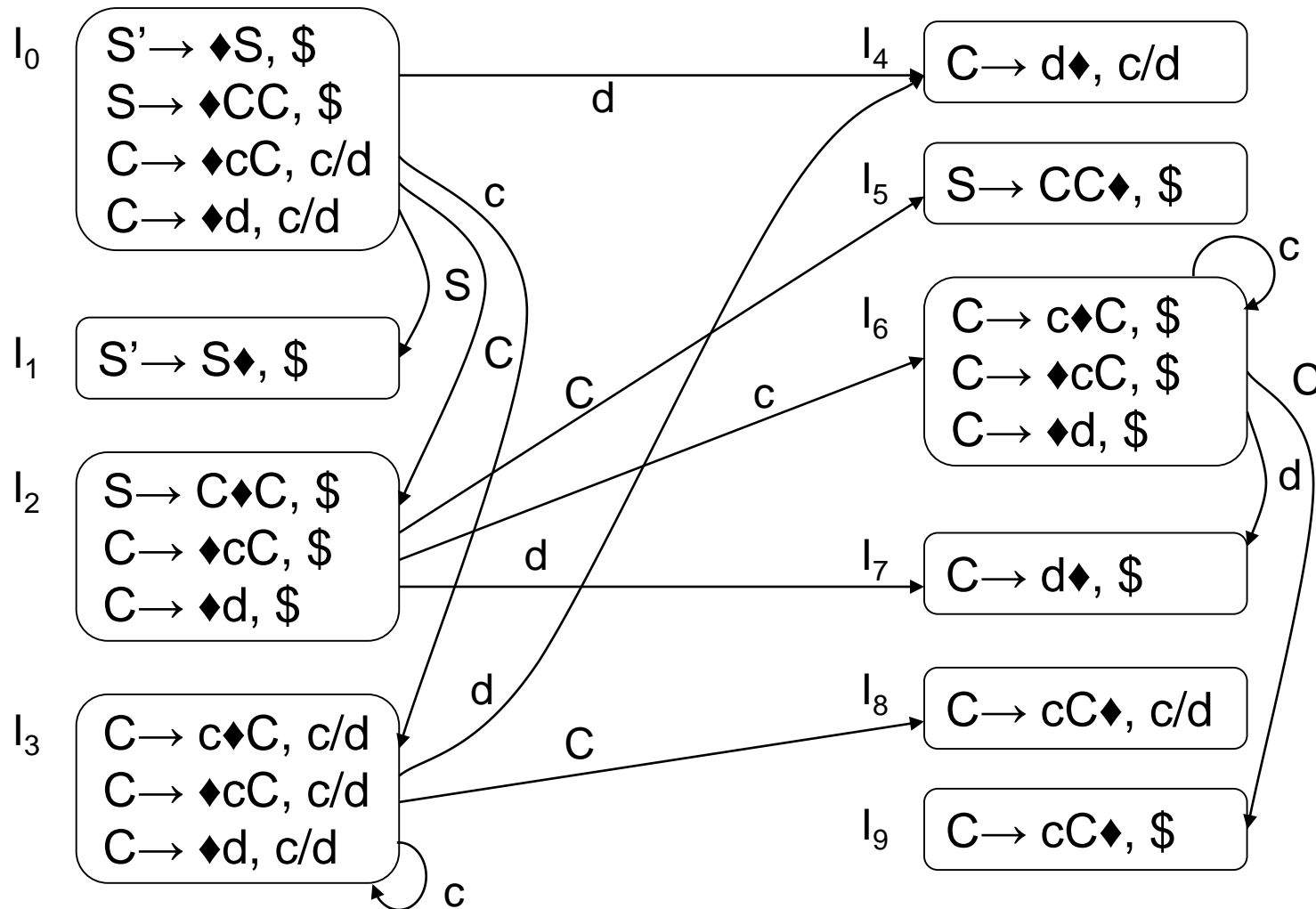
- $S' \rightarrow S$
- $S \rightarrow CC$
- $C \rightarrow cC$
- $C \rightarrow d$

# Example of closure construction for LR(1) items



- $I = \{[S' \rightarrow \blacklozenge S, \$]\}$
- $\text{CLOSURE}_1(I) =$ 
  - $S' \rightarrow \blacklozenge S, \$ \quad \beta = \Lambda, \text{FIRST}(\beta \$) = \text{FIRST}(\$) = \{\$\}$
  - $S \rightarrow \blacklozenge CC, \$ \quad \beta = C, \text{FIRST}(C \$) = \{c, d\}$
  - $C \rightarrow \blacklozenge cC, c/d$
  - $C \rightarrow \blacklozenge d, c/d$

# Example of construction of canonical collection of LR(1) items





# LR(1) parser construction

- We have an augmented grammar  $G'$ . LR(1) automaton tables are constructed by following algorithm
  - Construct a canonical collection  $C$  of sets of LR(1) items
  - State  $i$  is constructed from  $I_i$ . The parsing actions for state  $i$  are determined as follows
    - $[A \rightarrow \alpha \blacklozenge a \beta, b] \in I_i, a \in T \wedge \text{GOTO1}(I_i, a) = I_j$ , then  $\text{action}[i, a] = \text{shift } j$
    - $[A \rightarrow \alpha \blacklozenge, a] \in I_i \wedge A \neq S'$ , then  $\text{action}[i, a] = \text{reduce } A \rightarrow \alpha$
    - $[S' \rightarrow S \blacklozenge, \$] \in I_i$ , then  $\text{action}[i, \$] = \text{accept}$
  - If there is a conflict in the previous step, the grammar is not a LR(1) grammar and the automaton cannot be constructed
  - Table goto is indexed by state  $i$  and  $A \in N'$ : whenever  $\text{GOTO1}(I_i, A) = I_j$ , then  $\text{goto}[i, A] = j$
  - All empty cells are filled by error instruction
  - The initial state of the parser is the state, which contains LR(1) item  $[S' \rightarrow \blacklozenge S, \$]$



# LALR

- LALR=LookAhead-LR
- Often used in practice
  - Bison
  - Most common programming languages can be expressed by an LALR grammar
  - Parser tables are considerably smaller than LR(1) tables
- SLR and LALR parsers have the same number of states, LR parsers have greater number of states
  - Common languages have hundreds of states
  - LR(1) parsers have thousands of states for the same grammar



# How to make smaller tables?

- Idea: merge sets with the same core into one set including GOTO1 merge
  - Core: a set of LR(0) items (no lookahead)
  - Merge cannot produce shift/reduce conflict
    - Suppose in the union there is a conflict on lookahead **a** for LR(1) items  $[A \rightarrow \alpha \blacklozenge, a]$  and  $[B \rightarrow \beta \blacklozenge a \gamma, b]$
    - Cores are same, therefore in the set with  $[A \rightarrow \alpha \blacklozenge, a]$  must be  $[B \rightarrow \beta \blacklozenge a \gamma, c]$  as well for some **c**. There was already a shift/reduce conflict before merge
  - Merge can produce reduce/reduce conflict

# Easy LALR(1) table construction



- We have an augmented grammar  $G'$ . LALR(1) automaton tables are constructed by following algorithm
  - Construct a canonical collection  $C$  of sets of LR(1) items
  - For each core in collection  $C$ , find all sets having that core, and replace these sets by their union
  - Let  $C' = \{ J_0, J_1, \dots, J_m \}$  be the resulting collection of LR(1) items
  - Table action is constructed for  $C'$  in the same manner as for full LR(1) parser
  - If there is a conflict, the grammar is not LALR(1) grammar
  - If  $J \in C'$  is the union of sets of LR(1) items  $I_i$  ( $J = I_1 \cup I_2 \cup \dots \cup I_k$ ), then cores  $GOTO1(I_1, X), \dots, GOTO1(I_k, X)$  are the same, since  $I_1, \dots, I_k$  all have the same core. Let  $K$  be the union of all sets of items having the same core as  $goto(I_1, X)$ . Then  $GOTO1(J, X) = K$
- Important disadvantage – we need to construct full LR(1)

# Compiler principles

---

Syntax analysis

Jakub Yaghob

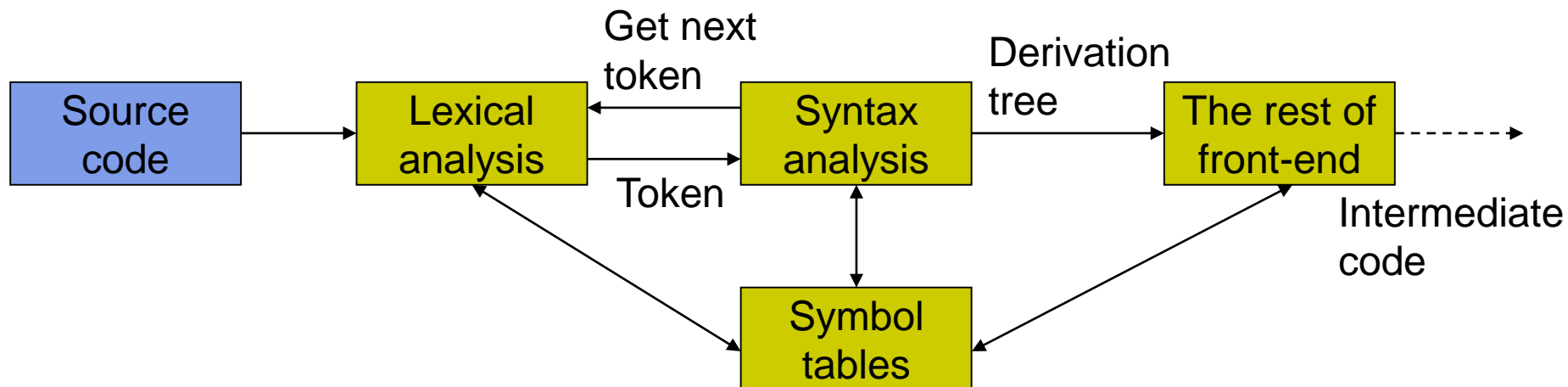






# Syntax analysis

- The main task
  - Decide, whether an input word is a word from an input language
- Other important tasks
  - Syntax-directed translation is the main loop of the compiler
  - Build the derivation tree
- Automaton type
  - We are talking about (deterministic) context-free grammars, therefore we are using (deterministic) pushdown automata





# Our grammar

1.  $E \rightarrow E + T$
2.  $E \rightarrow T$
3.  $T \rightarrow T * F$
4.  $T \rightarrow F$
5.  $F \rightarrow ( E )$
6.  $F \rightarrow \text{id}$

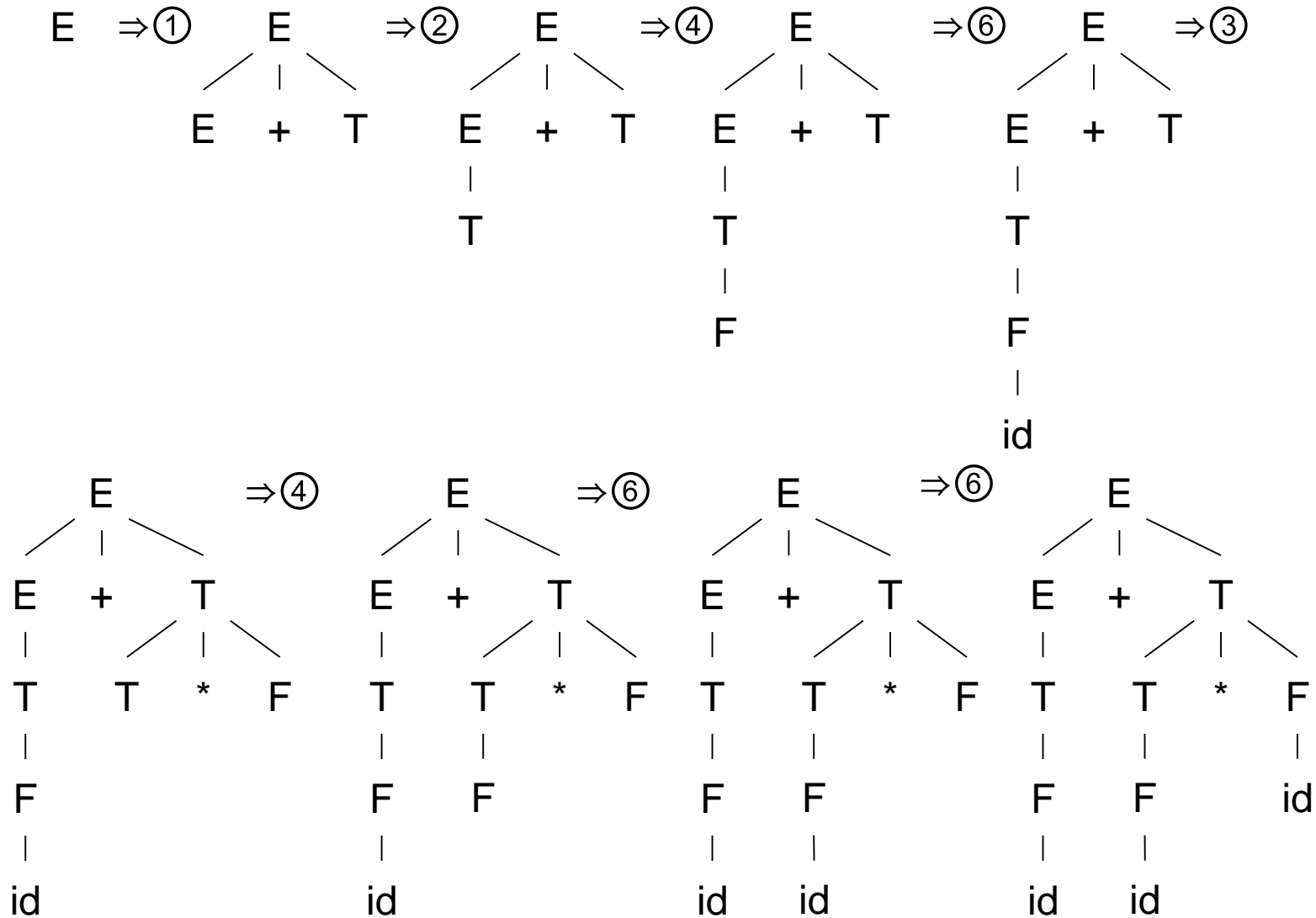


# Derivation (parse, syntax) tree

- Graphical representation of derivations using trees
  - Vertices are both non-terminals and terminals
  - Edges from inner vertex representing a non-terminal on the left side of a production rule to all symbols from the right side of a production rule
- $E \Rightarrow \textcircled{1} E+T \Rightarrow \textcircled{2} T+T \Rightarrow \textcircled{4} F+T \Rightarrow \textcircled{6} \text{id}+T \Rightarrow \textcircled{3} \text{id}+T^*F \Rightarrow \textcircled{4} \text{id}+F^*F \Rightarrow \textcircled{6} \text{id}+\text{id}^*F \Rightarrow \textcircled{6} \text{id}+\text{id}^*\text{id}$



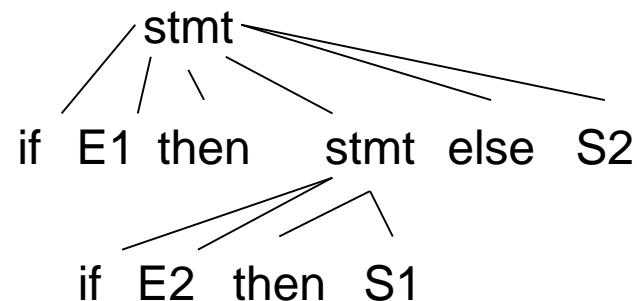
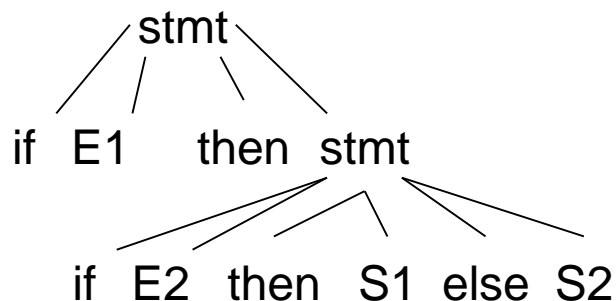
# Example





# Ambiguous grammar

- We can construct distinct derivation trees for the same input word
- Real-life example (dangling else):
  - $\text{stmt} \rightarrow$  **if** expr **then** stmt  
| **if** expr **then** stmt **else** stmt  
| **while** expr **do** stmt  
| **goto** num
  - Input word: **if**  $E_1$  **then** **if**  $E_2$  **then**  $S_1$  **else**  $S_2$





# Disambiguation

- Clarify, which tree is the right one
- In our case: **else** pairs with nearest “free” **if** (without **else**)
- Idea: “paired” statement is always between **if** and **else**
  - $\text{stmt} \rightarrow \text{m\_stmt}$
  - $\quad \quad \quad | \text{u\_stmt}$
  - $\text{m\_stmt} \rightarrow \text{if expr then m\_stmt else m\_stmt}$ 
    - $\quad \quad \quad | \text{while expr do m\_stmt}$
    - $\quad \quad \quad | \text{goto num}$
  - $\text{u\_stmt} \rightarrow \text{if expr then stmt}$ 
    - $\quad \quad \quad | \text{if expr then m\_stmt else u\_stmt}$
    - $\quad \quad \quad | \text{while expr do u\_stmt}$



# Left recursion elimination

- A grammar is a left-recursive grammar, when there is a non-terminal  $A$  for which it is true that  $A \Rightarrow^+ A\alpha$  for a string  $\alpha$
- It is a problem for top-down parsing
- A simple solution for  $\beta\alpha^m$ :
  - $A \rightarrow A\alpha$
  - $A \rightarrow \beta$
  - $A \rightarrow \beta A'$
  - $A' \rightarrow \alpha A'$
  - $A' \rightarrow \Lambda$

# Removing left recursion from our grammar



1.  $E \rightarrow E + T$

2.  $E \rightarrow T$

3.  $T \rightarrow T * F$

4.  $T \rightarrow F$

5.  $F \rightarrow ( E )$

6.  $F \rightarrow id$

1.  $E \rightarrow TE'$

2.  $E' \rightarrow + TE'$

3.  $E' \rightarrow \Lambda$

4.  $T \rightarrow FT'$

5.  $T' \rightarrow * FT'$

6.  $T' \rightarrow \Lambda$

7.  $F \rightarrow ( E )$

8.  $F \rightarrow id$





# Left factoring

- It is not clear, which option we should choose
  - $A \rightarrow \alpha\beta_1$
  - $A \rightarrow \alpha\beta_2$
  - $A \rightarrow \alpha A'$
  - $A' \rightarrow \beta_1$
  - $A' \rightarrow \beta_2$

# Non-context-free language constructions



- $L_1 = \{ w c w \mid w = (a|b)^* \}$ 
  - Check, whether an identifier **w** is declared before using
- $L_2 = \{ a^n b^m c^n d^m \mid n \geq 1, m \geq 1 \}$ 
  - Check, whether number of parameters in function call confirms to the function declaration
- $L_3 = \{ a^n b^n c^n \mid n \geq 0 \}$ 
  - The problem of “underscoring” a word
    - **a** is a char, **b** is BS, **c** is underscore
  - $(abc)^*$  is a regular expression

# Operators FIRST and FOLLOW

## – definitions



- If  $\alpha$  is any string of grammar symbols, let  $\text{FIRST}(\alpha)$  be the set of terminals that begin the strings derived from  $\alpha$ . If  $\alpha$  can be derived to  $\Lambda$ , then  $\Lambda$  is also in  $\text{FIRST}(\alpha)$
- Define  $\text{FOLLOW}(A)$ , for nonterminal  $A$ , to be the set of terminals that can appear immediately to the right of  $A$  in some string, where exists a derivation of the form  $S \Rightarrow^* \alpha A a \beta$  for some  $\alpha$  and  $\beta$ . If  $A$  can be the rightmost symbol in some sentential form, then  $\$$  is in  $\text{FOLLOW}(A)$ .

# Construction of the FIRST operator



- Construction for a grammar symbol  $X$ 
  - If  $X$  is terminal, then  $\text{FIRST}(X) = \{X\}$
  - If  $X \rightarrow \Lambda$  is a production, then add  $\Lambda$  to  $\text{FIRST}(X)$
  - If  $X$  is nonterminal and  $X \rightarrow Y_1 Y_2 \dots Y_k$  is a production, then place  $a$  in  $\text{FIRST}(X)$ , if for some  $i$ ,  $a$  is in  $\text{FIRST}(Y_i)$  and  $\Lambda \in \text{FIRST}(Y_j) \forall j < i$ . If  $\Lambda \in \text{FIRST}(Y_j) \forall j$ , then add  $\Lambda$  to  $\text{FIRST}(X)$
- Construction for any string
  - The construction of the FIRST operator for a string  $X_1 X_2 \dots X_n$  is similar as for nonterminal.

# Construction of the FOLLOW operator



- Construction for a nonterminal  $A$ 
  - Place  $\$$  in  $\text{FOLLOW}(S)$ , where  $S$  is the start symbol of a grammar and  $\$$  is EOS
  - If there is a production  $A \rightarrow \alpha B \beta$ , then everything in  $\text{FIRST}(\beta)$  except for  $\Lambda$  is placed in  $\text{FOLLOW}(B)$
  - If there is a production  $A \rightarrow \alpha B$  or  $A \rightarrow \alpha B \beta$  where  $\Lambda \in \text{FIRST}(\beta)$ , then everything in  $\text{FOLLOW}(A)$  is in  $\text{FOLLOW}(B)$

# FIRST and FOLLOW – an example for our grammar



- $\text{FIRST}(E) = \{ (, \text{id} \}$
- $\text{FIRST}(T) = \{ (, \text{id} \}$
- $\text{FIRST}(F) = \{ (, \text{id} \}$
- $\text{FIRST}(E') = \{ +, \Lambda \}$
- $\text{FIRST}(T') = \{ *, \Lambda \}$
- $\text{FOLLOW}(E) = \{ ), \$ \}$
- $\text{FOLLOW}(E') = \{ ), \$ \}$
- $\text{FOLLOW}(T) = \{ +, ), \$ \}$
- $\text{FOLLOW}(T') = \{ +, ), \$ \}$
- $\text{FOLLOW}(F) = \{ +, *, ), \$ \}$



# Top-down parsing

- An attempt to find a leftmost derivation for an input string
- An attempt to construct a parse tree for the input starting from the root and creating the nodes of the tree in preorder
- Recursive-descent parsing
  - Recursive descent using procedures
- Nonrecursive predictive parsing
  - An automaton with an explicit stack
- Both solutions have a problem with left recursion in a grammar
- Many current parser generators use top-down parsing
  - ANTLR, CocoR – LL(1) grammars with conflict resolution using dynamic look-ahead expansion to LL(k)



# Recursive-descent parsing

- One procedure/function for each nonterminal of a grammar
- Each procedure does two things
  - It decides, which grammar production with given nonterminal on the left side will be used using look-ahead. A production with right side  $\alpha$  will be used, when the look-ahead is in  $FIRST(\alpha)$ . If there is a conflict for the look-ahead among some production right sides, the grammar is not suitable for recursive-descent parsing. A production with  $\Lambda$  on the right side will be used, if the look-ahead is not in  $FIRST$  of any right side.
  - Procedure code copies the right side of a production. Nonterminal means calling a procedure for this nonterminal. Terminal is compared with the look-ahead. If they are equal, a next terminal is read. If they are not equal, it is an error.



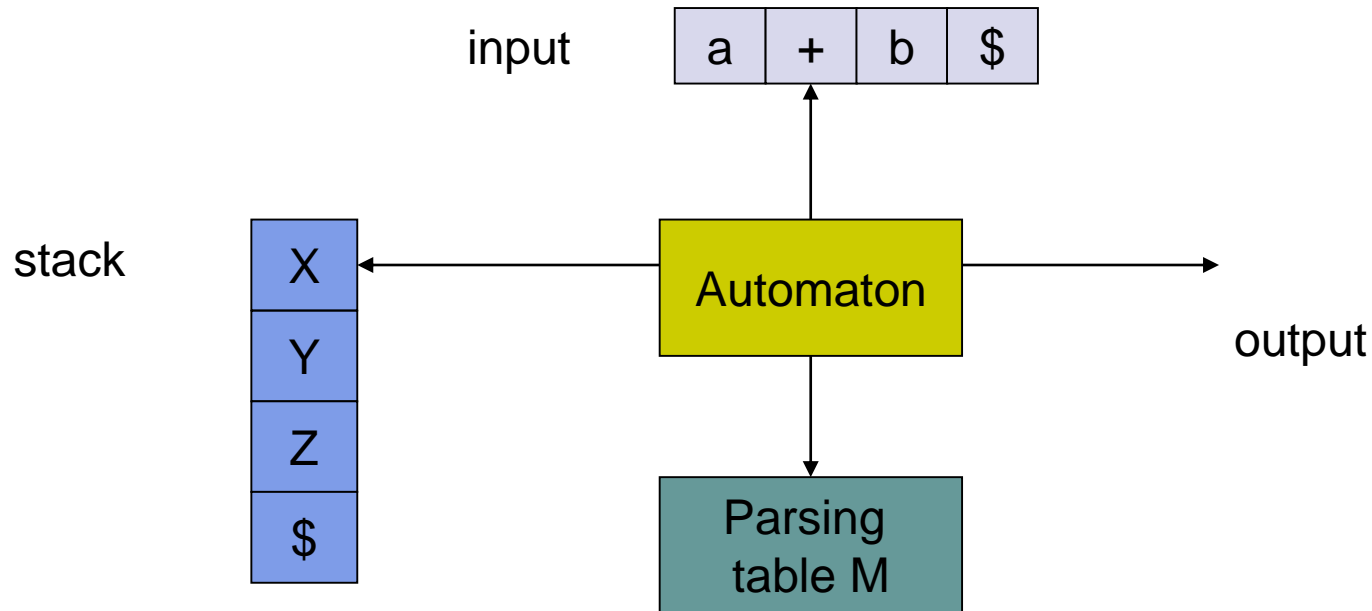
# Recursive-descent parsing – example for our grammar



```
void match(token t) {
    if(lookahead==t)
        lookahead = nexttoken();
    else error();
}
void E(void) {
    T(); Eap();
}
void Eap(void) {
    if(lookahead=='+') {
        match('+'); T(); Eap();
    }
}
void T(void) {
    F(); Tap();
}
```

```
void Tap(void) {
    if(lookahead=='*') {
        match('*'); F(); Tap();
    }
}
void F(void) {
    switch(lookahead) {
        case '(': match('('); E();
                    match(')'); break;
        case 'id':
                    match('id'); break;
        default:
                    error();
    }
}
```

# Nonrecursive predictive parsing



- Parsing table  $M[A, a]$ , where  $A$  is nonterminal and  $a$  is terminal
- The stack contains grammar symbols



# LL(1) automaton behavior

- Initial configuration
  - Input pointer points to the first terminal in the input string
  - The stack contains the start symbol of the grammar on top of \$
- In each step, the automaton decides, what to do, using a symbol  $X$  on top of the stack and a terminal  $a$ , pointed by input pointer
  - If  $X=a=\$$ , the parser halts, parsing finished successfully
  - If  $X=a\neq \$$ , the parser pops  $X$  from the stack and advances the input pointer to the next input symbol
  - If  $X\neq a$  and  $X\in T$ , the parser reports error
  - If  $X$  is a nonterminal, the parser uses entry  $M[X, a]$ . If this entry is a production, the parser replaces  $X$  on top of the stack by the right side (leftmost symbol on top of the stack). At the same time, the parser generates an output about using the production. If the entry is **error**, the parser informs about a syntax error.

# Construction of predictive parsing tables



- For each production  $A \rightarrow \alpha$  do following steps
  - For  $\forall a \in \text{FIRST}(\alpha)$  add  $A \rightarrow \alpha$  to  $M[A, a]$
  - If  $\Lambda \in \text{FIRST}(\alpha)$ , add  $A \rightarrow \alpha$  to  $M[A, b]$   
 $\forall b \in \text{FOLLOW}(A)$ . Moreover, if  $\$ \in \text{FOLLOW}(A)$ ,  
add  $A \rightarrow \alpha$  to  $M[A, \$]$
- Mark each empty entry in  $M$  as **error**

# Example of table construction for our grammar



	id	+	*	(	)	\$
E	$E \rightarrow TE'$			$E \rightarrow TE'$		
E'		$E' \rightarrow +TE'$			$E' \rightarrow \Lambda$	$E' \rightarrow \Lambda$
T	$T \rightarrow FT'$			$T \rightarrow FT'$		
T'		$T' \rightarrow \Lambda$	$T' \rightarrow *FT'$		$T' \rightarrow \Lambda$	$T' \rightarrow \Lambda$
F	$F \rightarrow \text{id}$			$F \rightarrow (E)$		

# Example of parser behavior for our grammar



Stack	Input	Output
\$E	id+id*id\$	
\$E'T	id+id*id\$	$E \rightarrow TE'$
\$E'T'F	id+id*id\$	$T \rightarrow FT'$
\$E'T'id	id+id*id\$	$F \rightarrow id$
\$E'T'	+id*id\$	
\$E'	+id*id\$	$T' \rightarrow \Lambda$
\$E'T+	+id*id\$	$E' \rightarrow +TE'$
\$E'T	id*id\$	
\$E'T'F	id*id\$	$T \rightarrow FT'$

Stack	Input	Output
\$E'T'id	id*id\$	$F \rightarrow id$
\$E'T'	*id\$	
\$E'T'F*	*id\$	$T' \rightarrow *FT'$
\$E'T'F	id\$	
\$E'T'id	id\$	$F \rightarrow id$
\$E'T'	\$	
\$E'	\$	$T' \rightarrow \Lambda$
\$	\$	$E' \rightarrow \Lambda$



# LL(1) grammar

- Context-free grammar  $G=(T,N,S,P)$  is a LL(1) grammar, if and only if whenever  $A \rightarrow \alpha$ ,  $A \rightarrow \beta \in P$  are two distinct ( $\alpha \neq \beta$ ) productions of  $G$  and we have any left sentential forms  $uA\gamma$ ,  $vA\delta$ , where  $u,v \in T^*$  and  $\gamma, \delta \in (T \cup N)^*$ , the following condition holds:
  - $\text{FIRST}(\alpha\gamma) \cap \text{FIRST}(\beta\delta) = \emptyset$
- Simplified detection: no ambiguous or left-recursive grammar can be LL(1)



# Grammar terminology

- $PXY(k)$
- $X$  – direction of the input reading
  - In our case always  $L$ , i.e. from left to right
- $Y$  – kind of derivation
  - $L$  – left derivation
  - $R$  – right derivation
- $P$  – prefix
  - Subtle division of some grammar classes
- $k$  – look-ahead
  - An integer, usually 1, can be 0 or more generally  $k$
- Examples
  - $LL(1)$ ,  $LR(0)$ ,  $LR(1)$ ,  $LL(k)$ ,  $SLR(1)$ ,  $LALR(1)$



# Expanding definition of FIRST and FOLLOW on k



- If  $\alpha$  is a string from grammar symbols, then  $\text{FIRST}_k(\alpha)$  is a set of terminal words with maximal length  $k$ , which are on the beginning of at least one string derived from  $\alpha$ . If  $\alpha$  can be derived on  $\Lambda$ , then  $\Lambda$  is in  $\text{FIRST}_k(\alpha)$ .
- $\text{FOLLOW}_k(A)$  for nonterminal  $A$  is a set of terminal words with maximal length  $k$ , which can be on the right side of  $A$  in any string derived from the start nonterminal ( $S \Rightarrow^* \alpha A u \beta$  for some  $\alpha$  and  $\beta$ ). If  $A$  is the right-most symbol in any sentential form, then  $\$$  is in  $\text{FOLLOW}_k(A)$ .



# LL(k) grammar

- Context-free grammar  $G=(T,N,S,P)$  is a strong LL(k) grammar for  $k \geq 1$ , if and only if whenever  $A \rightarrow \alpha, A \rightarrow \beta \in P$  are two distinct ( $\alpha \neq \beta$ ) productions and we have any left sentential forms  $uA\gamma, vA\delta$ , where  $u, v \in T^*$  and  $\gamma, \delta \in (T \cup N)^*$ , the following condition holds:
  - $\text{FIRST}_k(\alpha\gamma) \cap \text{FIRST}_k(\beta\delta) = \emptyset$ .
- LL(k) (not strong)
  - $u=v, \gamma=\delta$

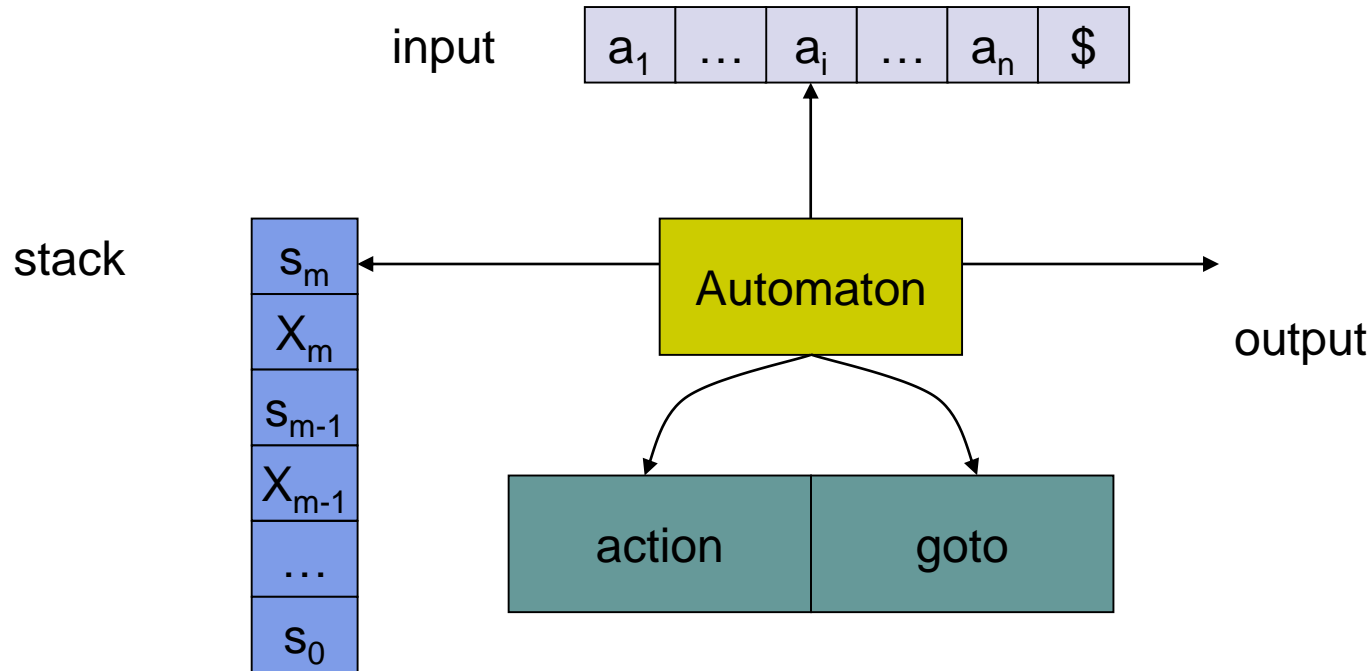


# Bottom-up analysis

- Attempts to find in reverse the rightmost derivation for an input string
- Attempts to construct a parse tree for an input string beginning at the leaves and working up towards the root
- Replace a substring corresponding to a right side of a production by a nonterminal from the left side of the production in each reduce step
- Used in parser generators
  - Bison – LALR(1), GLR(1)
- Advantages against LL(1) parsers
  - It can be implemented with the same efficiency as top-down parsing
  - The class of decidable languages LR(1) is a proper superset of LL(1)
- SLR(1), LR(1), LALR(1)



# LR parser automaton



- $s_i$  are states
  - A state on the top of the stack is the current state of the automaton
- $x_i$  are grammar symbols



# LR(1) automaton behavior

- Initial configuration
  - Input pointer points to the first terminal in the input string
  - Initial state  $s_0$  is on the stack
- In each step address table action[ $s_m, a_i$ ] using  $s_m$  and  $a_i$ 
  - Shift  $s$ , where  $s$  is a new state
    - It shifts the input tape by 1 terminal and add  $a_i$  and  $s$  on the top of the stack
  - Reduce using production  $A \rightarrow \alpha$ 
    - Remove  $r=|\alpha|$  pairs  $(s_k, X_k)$  from the top of the stack, add  $A$  on the top of the stack and then goto[ $s_{m-r}, A$ ] ( $s_{m-r}$  is a state on the top of the stack after erasing pairs)
    - Generate an output
  - Accept
    - The input string is accepted
    - Generate an output
  - Error
    - The input string is not in the input language

# LR automaton tables for our grammar



state	action						goto		
	id	+	*	(	)	\$	E	T	F
0	s5			s4			1	2	3
1		s6				acc			
2		r2	s7		r2	r2			
3		r4	r4		r4	r4			
4	s5			s4			8	2	3
5		r6	r6		r6	r6			
6	s5			s4				9	3
7	s5			s4					10
8		s6			s11				
9		r1	s7		r1	r1			
10		r3	r3		r3	r3			
11		r5	r5		r5	r5			



# Example of LR parser behavior

Stack	Input	Action
0	<b>id+id*id\$</b>	s5
0 <b>id</b> 5	<b>+id*id\$</b>	r6: $F \rightarrow id$
0 F 3	<b>+id*id\$</b>	r4: $T \rightarrow F$
0 T 2	<b>+id*id\$</b>	r2: $E \rightarrow T$
0 E 1	<b>+id*id\$</b>	s6
0 E 1 <b>+</b> 6	<b>id*id\$</b>	s5
0 E 1 <b>+</b> 6 <b>id</b> 5	<b>*id\$</b>	r6: $F \rightarrow id$
0 E 1 <b>+</b> 6 F 3	<b>*id\$</b>	r4: $T \rightarrow F$
0 E 1 <b>+</b> 6 T 9	<b>*id\$</b>	s7
0 E 1 <b>+</b> 6 T 9 <b>*</b> 7	<b>id\$</b>	s5
0 E 1 <b>+</b> 6 T 9 <b>*</b> 7 <b>id</b> 5	<b>\$</b>	r6: $F \rightarrow id$
0 E 1 <b>+</b> 6 T 9 <b>*</b> 7 F 10	<b>\$</b>	r3: $T \rightarrow T * F$
0 E 1 <b>+</b> 6 T 9	<b>\$</b>	r1: $E \rightarrow E + T$
0 E 1	<b>\$</b>	acc



# LR(k) grammar

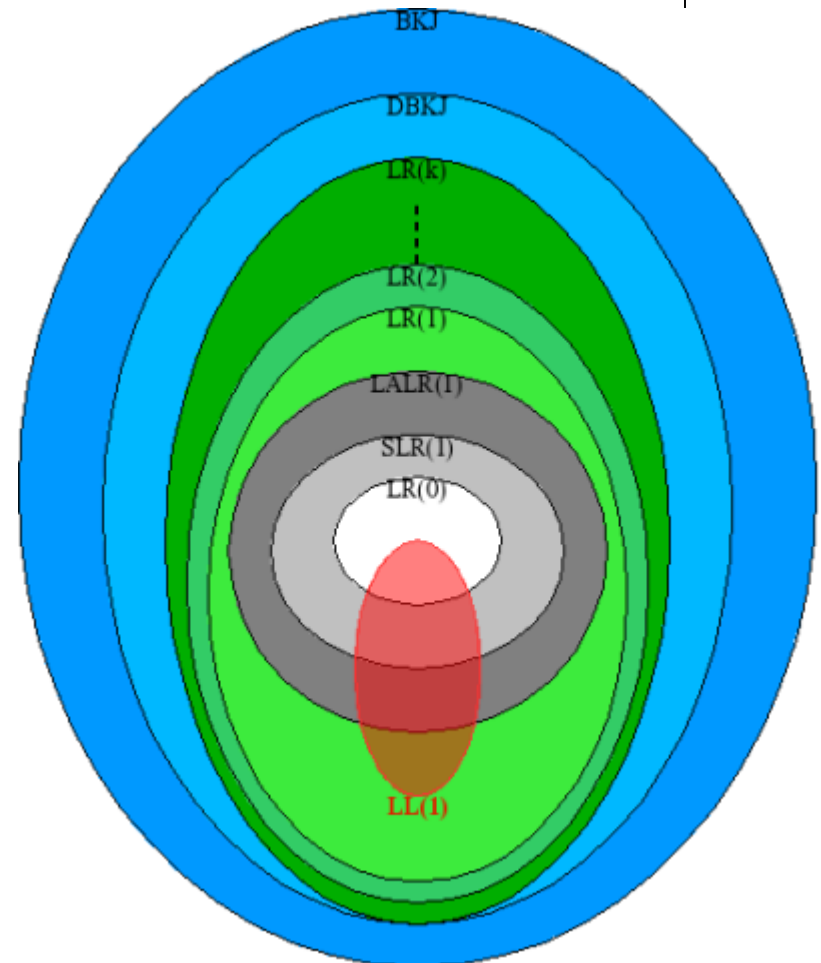
- Context-free grammar  $G=(T,N,S,P)$  is LR(k) grammar for  $k \geq 1$ , if and only if whenever  $A \rightarrow \alpha, A \rightarrow \beta \in P$  are two distinct ( $\alpha \neq \beta$ ) productions of  $G$  and we have any two right sentential forms  $\gamma Au, \delta Av$ , where  $u, v \in T^*$  and  $\gamma, \delta \in (T \cup N)^*$ , the following condition holds:
  - $\text{FIRST}_k(u) \cap \text{FIRST}_k(v) = \emptyset$



# Grammars (languages) strength



- Union of all  $LR(k)$  are deterministic context-free languages (DBKJ) and it is a proper subset of all context-free languages (BKJ)





# Grammar augmentation

- Augmentation of a grammar  $G=(T,N,S,P)$  is a grammar  $G'=(T,N',S',P')$ , where  $N'=N\cup\{S'\}$  and  $P'=P\cup\{S'\rightarrow S\}$
- The augmentation is not necessary whenever  $S$  is on the left side of one production and it isn't on any right side of grammar productions
- It helps recognize the end of parsing
- For our grammar:
  - $S'\rightarrow E$



# LR(0) items

- LR(0) item of a grammar  $G$  is a production with a special symbol dot on the right side
  - Special symbol is a valid symbol for comparison of two LR(0) items of a same production. LR(0) items of the same production are different, whenever the dot is on different position. Moreover, the dot is not a grammar symbol
- An example for production  $E \rightarrow E + T$ :

$E \rightarrow \blacklozenge E + T$	$E \rightarrow E + \blacklozenge T$
$E \rightarrow E \blacklozenge + T$	$E \rightarrow E + T \blacklozenge$



# The closure operation

- If  $I$  is a set of LR(0) items for a grammar  $G$ , then  $\text{CLOSURE}(I)$  is a set of LR(0) items constructed from  $I$  by following rules:
  - Add  $I$  to the  $\text{CLOSURE}(I)$
  - $\forall A \rightarrow \alpha \blacklozenge B \beta \in \text{CLOSURE}(I)$ , where  $B \in N$ , add  $\forall B \rightarrow \gamma \in P$  to  $\text{CLOSURE}(I)$  LR(0) item  $B \rightarrow \blacklozenge \gamma$ , if it is not already there. Apply this rule until no more new LR(0) items can be added to  $\text{CLOSURE}(I)$

# Example of closure for our grammar



- $I = \{S' \rightarrow \blacklozenge E\}$
- $\text{CLOSURE}(I) =$ 
  - $S' \rightarrow \blacklozenge E$
  - $E \rightarrow \blacklozenge E + T$
  - $E \rightarrow \blacklozenge T$
  - $T \rightarrow \blacklozenge T * F$
  - $T \rightarrow \blacklozenge F$
  - $F \rightarrow \blacklozenge ( E )$
  - $F \rightarrow \blacklozenge \text{id}$



# GOTO operation

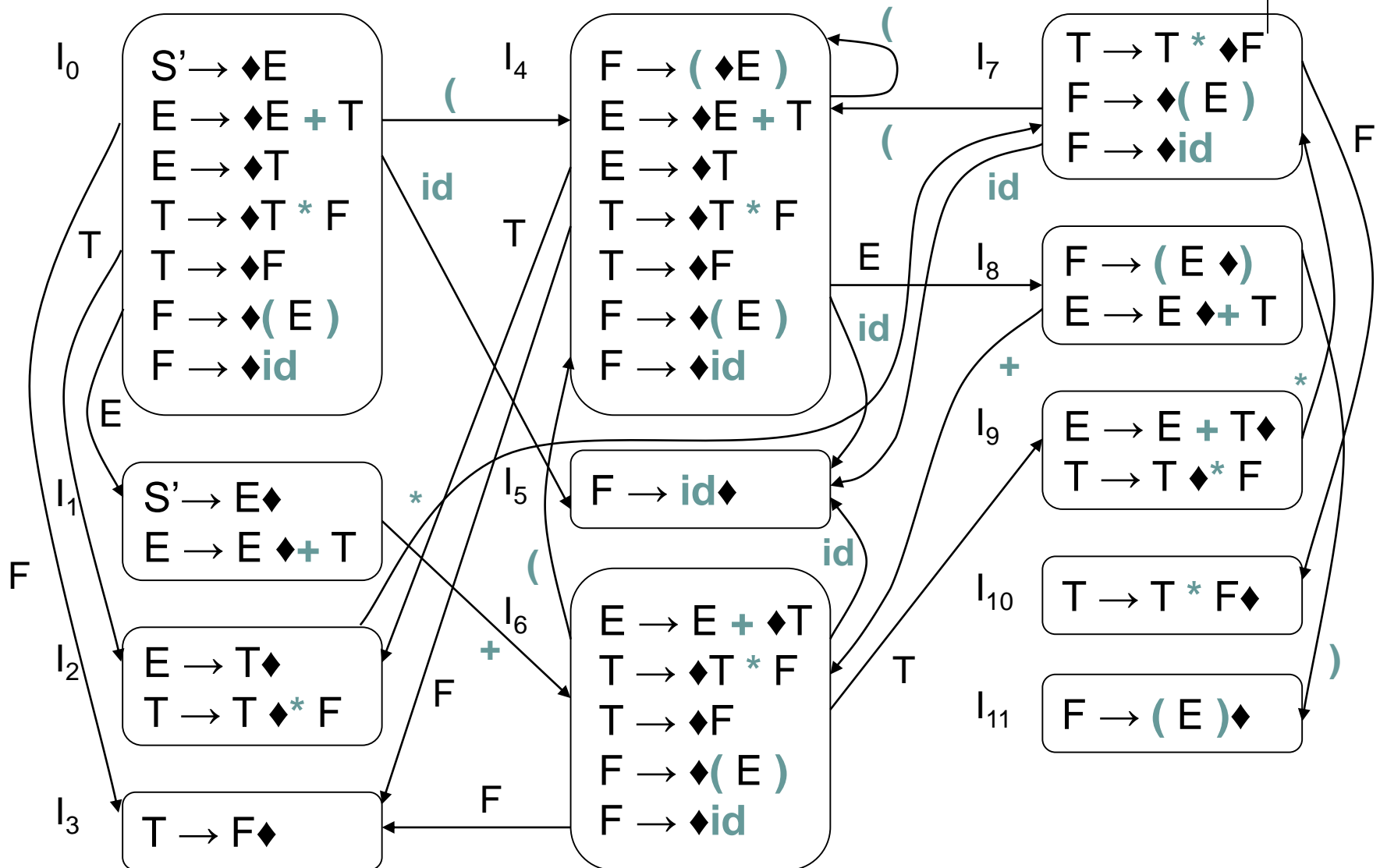
- GOTO( $I$ ,  $X$ ) operation for a set  $I$  of LR(0) items and a grammar symbol  $X$  is defined to be the closure of the set of all LR(0) items  $A \rightarrow \alpha X \diamond \beta$  such that  $A \rightarrow \alpha \diamond X \beta \in I$

# Construction of canonical collection of sets of LR(0) items



- We have an augmented grammar  $G'=(T,N',S',P')$
- Construction of canonical collection  $C$  of sets of LR(0) items:
  - We start with  $C=\{ \text{CLOSURE}(\{S' \rightarrow \blacklozenge S\}) \}$
  - $\forall I \in C$  and  $\forall X \in T \cup N'$  such as  $\text{GOTO}(I, X) \notin C \wedge \text{GOTO}(I, X) \neq \emptyset$ , add  $\text{GOTO}(I, X)$  to  $C$ . Repeat this step, until something new is added to  $C$ .

# Construction of canonical collection for our grammar







# Valid items

- LR(0) item  $A \rightarrow \beta_1 \blacklozenge \beta_2$  is a valid item for a viable prefix  $\alpha\beta_1$ , if there is a rightmost derivation  $S' \Rightarrow^+ \alpha A w \Rightarrow \alpha \beta_1 \beta_2 w$
- It is a great hint for a parser. It helps to decide, if the parser should make a shift or a reduction, if  $\alpha\beta_1$  is on top of the stack
- Basic LR parsing theorem: A set of valid items for a viable prefix  $\gamma$  is exactly a set of items reachable from the initial state through the prefix  $\gamma$  by deterministic finite automaton constructed from canonical collection with GOTO transitions.

# SLR(1) automaton construction



- We have an augmented grammar  $G'$ . Tables of SLR(1) automaton are constructed by following algorithm
  - Construct a canonical collection  $C$  of sets of LR(0) items
  - State  $i$  is constructed from  $I_i$ . The parsing actions for state  $i$  are determined as follows
    - $A \rightarrow \alpha \diamond a \beta \in I_i, a \in T \wedge \text{GOTO}(I_i, a) = I_j$ , then  $\text{action}[i, a] = \text{shift } j$
    - $A \rightarrow \alpha \diamond \in I_i$ , then  $\forall a \in \text{FOLLOW}(A) \wedge A \neq S'$  is  $\text{action}[i, a] = \text{reduce } A \rightarrow \alpha$
    - $S' \rightarrow S \diamond \in I_i$ , then  $\text{action}[i, \$] = \text{accept}$
  - If there is a conflict in the previous step, the grammar is not a SLR(1) grammar and the automaton cannot be constructed
  - Table goto is indexed by state  $i$  and  $A \in N'$ : whenever  $\text{GOTO}(I_i, A) = I_j$ , then  $\text{goto}[i, A] = j$
  - All empty cells are filled by error instruction
  - The initial state of the parser is the state, which contains LR(0) item  $S' \rightarrow \diamond S$



# Full LR(1) automata

- $\text{action}[i, a]$  is set to reduction  $A \rightarrow \alpha$ , when  $A \rightarrow \alpha \diamond \in I_i$ ,  $\forall a \in \text{FOLLOW}(A)$  for a state  $i$  during SLR(1) construction
- In some situation, when  $i$  is on top of the stack, the viable prefix  $\beta\alpha$  is in form, where  $\beta A$  cannot be followed by a terminal  $a$  in any right sentential form. Therefore reduction  $A \rightarrow \alpha$  is for lookahead  $a$  invalid.
- Solution: add more information to states, so we can avoid invalid reductions.



# LR(1) items

- The added information is stored as an additional terminal for each LR(0) item. Such item has a form  $[A \rightarrow \alpha \blacklozenge \beta, a]$ , where  $A \rightarrow \alpha \beta \in P$ ,  $a \in T$ , and we call it LR(1) item. The terminal **a** is called lookahead.
  - The lookahead has no meaning for  $A \rightarrow \alpha \blacklozenge \beta$ , where  $\beta \neq \Lambda$
  - Reduction  $A \rightarrow \alpha$  is set only when  $[A \rightarrow \alpha \blacklozenge, a] \in I_i$  for current state  $i$  and a terminal **a** on the input
  - A set of terminals created from lookaheads of LR(1) items  $\subseteq \text{FOLLOW}(A)$
- LR(1) item  $[A \rightarrow \alpha \blacklozenge \beta, a]$  is valid for viable prefix  $\gamma$ , whenever  $\exists$  right derivation  $S \Rightarrow^+ \delta A w \Rightarrow \delta \alpha \beta w$ , where
  - $\gamma = \delta \alpha$
  - Either **a** is the first symbol of  $w$  or  $w = \Lambda$  and **a** is \$



# Closure for LR(1) items

- We have a set of LR(1) items  $I$  for a grammar  $G$ . We define  $CLOSURE_1(I)$  as a set of LR(1) items constructed from  $I$  by following procedure:
  - Add set  $I$  to  $CLOSURE_1(I)$
  - $\forall [A \rightarrow \alpha \blacklozenge B \beta, a] \in CLOSURE_1(I)$ , where  $B \in N$ , add LR(1) item  $[B \rightarrow \blacklozenge \gamma, b] \forall B \rightarrow \gamma \in P$  and  $\forall b \in FIRST(\beta a)$  to  $CLOSURE_1(I)$ , if it isn't there already. Repeat this step, until something is added to  $CLOSURE_1(I)$ .

# GOTO operation for LR(1) items



- We define  $\text{GOTO}_1(I, X)$  operation for a set  $I$  of LR(1) items and a grammar symbol  $X$  as a  $\text{CLOSURE}_1$  of a set of all items  $[A \rightarrow \alpha X \diamond \beta, a]$  where  $[A \rightarrow \alpha \diamond X \beta, a] \in I$

# Construction of canonical collection of sets of LR(1) items



- We have an augmented grammar  $G'=(T,N',S',P')$
- Construction of canonical collection  $C$  of LR(1) items:
  - We start with  $C=\{ \text{CLOSURE}_1(\{[S' \rightarrow \diamond S, \$]\}) \}$
  - Add  $\text{GOTO}_1(I, X)$  to  $C \forall I \in C$  and  $\forall X \in T \cup N'$ , where  $\text{GOTO}_1(I, X) \notin C \wedge \text{GOTO}_1(I, X) \neq \emptyset$ . Repeat this step, until something new is added to  $C$ .

# Example of LR(1) grammar, which is not SLR(1)



- $S' \rightarrow S$
- $S \rightarrow CC$
- $C \rightarrow cC$
- $C \rightarrow d$

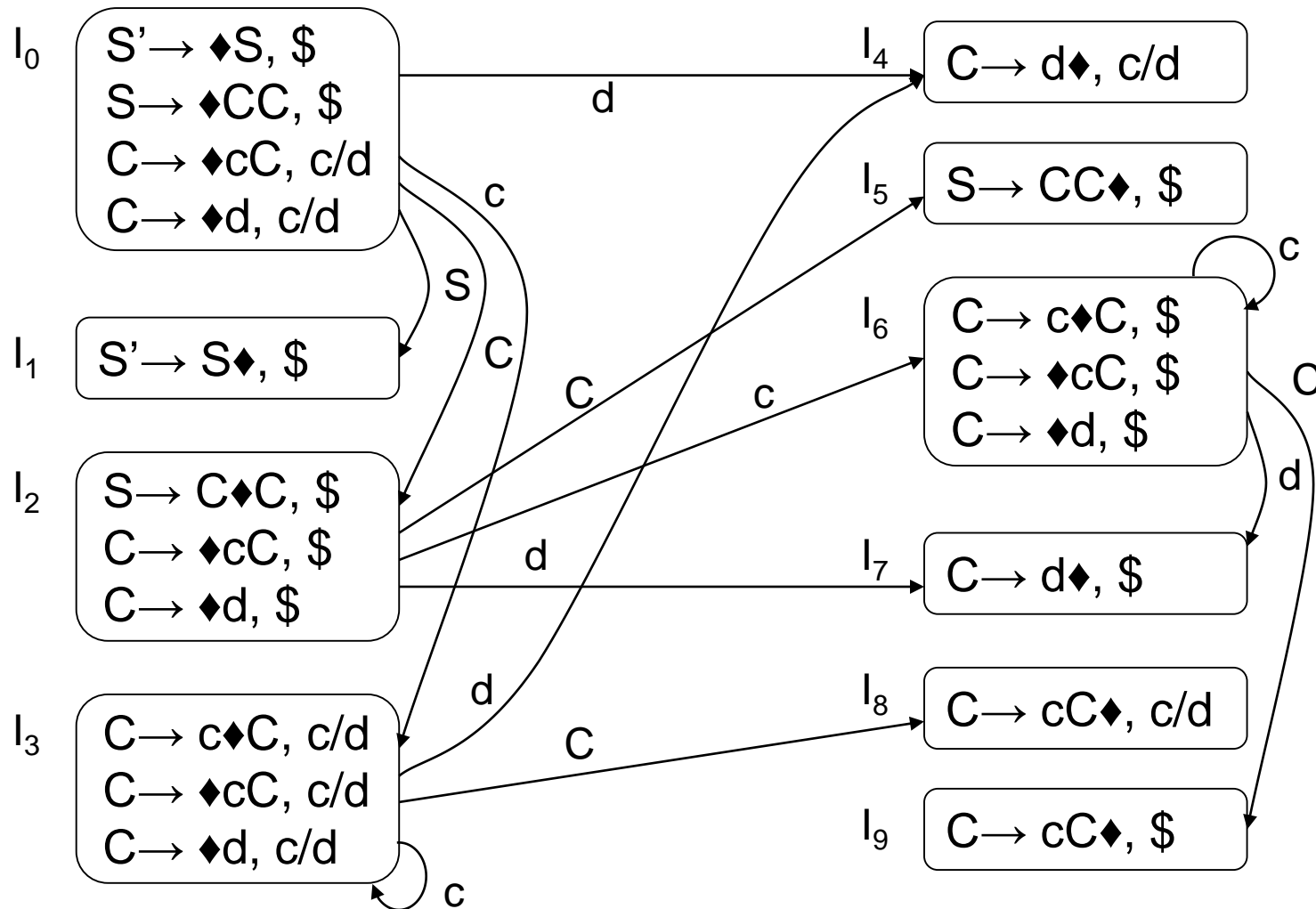


# Example of closure construction for LR(1) items



- $I = \{[S' \rightarrow \blacklozenge S, \$]\}$
- $\text{CLOSURE}_1(I) =$ 
  - $S' \rightarrow \blacklozenge S, \$ \quad \beta = \Lambda, \text{FIRST}(\beta \$) = \text{FIRST}(\$) = \{\$\}$
  - $S \rightarrow \blacklozenge CC, \$ \quad \beta = C, \text{FIRST}(C \$) = \{c, d\}$
  - $C \rightarrow \blacklozenge cC, c/d$
  - $C \rightarrow \blacklozenge d, c/d$

# Example of construction of canonical collection of LR(1) items





# LR(1) parser construction

- We have an augmented grammar  $G'$ . LR(1) automaton tables are constructed by following algorithm
  - Construct a canonical collection  $C$  of sets of LR(1) items
  - State  $i$  is constructed from  $I_i$ . The parsing actions for state  $i$  are determined as follows
    - $[A \rightarrow \alpha \blacklozenge a \beta, b] \in I_i, a \in T \wedge \text{GOTO1}(I_i, a) = I_j$ , then  $\text{action}[i, a] = \text{shift } j$
    - $[A \rightarrow \alpha \blacklozenge, a] \in I_i \wedge A \neq S'$ , then  $\text{action}[i, a] = \text{reduce } A \rightarrow \alpha$
    - $[S' \rightarrow S \blacklozenge, \$] \in I_i$ , then  $\text{action}[i, \$] = \text{accept}$
  - If there is a conflict in the previous step, the grammar is not a LR(1) grammar and the automaton cannot be constructed
  - Table goto is indexed by state  $i$  and  $A \in N'$ : whenever  $\text{GOTO1}(I_i, A) = I_j$ , then  $\text{goto}[i, A] = j$
  - All empty cells are filled by error instruction
  - The initial state of the parser is the state, which contains LR(1) item  $[S' \rightarrow \blacklozenge S, \$]$



# LALR

- LALR=LookAhead-LR
- Often used in practice
  - Bison
  - Most common programming languages can be expressed by an LALR grammar
  - Parser tables are considerably smaller than LR(1) tables
- SLR and LALR parsers have the same number of states, LR parsers have greater number of states
  - Common languages have hundreds of states
  - LR(1) parsers have thousands of states for the same grammar



# How to make smaller tables?

- Idea: merge sets with the same core into one set including GOTO1 merge
  - Core: a set of LR(0) items (no lookahead)
  - Merge cannot produce shift/reduce conflict
    - Suppose in the union there is a conflict on lookahead **a** for LR(1) items  $[A \rightarrow \alpha \blacklozenge, a]$  and  $[B \rightarrow \beta \blacklozenge a \gamma, b]$
    - Cores are same, therefore in the set with  $[A \rightarrow \alpha \blacklozenge, a]$  must be  $[B \rightarrow \beta \blacklozenge a \gamma, c]$  as well for some **c**. There was already a shift/reduce conflict before merge
  - Merge can produce reduce/reduce conflict

# Easy LALR(1) table construction



- We have an augmented grammar  $G'$ . LALR(1) automaton tables are constructed by following algorithm
  - Construct a canonical collection  $C$  of sets of LR(1) items
  - For each core in collection  $C$ , find all sets having that core, and replace these sets by their union
  - Let  $C' = \{ J_0, J_1, \dots, J_m \}$  be the resulting collection of LR(1) items
  - Table action is constructed for  $C'$  in the same manner as for full LR(1) parser
  - If there is a conflict, the grammar is not LALR(1) grammar
  - If  $J \in C'$  is the union of sets of LR(1) items  $I_i$  ( $J = I_1 \cup I_2 \cup \dots \cup I_k$ ), then cores  $GOTO1(I_1, X), \dots, GOTO1(I_k, X)$  are the same, since  $I_1, \dots, I_k$  all have the same core. Let  $K$  be the union of all sets of items having the same core as  $goto(I_1, X)$ . Then  $GOTO1(J, X) = K$
- Important disadvantage – we need to construct full LR(1)

# Compiler principles

---

Intermediate code

Jakub Yaghob





# Intermediate code

- Intermediate representation of the source code
- Separates front end from back end
- Advantages
  - Different back ends for the same input language – support for different CPU architectures
    - gcc
  - Different front ends for the same output language – support more programming languages for the same CPU architecture
    - .NET
  - A machine-independent optimizations – HLO





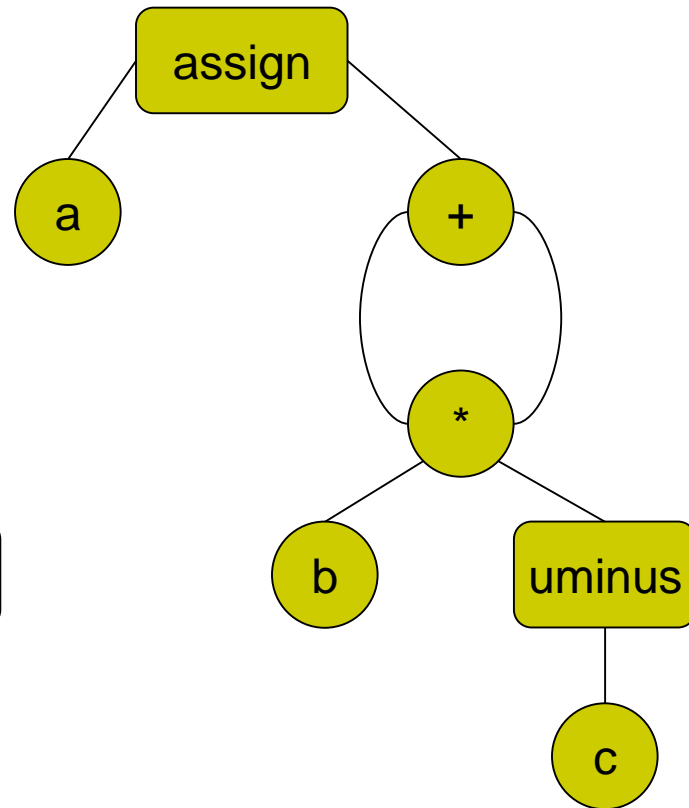
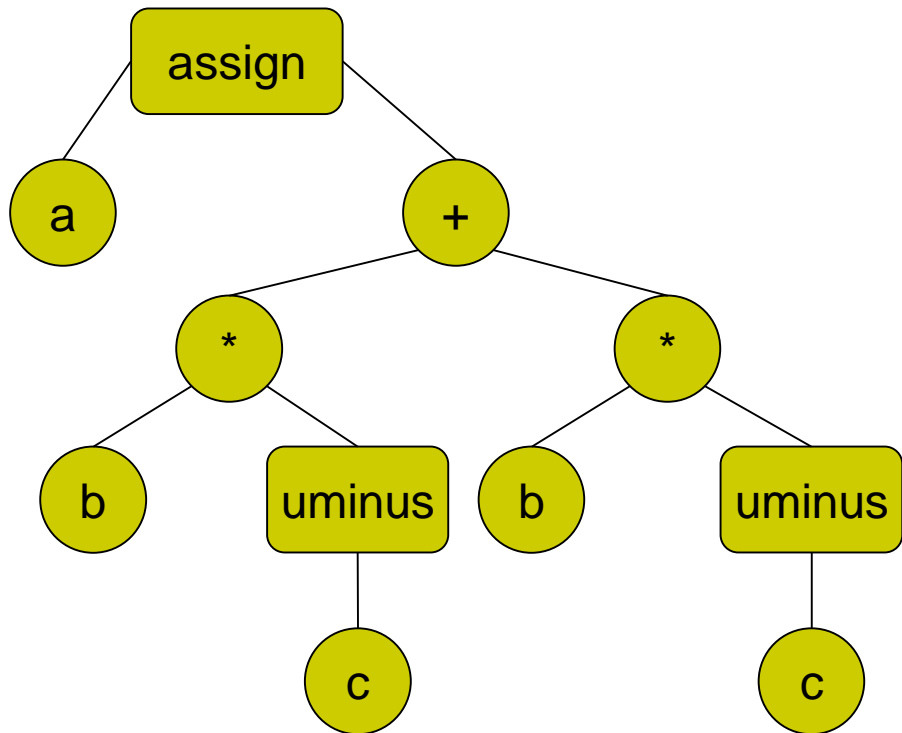
# Intermediate languages

- Syntax tree
  - Can be even DAG
- Postfix notation
  - Linearized representation of a syntax tree
  - Tree edges not in notation, can be reconstructed from the order and number of operands of an operator
- Three-address code
  - Linearized representation of a syntax tree as well
  - A sequence of statements in form
    - $x := y \text{ op } z$

# Examples – syntax tree and DAG



- $a := b^* - c + b^* - c$



# Examples – postfix notation and three-address code



- $a \ b \ c \ \text{uminus} \ * \ b \ c \ \text{uminus} \ * \ + \ \text{assign}$
  - $t1 := -c$
  - $t2 := b * t1$
  - $t3 := -c$
  - $t4 := b * t3$
  - $t5 := t2 + t4$
  - $a := t5$
- $t1 := -c$
  - $t2 := b * t1$
  - $t5 := t2 + t2$
  - $a := t5$



# Three-address code operands

- A name
  - A variable
  - A type
  - Other names
- A constant
  - Different literals (UINT, string, ...)
- Temporary variable
  - Generated by a compiler
  - Easily they can be thought of as a CPU registers

# Types of three-address statements



- Binary arithmetic and logical operation
- Unary operations
- Assignment/copy
- Unconditional jump
- Conditional jump
- Procedure/function call mechanism
  - Parameters, call, return
- Array indexation
- Address operators
  - Address of an object, dereference
- Declaration

# Implementation of three-address code – quadruples



- A record with four fields
  - *op*, *arg1*, *arg2*, *res*
- Some statements don't use an *arg* or even *res*
- Operands are references to symbol tables

	op	arg1	arg2	res
(0)	uminus	c		t1
(1)	*	b	t1	t2
(2)	uminus	c		t3
(3)	*	b	t3	t4
(4)	+	t2	t4	t5
(5)	:=	t5		a

# Implementation of three-address code – triples



- Avoid generating temporary variables
- A record with three fields
  - *op*, *arg1*, *arg2*
- Operands are references to the symbol tables (constants or variables) or a position of the statement that compute a value

	op	arg1	arg2
(0)	uminus	c	
(1)	*	b	(0)
(2)	uminus	c	
(3)	*	b	(2)
(4)	+	(1)	(3)
(5)	:=	a	(4)

# Implementation of three-address code – indirect triples



- One array with triples
- One array with references

	op	arg1	arg2
(0)	uminus	c	
(1)	*	b	[0]
(2)	uminus	c	
(3)	*	b	[2]
(4)	+	[1]	[3]
(5)	:=	a	[4]

	stmt
[0]	(0)
[1]	(1)
[2]	(2)
[3]	(3)
[4]	(4)
[5]	(5)





# Implementation comparison

- Quadruples
  - Better for intermediate code optimization
  - They can be easily moved, names remain the same
- Triples
  - Tighter
  - Poor handling while optimizing – renumbering in the whole intermediate code
- Indirect triples
  - Good handling while optimizing – the change is only in the reference array
  - About the same memory size as quadruples

# Compiler principles

---

Semantic analysis

Jakub Yaghob





# Syntax-directed definitions

- Each grammar symbol has an associated set of attributes
  - Like a record
  - Two kinds of attributes
    - Synthesized
    - Inherited
  - Attributes can represent anything
- Attribute values defined by semantic rules assigned to grammar productions
  - The order of evaluation of semantic rules is determined by the dependency graph
  - Evaluation of semantic rules defines values of attributes



# Kinds of attributes

- Each production  $A \rightarrow \alpha$  has associated with it a set of semantic rules of the form  $b = f(c_1, \dots, c_k)$ , where  $f$  is a function,  $c_i$  are grammar symbols attributes from the given production, and either
  - $b$  is a synthesized attribute of nonterminal  $A$
  - $b$  is an inherited attribute of a grammar symbol on the right side of the production



# Attribute grammar

- Syntax tree with attribute values is called annotated (colored) syntax tree
- Attributed grammar is a syntax directed definition, where functions (semantic rules) don't have side effects

# Attributed grammar for our grammar



1.  $E \rightarrow E_R + T$

$E.val = E_R.val + T.val$

2.  $E \rightarrow T$

$E.val = T.val$

3.  $T \rightarrow T_R * F$

$T.val = T_R.val * F.val$

4.  $T \rightarrow F$

$T.val = F.val$

5.  $F \rightarrow ( E )$

$F.val = E.val$

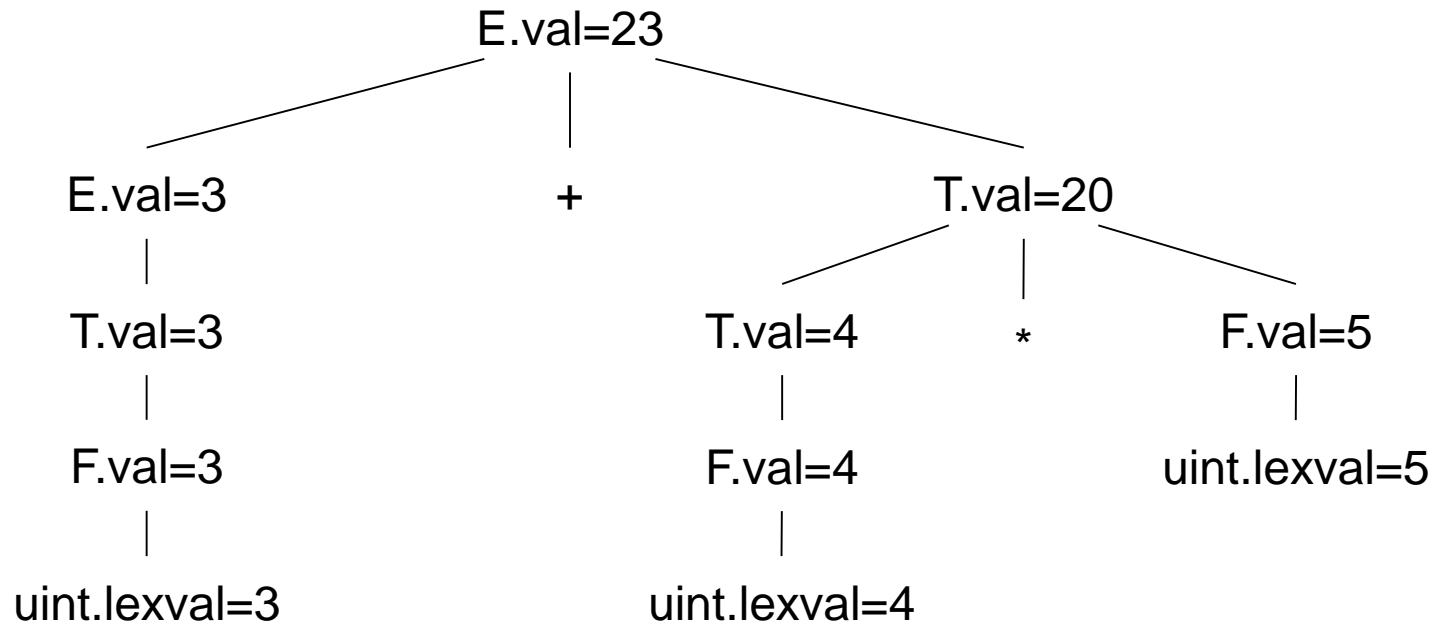
6.  $F \rightarrow \text{uint}$

$F.val = \text{uint.lexval}$



# Example of annotated tree

- $3+4*5$





# Synthesized attributes

- An attribute of a nonterminal from the left side of a production is evaluated based on attributes of symbols from the right side of the production
- Extensively used in practice
  - E.g. expression evaluation
- Attributed grammar, which uses only synthesized attributes, is called S-attributed grammar (purely synthesized attributed grammar)
- S-attributed grammar is easily used in bottom-up analysis
  - Evaluation during reduction
  - Attributes for grammar symbols lie on parser stack





# Inherited attributes

- The value of an inherited attribute is evaluated based on parent and/or siblings attributes
- They are used for expressing dependency of a syntax construction on a context
  - E.g. whether an identifier is on left or right side of an assignment
- We can always rewrite an attribute grammar to a S-attributed grammar



# Example of inherited attributes

- |                                    |   |
|------------------------------------|---|
| 1. $D \rightarrow T L ;$           | $L.in = T.typ$                            |
| 2. $T \rightarrow \text{int}$      | $T.typ = \text{int}$                      |
| 3. $T \rightarrow \text{double}$   | $T.typ = \text{double}$                   |
| 4. $L \rightarrow L_R , \text{id}$ | $L_R.in = L.in$<br>$\text{id}.typ = L.in$ |
| 5. $L \rightarrow \text{id}$       | $\text{id}.typ = L.in$                    |



# Dependency graph

- Construction for a syntax tree
  - Nodes are created for each attribute of each node of the syntax tree
  - For each semantic rule  $b=f(c_1, \dots, c_k)$  construct directed edges from a node of the dependency graph representing  $c_i$  to the node representing  $b$



# Evaluation order

- Any topological sort of a dependency graph gives a valid order in which the semantic rules associated with the nodes in a syntax tree can be evaluated
  - If the dependency graph contains a circle, we are not able to determine any evaluation order



# L-attributed grammar

- Czech alias “jednoduše zleva-doprava 1-průchodová gramatika”
- Attributed grammar is L-attributed, if each inherited attribute of a symbol  $X_j$  on the right side of  $A \rightarrow X_1 \dots X_n$  depends only on
  - The attributes of the symbols  $X_1, \dots, X_{j-1}$  (to the left of  $X_j$ )
  - The inherited attributes of  $A$
- Used for direct attribute evaluation in top-down analysis
- Every S-attributed grammar is L-attributed



# Syntax tree traversal

- If the attribute grammar isn't L-attributed grammar, it will not be possible to evaluate attributes directly during parsing
- We need to fully construct syntax tree. It will be traversed (possibly several times) during the semantic analysis. Several attributes, which we were unable to evaluate during syntax-directed translation, will be evaluated during the traversal

# Static checking during translation



- Type checking
  - Incorrect operand type of an operator
    - Pointer multiplication
- Checking control flow
  - If the change in control flow is legal
    - Break statement in C, if it is not in switch or a loop
- Uniqueness checking
  - Some objects can be defined only once
    - Labels in a function, global objects identifiers
- Name checking
  - Some constructions must have the same name at the start and at the end
    - Assembler procedures



# Symbolic tables

- Growing tables
  - Constants
- Stack tables
  - Identifier visibility in blocks
  - Simple implementation
    - Visibility linked list for a block (stack)
  - Used implementation
    - Identifiers „colored“ by a unique block number





# Error handling

- The compiler must find all errors in the input word and must not show non-existent errors
- Error reporting
  - Clearly and accurately
  - Recover from an error quickly enough and continue in translation
  - Do not significantly slow down the processing of a correct input
- Introduced errors
  - Imprecise recovery from a previous error causes inception of non-existent errors



# Types of errors

- Lexical errors
  - Malformed lexical elements
    - Unfinished string and the EOL
  - Error recovery by ignoring the error
- Syntax errors
  - The input word is not in an input language
    - Unpaired parenthesis
- Semantic errors
  - Static checks
    - Undeclared variable, wrong number of parameters in a function call, wrong type used with an operator
- Logical errors
  - Errors in programming
    - Indefinite loop, uninitialized variable



# Syntax errors recovery

- Panic mode
  - A set of skeletal symbols
  - When an error is encountered, skip all symbols until a symbol from the skeletal set is found
  - Then the parser is put into a known state
- Productions modifications
  - Insert, remove, replace a terminal in a production
- Intentional error production
  - Grammar augmentation with usual errors with specific error message
    - E.g. assignment in Pascal

# Compiler principles

---

Intermediate code generation

Jakub Yaghob



# Syntax-directed translation into three-address code



- Add several attributes to each grammar symbol
  - Placement – a name of an object holding a value of the object
  - Code – a sequence of three-address instructions evaluating the symbol
  - Label – an absolute or relative address to three-address code



# Example for our grammar

- $E \rightarrow E_R + T$   
 $E.p = \text{newtemp}$   
 $E.c = E_R.c \mid T.c \mid \text{gen}(E.p = E_R.p + T.p)$
- $E \rightarrow T$   
 $E.p = T.p$   
 $E.c = T.c$
- $T \rightarrow T_R * F$   
 $T.p = \text{newtemp}$   
 $T.c = T_R.c \mid F.c \mid \text{gen}(T.p = T_R.p * F.p)$
- $T \rightarrow F$   
 $T.p = F.p$   
 $T.c = F.c$
- $F \rightarrow ( E )$   
 $F.p = E.p$   
 $F.c = E.c$
- $F \rightarrow \text{id}$   
 $F.p = \text{id.p}$   
 $F.c = ''$

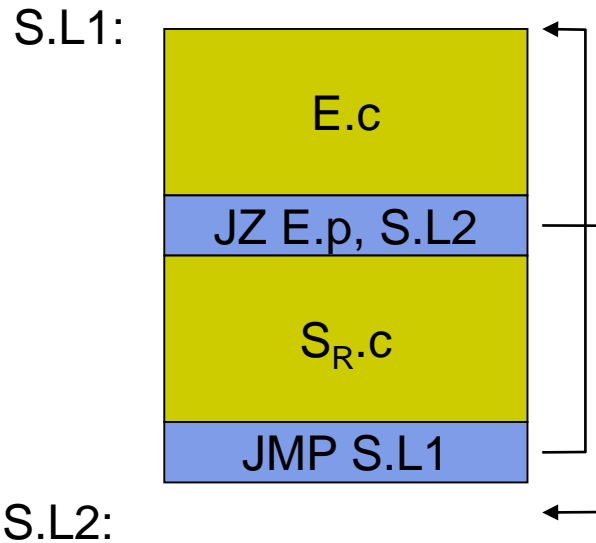


# Example for while (amateur)

- $S \rightarrow \text{while } E \text{ do } S_R$      $S.L1 = \text{curradr}$

$S.L2 = \text{curradr} + E.c.size + S_R.c.size + 2$

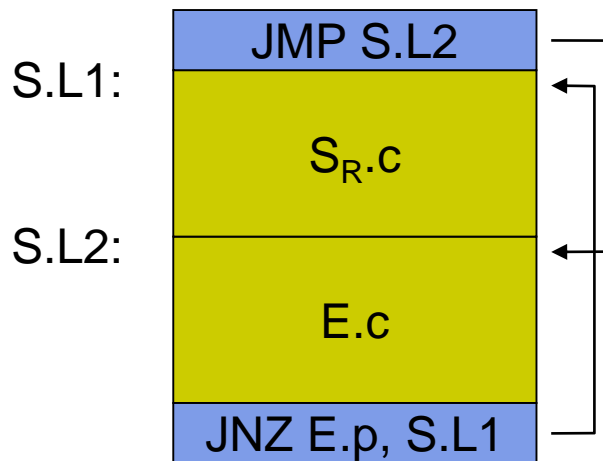
$S.c = E.c \mid \text{gen}(\text{JZ } E.p, S.L2) \mid S_R.c \mid \text{gen}(\text{JMP } S.L1)$



# Example for while (professional)



- $S \rightarrow \text{while } E \text{ do } S_R$



$S.L1 = \text{curradr} + 1$

$S.L2 = \text{curradr} + S_R.c.size + 1$

$S.c = \text{gen}(\text{JMP } S.L2) \mid$   
 $S_R.c \mid E.c \mid \text{gen}(\text{JNZ}$   
 $E.p, S.L1)$





# Declarations

- Global objects declarations
  - E.g. global variables in C
- Local objects declarations
  - E.g. local variables in a function
  - At the start of the function × at the start of a block × in the block body
- Moving declarations to a “good” place
  - Global declarations in one well-defined place
  - Local declarations at the start of a function
    - Live-variable in the function
- Determine the size of the object
- Determine the offset of fields in a structure
  - It does not have to be solved in intermediate code



# Assignment

- Type conversion
  - Usually explicit
    - Intermediate code does not know input language semantic
- Structure field access
  - Use calculated offset from declaration
  - Direct access using a pointer and added constant
- Array unfolding
  - Multidimensional-array are perceived as one-dimensional (like a memory) – stored in one of the two forms, either row-major (row-by-row) or column-major (column-by-column)
  - One-dimensional array  $A[lb..ub]$  of type with width  $w$ 
    - $base + (i-lb) * w$
    - $i * w + (base - lb * w)$
  - Two-dimensional array  $A[lb_1..ub_1, lb_2..ub_2]$ 
    - $base + ((i_1-lb_1) * (ub_2-lb_2+1) + i_2-lb_2) * w$



# Boolean expression

- Sometimes replaced FALSE=0, TRUE=1 (or anything !=0)
- Numeric (full) evaluation
  - All parts of the expression are evaluated
  - Pascal
- Short evaluation
  - if the result is already known during the evaluation, it won't be evaluated further
    - Jumps in the evaluation
  - C



# Switch

- What we need
  - Evaluate the expression
  - Find which value in the list of cases is the same as the value of the expression
  - Execute the statement associated with the value found
- Finding the case
  - Sequence of conditional branches
    - Small number of cases (<10)
  - Binary search in table [value,caseptr]
  - Binary tree of conditions
  - Table of pointers indexed by value
    - High density of values in the range of case values



# Backpatching

- Single pass translation: how to set a destination address in a jump to a forward label?
  - Forward jump
  - Calling a not yet defined function in a module
    - Can be resolved by a linker
- Generate branches with the targets of the jumps temporarily left unspecified
- For each label remember a set of instructions referencing it
- Determine the address of the label
- Go through the set of instructions and *backpatch* the determined address



# Procedure calls

- Evaluate parameters
- Pass parameters by value or by reference
  - VAR or non-VAR parameters in Pascal
- Binding real and formal parameters
  - Positional
  - By name
- Return value

# Compiler principles

---

High-level optimization

Jakub Yaghob





# Optimization

- Ideal situation – generated code is equal to manually written code
  - Is it valid today?
- Reality
  - Only in several defined situations
  - It is difficult
- Optimization
  - Program transformation
  - For speed or size
- High-level optimization
  - Intermediate code





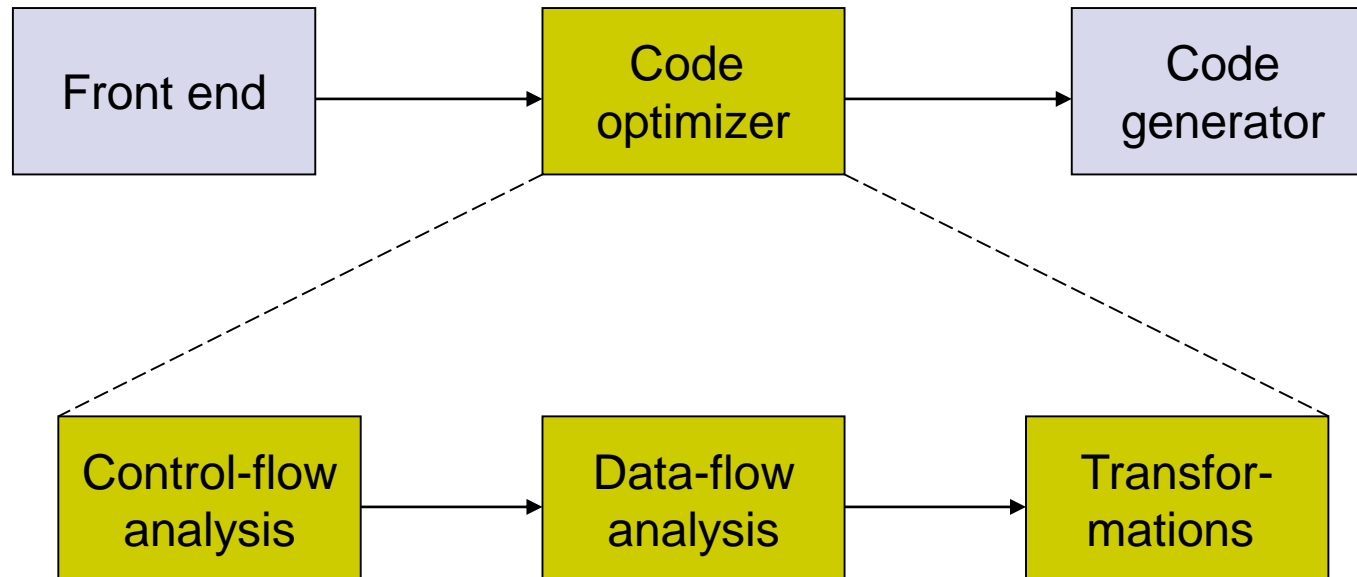
# Criteria for transformation

- Preserve the meaning of programs
  - Output and program behavior must not be changed by a transformation – always use safe approach
- Speed up programs by a measurable amount
  - Sometimes space reduction
  - Occasionally the optimization can slow down a program, on the average it improves it
- A transformation must be worth the effort
  - Complex/slow optimization with negligible effect are not worth the effort

# Organization of the code optimizer



- Control-flow analysis
  - Construct the control flow graph
- Data-flow analysis
  - Live variables





# Control flow graph

- Three-address code graph representation
- Nodes represent computation
- Directed edges represent control flow
- Important for optimization and representation



# Basic block

- A sequence of consecutive three-address statements in which flow of control enters at the beginning and leaves at the end without halt or possibility of branching except at the end
- Algorithm for partitioning a sequence of three-address statements into basic blocks
  - Determine the set of leaders
    - The first statement is a leader
    - Any statement that is the target of a jump is a leader
    - Any statement that immediately follows a jump is a leader
  - For each leader, its basic block consists of the leader and all statements up to but not including the next leader or the end of the sequence

# Construction of control flow graph

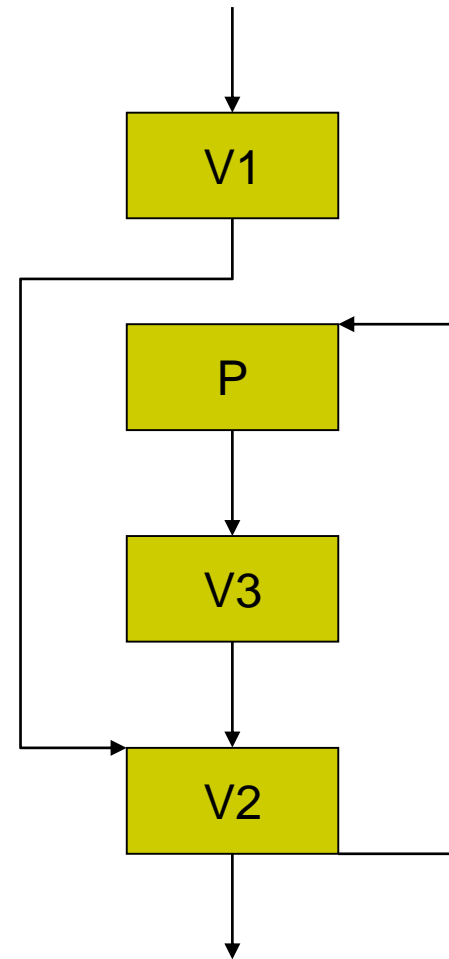
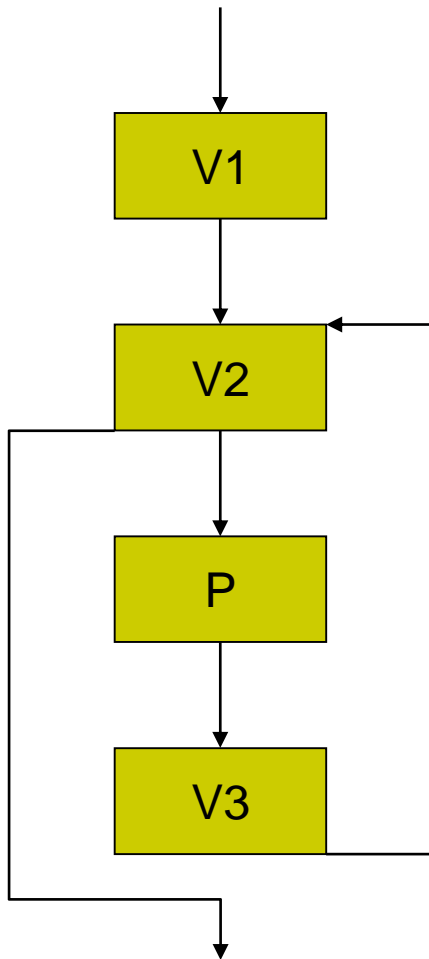


- A directed graph
- The nodes are the basic blocks
- One node is initial – a block whose leader is the first statement in the sequence
- There is a directed edge from block  $B_1$  to block  $B_2$ , if one of the following condition holds
  - The last statement of  $B_1$  is a jump to the leader of  $B_2$
  - $B_2$  immediately follows  $B_1$  and  $B_1$  does not end in an unconditional jump



# Example of control flow graph

- for(V1;V2;V3) P





# Data-flow analysis

- Assign live variables to nodes of the control-flow graph
  - Variables need not to be live in the whole basic block
  - Compute exact lifespan inside the basic block
- Evaluation order
  - Topological sort of DAG
    - Nodes are definitions and uses of a variable
    - Oriented edges are from a use to a definition in a computation
- Pointer modeling
  - Pointer aliasing
- Function calls modeling
  - Interprocedural optimization



# Live variable analysis

- Three-address instruction  $x := y \text{ op } z$  *defines*  $x$  and *uses*  $y$  and  $z$
- A variable is *alive* at some point of the intermediate code, if the point lies on a path from the point, where it is defined, to a point, where it is used, at the same time there are no other definitions on the path





# Kinds of optimizations

- Inside a basic block
  - No jumps
- Inside a function
  - Code movement between basic blocks
  - Create, remove basic blocks
- Whole program (interprocedural)
  - Speed up calls between functions
  - During linker phase with delayed compilation



# Local optimization

- Inside a one basic block
  - Common subexpression elimination (CSE)
  - Copy propagation
  - Dead-code elimination
    - $x = y + z$ , and  $x$  is not used any more
  - Constant folding
  - Algebraic transformations/identities
    - $x = x + 0$ ,  $x = x * 1$



# Local CSE

- $a = b + c$
- $b = a - d$
- $c = b + c$
- $d = a - d$

- $a = b + c$
- $b = a - d$
- $c = b + c$
- $d = b$

- There is a hidden problem
  - Pointer aliasing
  - $ap = \&a$
  - $b = a + c$
  - $*ap = d$
  - $e = a + c$



# Global optimization

- Among more basic blocks, possible movement of a code between BB, control-flow graph modification
  - Common subexpression elimination
  - Copy propagation
  - Dead-code elimination
  - Loop optimization
    - Invariant code motion
    - Reduction in strength of operation
    - Removing induction variable



# Loop optimization – 1

- Invariant code motion
  - Expression yields the same result independent of the number of times a loop is executed

```
while(i<limit-5) S           t = limit - 5;  
                             while(i<t) S
```



# Loop optimization – 2

- Reduction in strength of operation
  - Transform multiplication to addition
- Removing induction variable
  - Only one induction variable for one loop
  - Usually removed during reduction in strength

```
for (i=3; i<8; i+=3)
{
    j = i*2;
    a[j] = j;
}
```

```
for (t=6; t<16; t+=6)
    a[t] = t;
```



## Example – C code

```
i=m-1; j=n; v=a[n];  
for(;;) {  
    do ++i; while(a[i]<v);  
    do --j; while(a[j]>v);  
    if(i>=j) break;  
    x=a[i]; a[i]=a[j]; a[j]=x;  
}  
x=a[i]; a[i]=a[n]; a[n]=x;
```



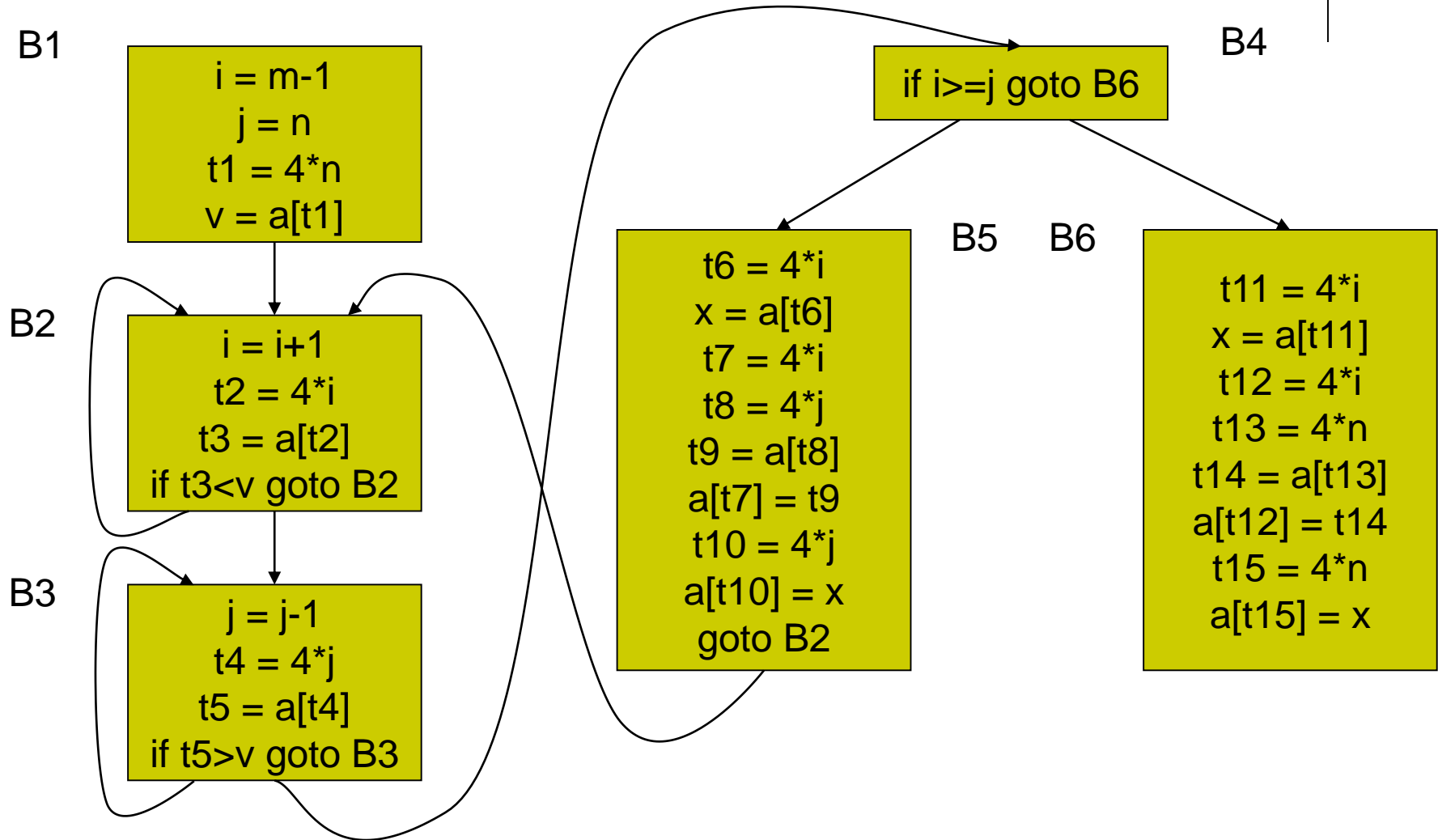
# Example – three-address code

1) $i = m - 1$	11) $t5 = a[t4]$	21) $a[t10] = x$
2) $j = n$	12) if $t5 > v$ goto 9	22) goto 5
3) $t1 = 4 * n$	13) if $i \geq j$ goto 23	23) $t11 = 4 * i$
4) $v = a[t1]$	14) $t6 = 4 * i$	24) $x = a[t11]$
5) $i = i + 1$	15) $x = a[t6]$	25) $t12 = 4 * i$
6) $t2 = 4 * i$	16) $t7 = 4 * i$	26) $t13 = 4 * n$
7) $t3 = a[t2]$	17) $t8 = 4 * j$	27) $t14 = a[t13]$
8) if $t3 < v$ goto 5	18) $t9 = a[t8]$	28) $a[t12] = t14$
9) $j = j - 1$	19) $a[t7] = t9$	29) $t15 = 4 * n$
10) $t4 = 4 * j$	20) $t10 = 4 * j$	30) $a[t15] = x$





# Example – control-flow graph





# Example - LCSE

B5

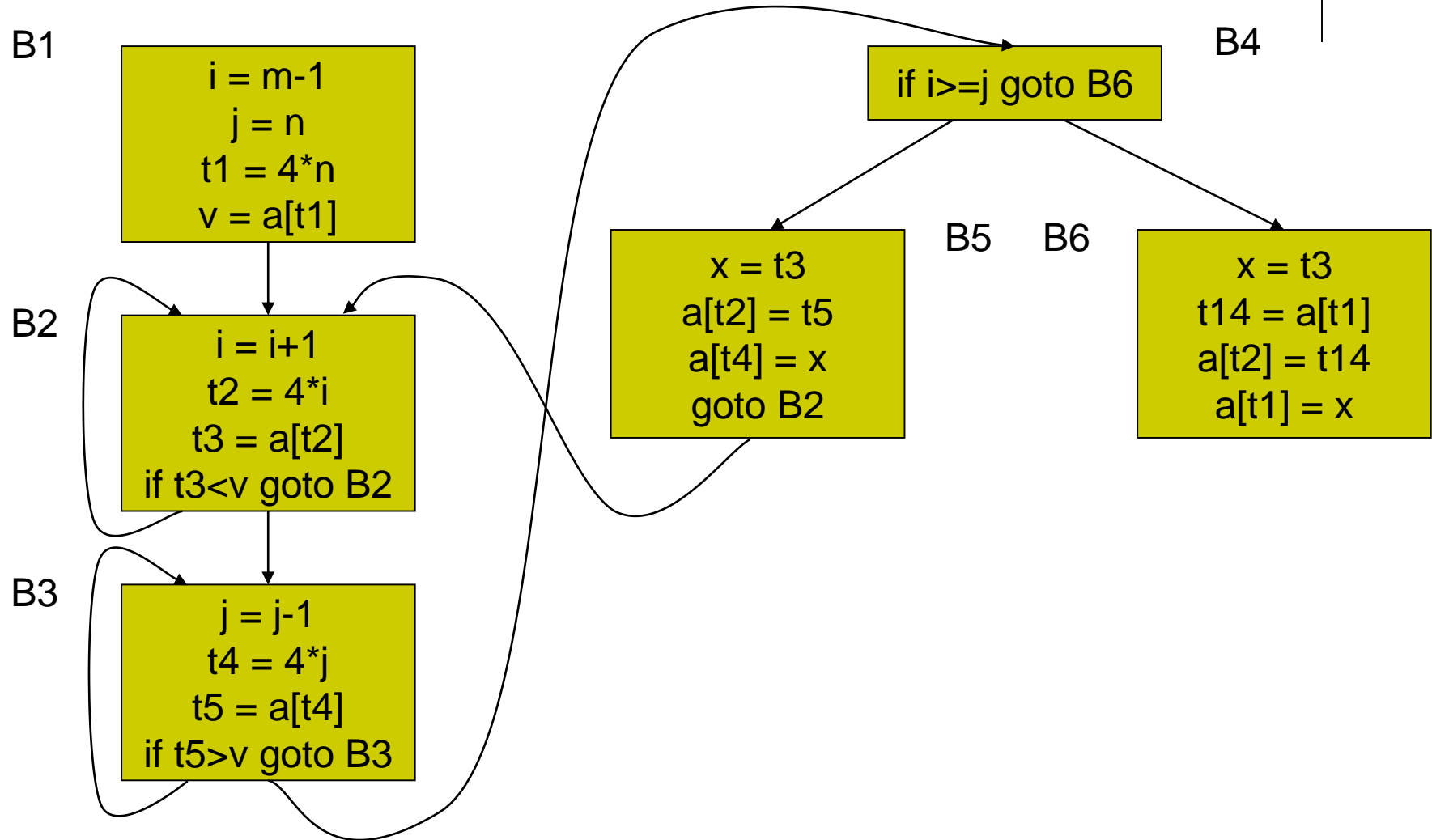
```
t6 = 4*i
x = a[t6]
t7 = 4*i
t8 = 4*j
t9 = a[t8]
a[t7] = t9
t10 = 4*j
a[t10] = x
goto B2
```

B5

```
t6 = 4*i
x = a[t6]
t8 = 4*j
t9 = a[t8]
a[t6] = t9
a[t8] = x
goto B2
```



# Example – GCSE



# Example – copy propagation and dead-code elimination



```
x = t3  
a[t2] = t5  
a[t4] = x  
goto B2
```

B5

```
x = t3  
a[t2] = t5  
a[t4] = t3  
goto B2
```

B5

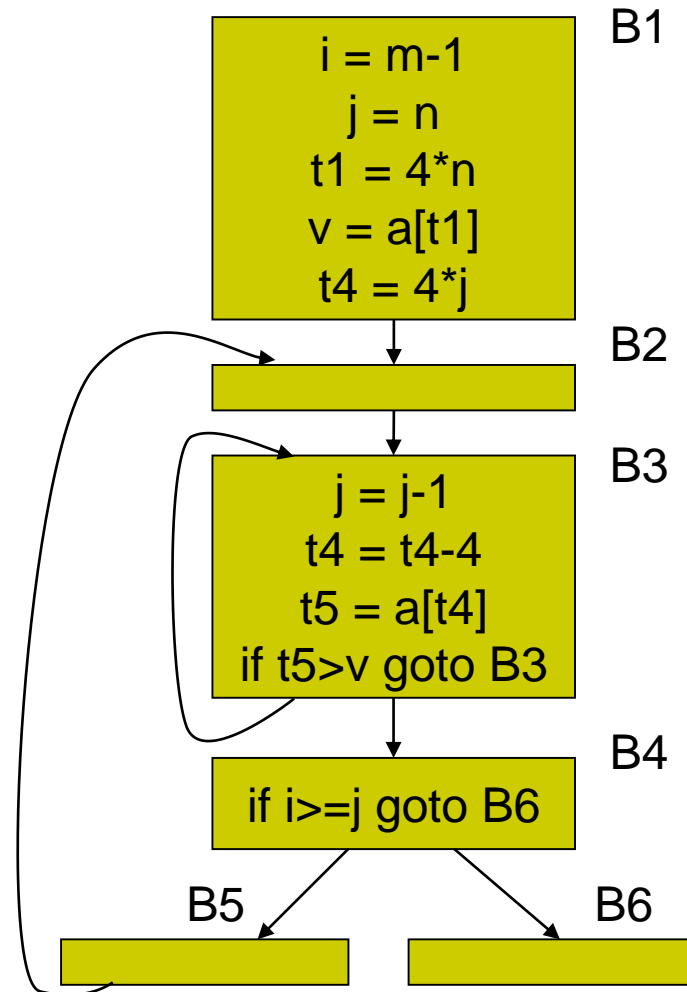
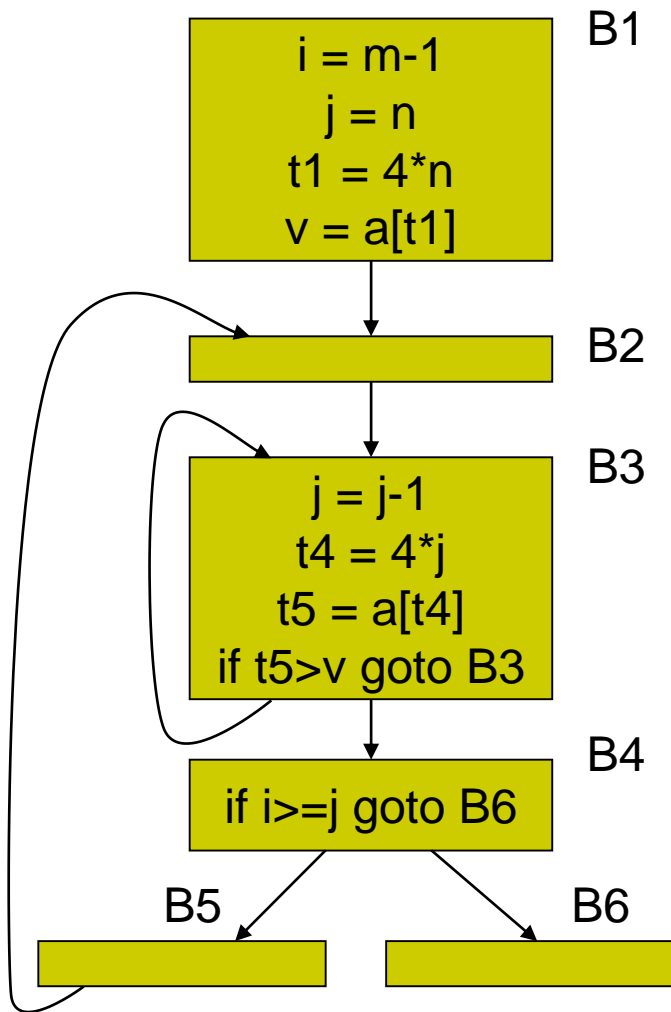
```
x = t3  
a[t2] = t5  
a[t4] = t3  
goto B2
```

B5

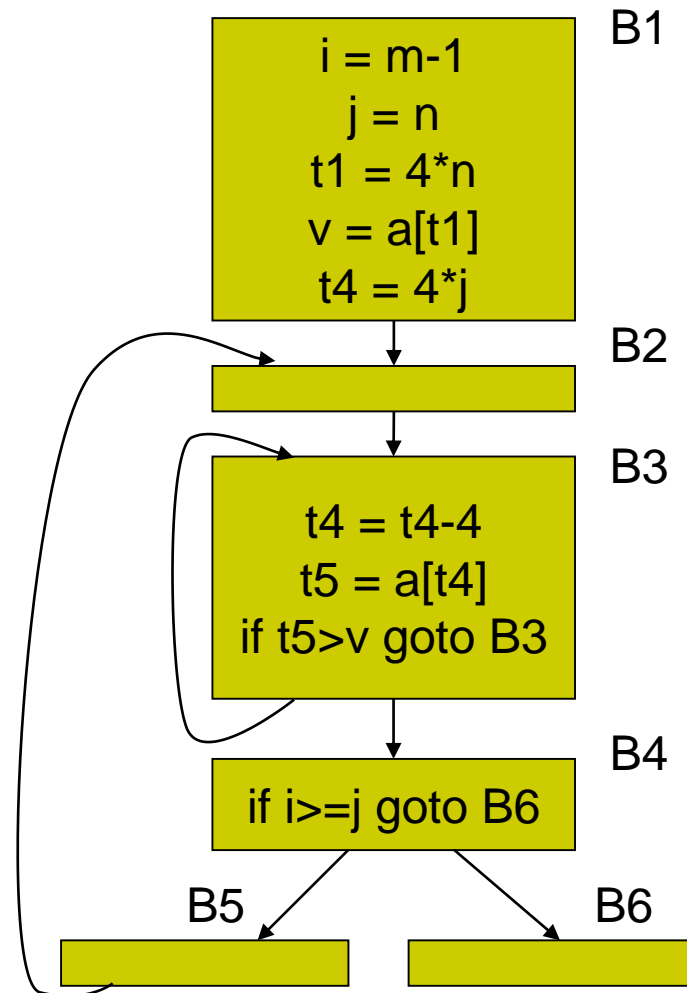
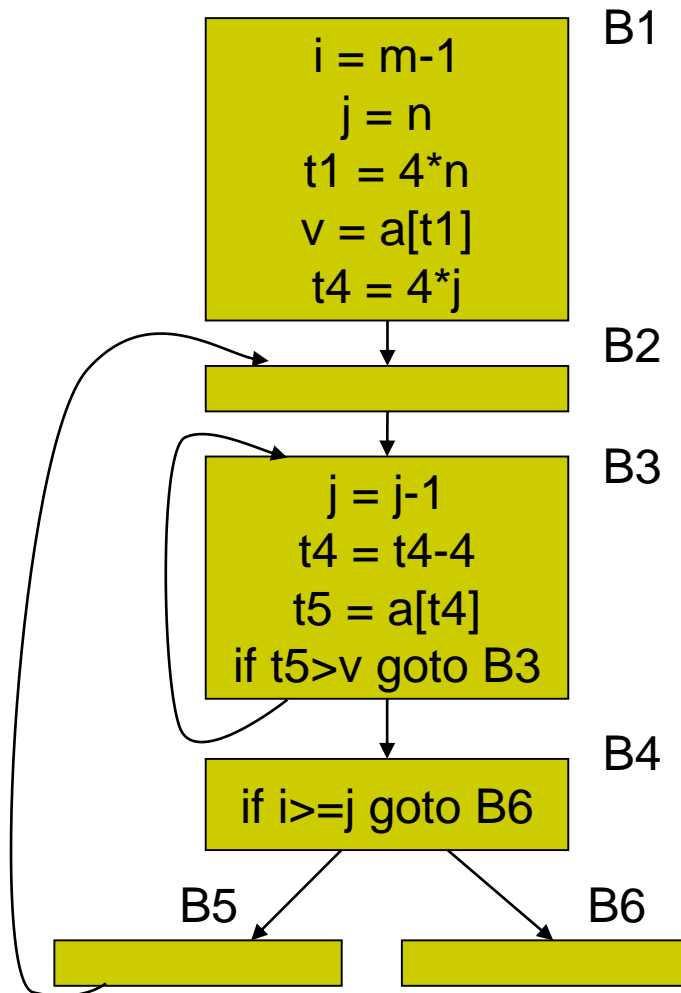
```
a[t2] = t5  
a[t4] = t3  
goto B2
```

B5

# Example – reduction in strength

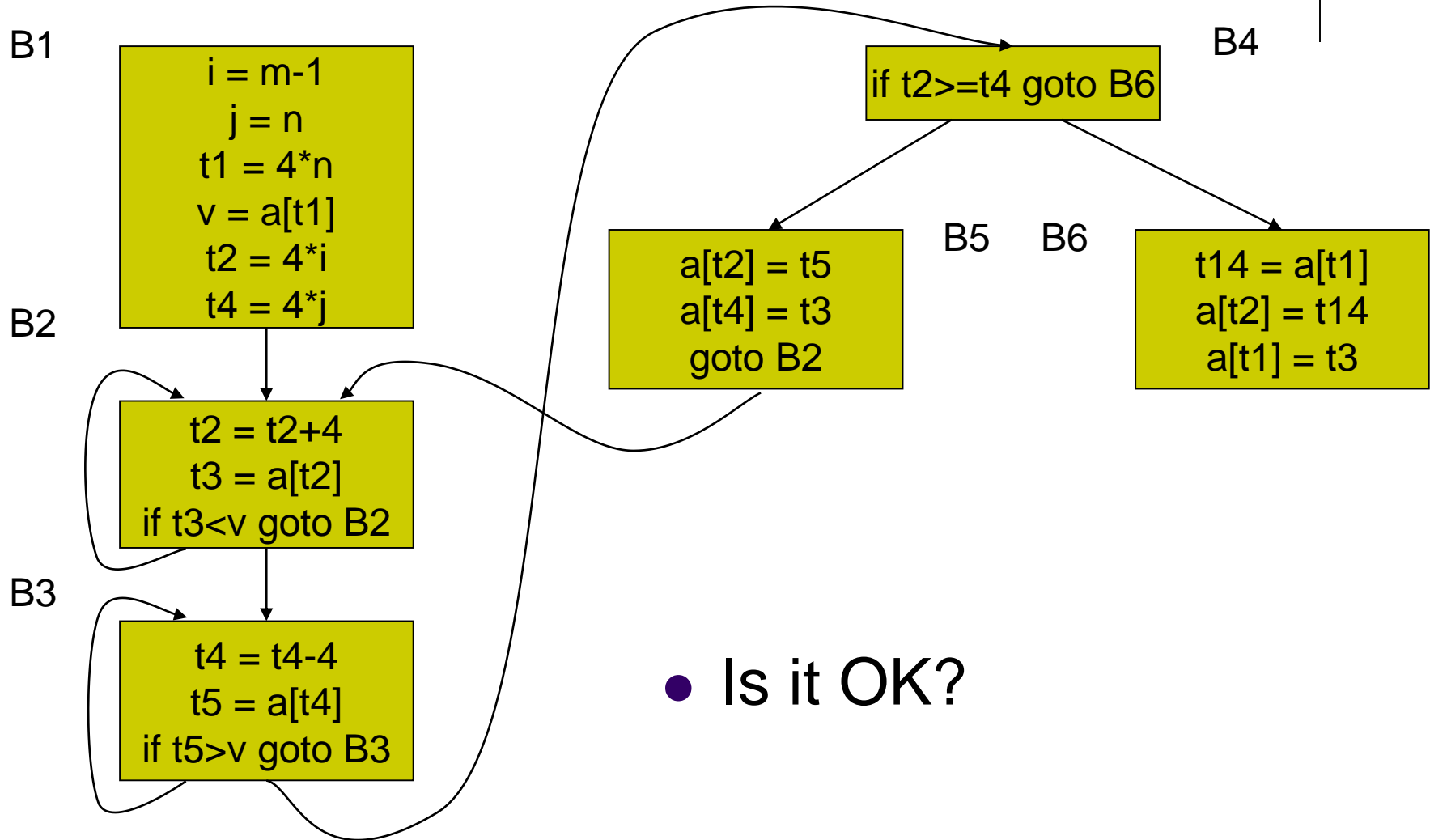


# Example – removing induction variable



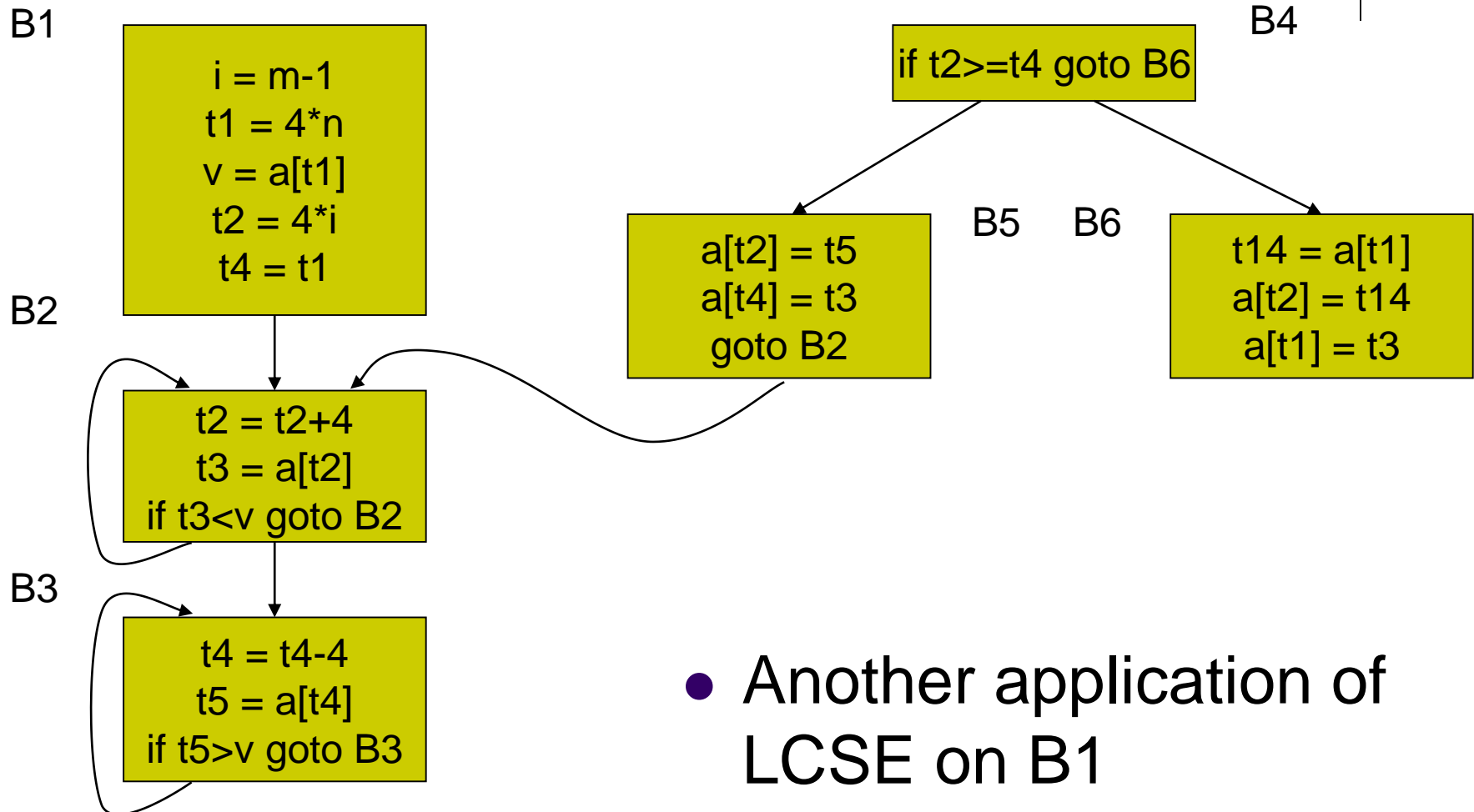


# Example – result





# Example – final result





# Parallelization and vectorization



- Parallelization
  - Implicit × explicit
  - Code parts allowing concurrent execution
  - Prefetch variables to a cache
- Vectorization
  - Expression parallelization using SIMD operations



# Profile Guided Optimization

- Three phases
  - Instrumentation
    - Specially compiled code with calls to “collecting” functions at the start and the end of each basic block
  - Profiling
    - Run the instrumented code with a “typical” input
    - Collect statistical data about visiting collecting functions
  - Optimization
    - Collected data influence optimization
    - Used in optimization upon intermediate code and code
    - Adding weights to edges in control-flow graph
    - Looking for most important paths in control-flow graph
    - Register assignment

# Compiler principles

---

Processor architectures

Jakub Yaghob





# Processor architectures

- Processor architecture is a target language
  - For compilers targeting a processor code
  - It influences the compiler backend, namely code generator
  - It can affect an intermediate code form and instructions, used high-level optimizations



# Registers

- The fastest memory
- Precious resource
  - x86 has 7 of 32-bits integer registers
  - IA-64 has 128 of 64-bits integer registers
- Different register types for different data types
  - Integer, floating point, address, vector
- Different access modes
  - Direct
  - Stack (FPU)



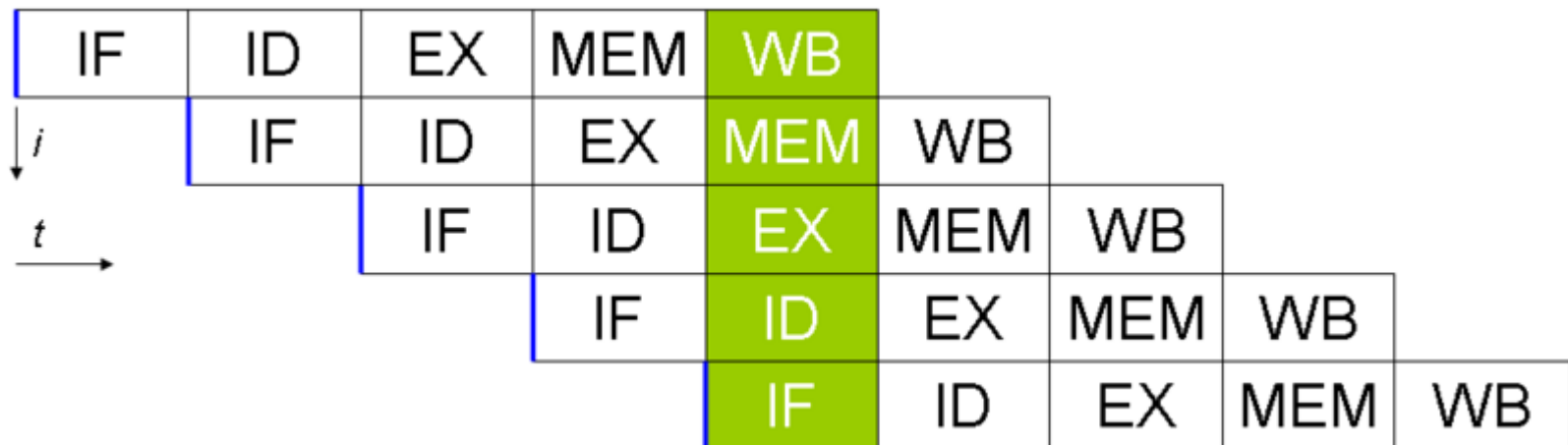
# Instruction set

- RISC
  - Simple instructions, small number (e.g. no division)
- CISC
  - Complex instructions, wide repertoire
- Load-Execute-Store
- Orthogonality
  - x86 has non-orthogonal instruction set



# Pipelining

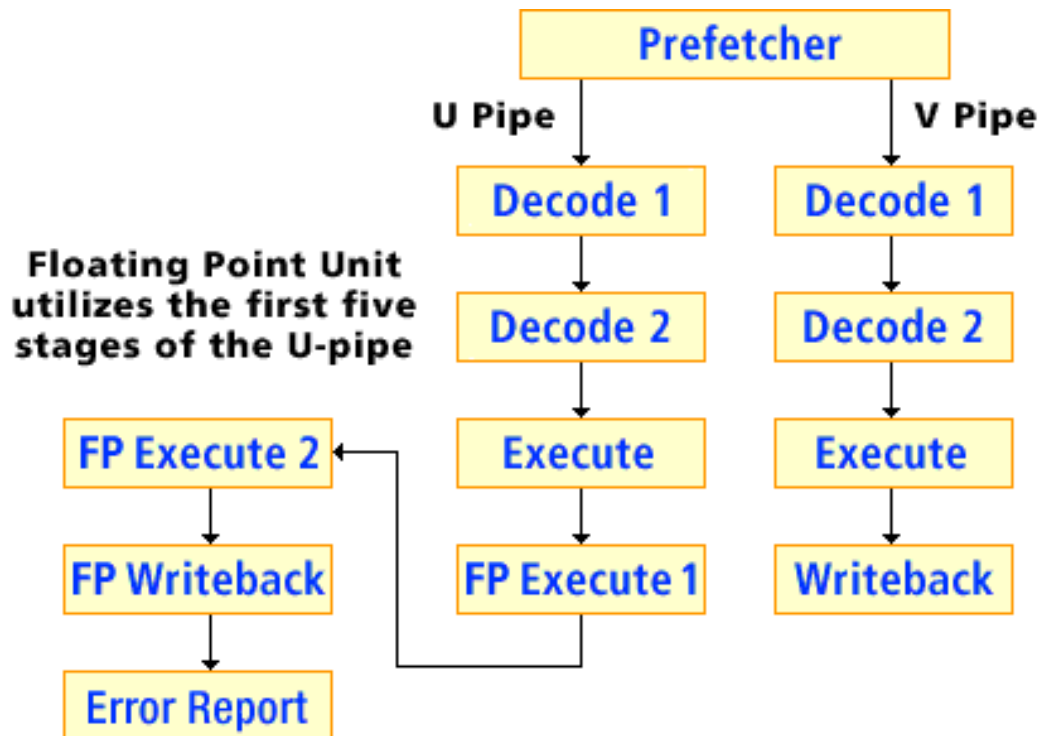
- Next instruction fetch and decode started before the previous instruction was finished
- Each stage and each instruction has a latency
  - Problem with operand dependency (RAW)



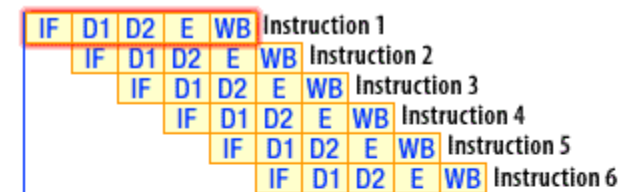


# Superscalar processor

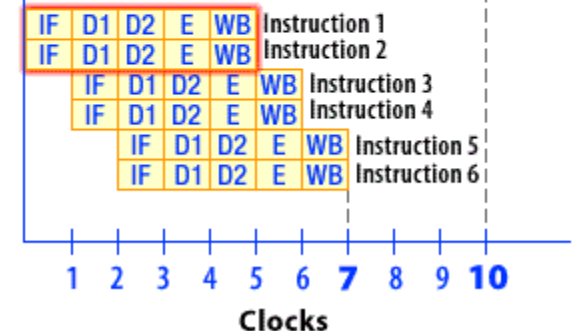
- More equal units capable of parallel instruction execution



Pipelined Scalar Execution



Pipelined Superscalar Execution (Assumes optimum pairing)



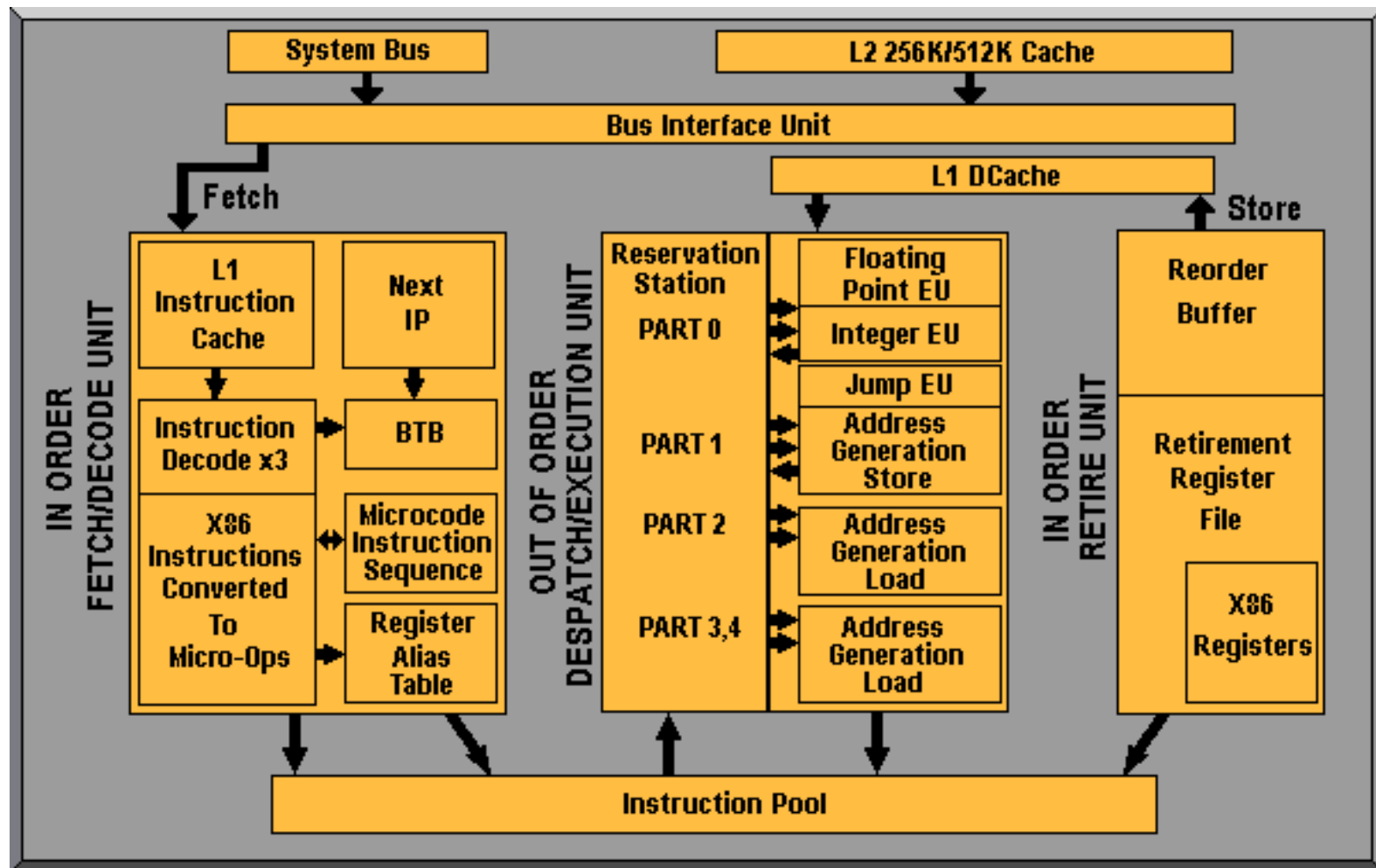




# Out-of-Order Execution

- Instruction fetch
- Instruction dispatch to an instruction queue
- The instruction waits until its input operands are available
- The instruction is issued to its execution unit
- The results are queued
- When all older instructions written their results to the register file, the instruction writes its result

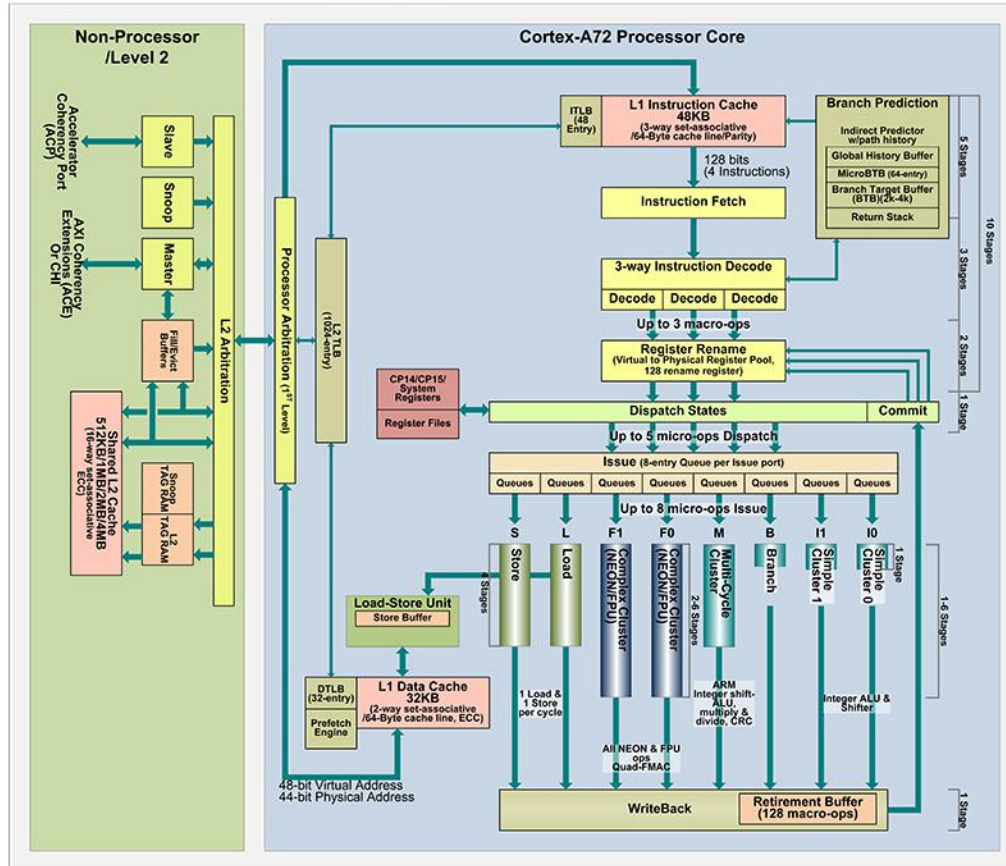
# Out-of-Order Execution – AMD



# Out-of-Order Execution – ARM Cortex-A72



ARM Cortex-A72 Block Diagram





# Branch prediction

- Deep pipelines have a problem with not taken conditional jumps
- Dynamic branch prediction
  - BTB
  - CPU uses history (patterns of some length) for branch prediction
  - Static prediction used when no history is available
- Static branch prediction
  - No hint
    - Forward jump is not taken, backward jump is taken
  - Hint
    - A compiler computes branch probability and generates appropriate branch hint
- Delay slot



# Speculative execution

- Execution of a code in advance which is not used later
- Significant disproportion between CPU and memory speed
- Usually used for read operations in advance
- CPU postpones write operations
- Memory barriers
- Code speculation
  - Check for successful execution
- Data speculation
  - Moved forward even with possible pointer-aliasing



# SIMD instructions

- Sometimes called multimedia instructions
- Integer and float data types
- Expression vectorization already in intermediate code or later in code generation
- Considerable (constant) execution speedup
- Sometimes difficult detection of vectorization opportunity



# VLIW

- Instruction encoding
  - Templates in IA-64
- ILP (Instruction Level Parallelism)
  - Concurrency encoded directly in instructions
  - Dependencies in concurrency group
    - RAW, WAW disallowed
    - WAR allowed

# SoC (System on Chip) + microcontrollers



- Small memory space
  - Optimization for size of code and data
    - One-bit variables
- Separated memory spaces
  - Code and several data address spaces
  - Different access
- Simple pipeline, no superscalarity





# Code generation

- Memory allocation
- Instruction selection
- Register allocation
- Instruction scheduling
- Output
  - Output file format
  - Objects
    - Code, data, relocations, external and public symbols, debug information



# Memory allocation

- Code
  - Jump resolving
- Static data
  - Global data
- Stack
  - Local variables and function parameters
- Size of data types
- Memory placement
- Data alignment
- Decide, what is in memory and what gets an allocated register



# Instruction selection

- 1:N
  - One intermediate code instruction corresponds with N target instructions
  - Simple, the generated code is considerably suboptimal
- M:N
  - M intermediate code instructions correspond with N target instructions
  - NP-complete problem of selection
  - Heuristics
- More possibilities of code generation
  - CISC, SIMD, VLIW
- Non-orthogonal instruction set problem



# Register allocation

- What is placed to the memory and what is placed into a register of a required register type
- Non-orthogonal instruction set problem
- Coloring of a graph of variable liveness
  - The graph is constructed only for a selected variable type (e.g. integer variables)
  - The number of colors is equal to the number of available registers of the selected type
  - When there is no graph coloring, remove “less important” graph nodes and try again
  - We don't need to find an optimal graph coloring, we are looking for any graph coloring



# Instruction scheduling

- Rearrange instructions keeping semantic and avoiding pipeline stalls
  - RAW, WAW, WAR dependencies
  - Construct graph of data dependency
    - Oriented graph, where an edge represents instruction ordering for data availability
    - Find any topological order
  - Instruction latency, speculations, superscalarity, out-of-order execution
  - Usually performed upon one basic block

# Compiler principles

---

Processor architectures

Jakub Yaghob





# Processor architectures

- Processor architecture is a target language
  - For compilers targeting a processor code
  - It influences the compiler backend, namely code generator
  - It can affect an intermediate code form and instructions, used high-level optimizations



# Registers

- The fastest memory
- Precious resource
  - x86 has 7 of 32-bits integer registers
  - IA-64 has 128 of 64-bits integer registers
- Different register types for different data types
  - Integer, floating point, address, vector
- Different access modes
  - Direct
  - Stack (FPU)





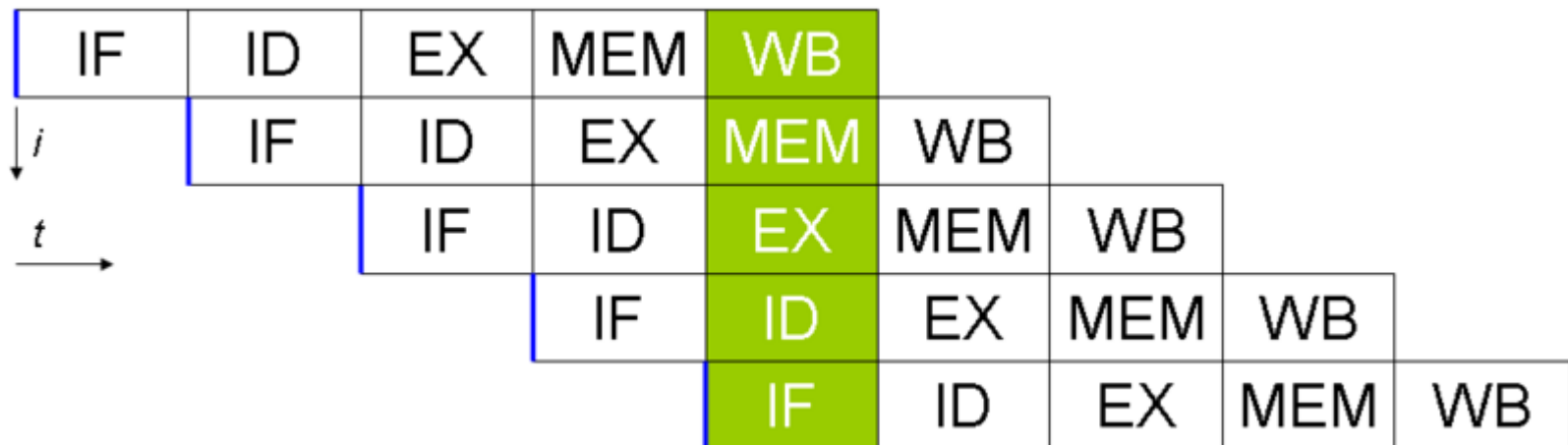
# Instruction set

- RISC
  - Simple instructions, small number (e.g. no division)
- CISC
  - Complex instructions, wide repertoire
- Load-Execute-Store
- Orthogonality
  - x86 has non-orthogonal instruction set



# Pipelining

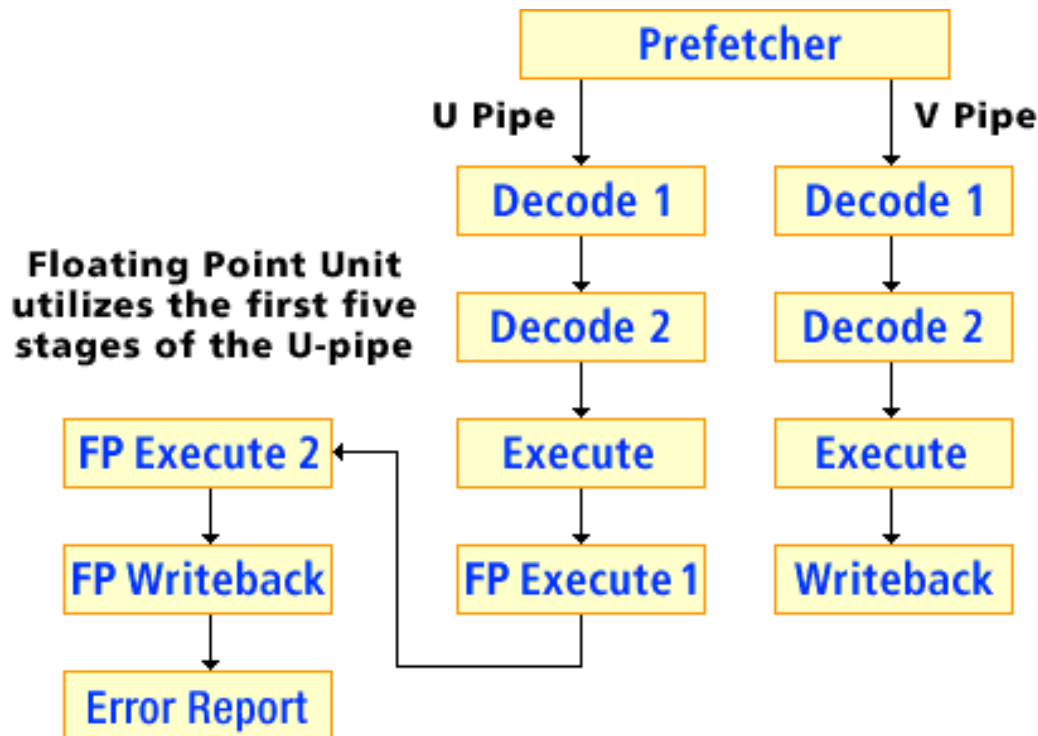
- Next instruction fetch and decode started before the previous instruction was finished
- Each stage and each instruction has a latency
  - Problem with operand dependency (RAW)





# Superscalar processor

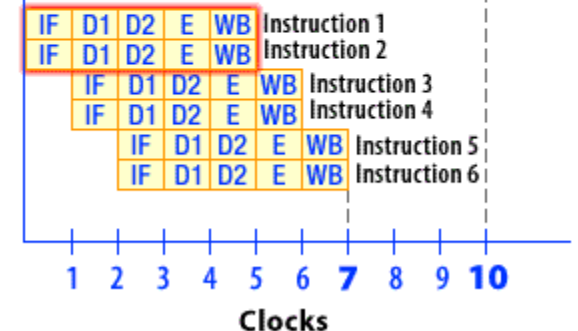
- More equal units capable of parallel instruction execution



Pipelined Scalar Execution



Pipelined Superscalar Execution (Assumes optimum pairing)

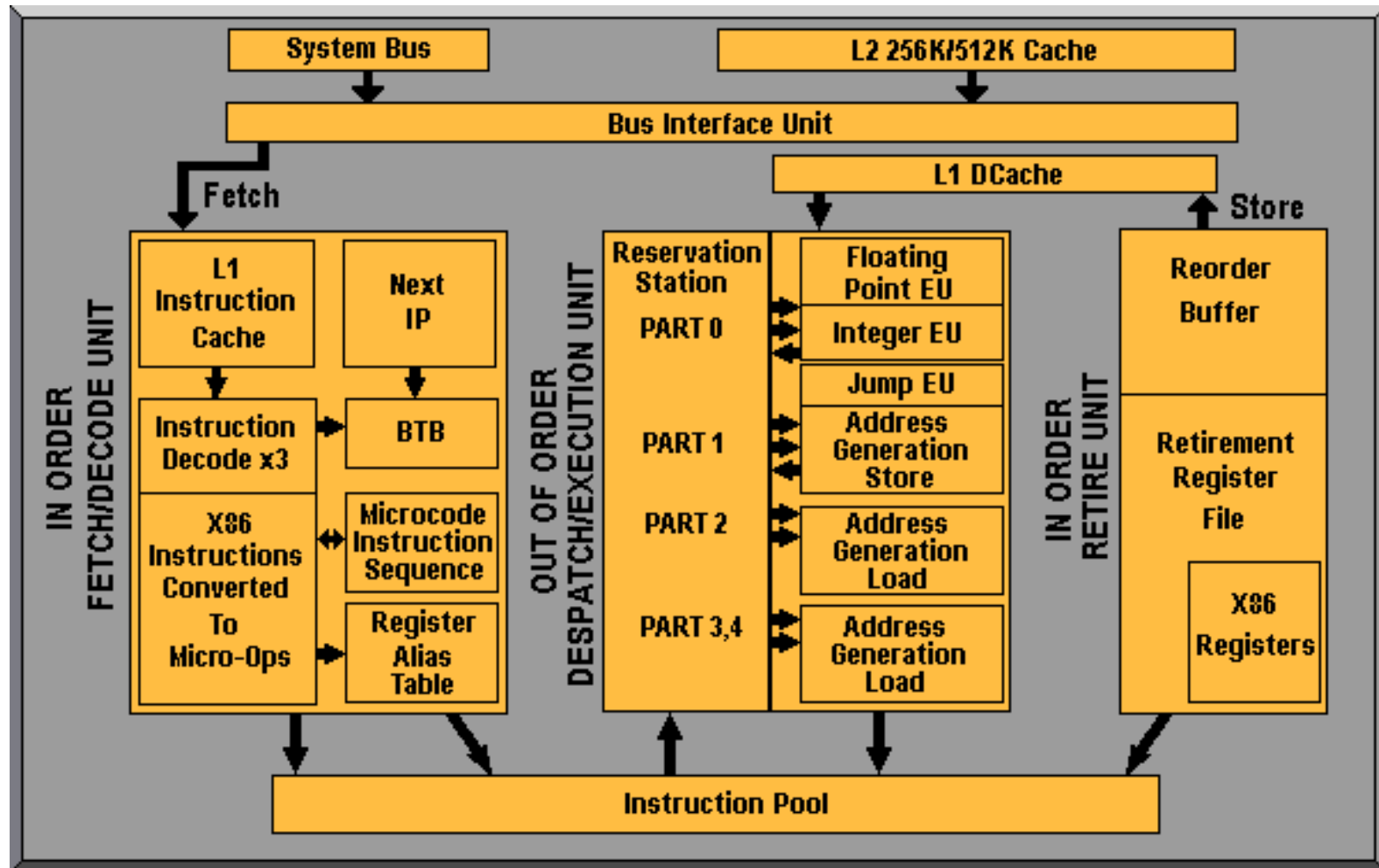




# Out-of-Order Execution

- Instruction fetch
- Instruction dispatch to an instruction queue
- The instruction waits until its input operands are available
- The instruction is issued to its execution unit
- The results are queued
- When all older instructions written their results to the register file, the instruction writes its result

# Out-of-Order Execution – AMD



The diagram illustrates the Cortex-A72 Processor Core architecture, divided into three main sections: Non-Processor/Level 2, the Processor Core itself, and the L2 Arbitration and Cache components.

### Non-Processor/Level 2

- Accelerator Controller (ACP) Port**: Connects to the **Slave** component.
- AXI Coherency Extensions (ACE)**: Connects to the **Master** component.
- Slave** and **Master** components are connected to the **L2 Arbitration** block.
- Shared L2 Cache** (512KB/1MB/2MB/4MB) with **512-entry ECC** is connected to the **L2 Arbitration** block.
- Snoop** and **TAG RWM** components are connected to the **Shared L2 Cache**.

### Cortex-A72 Processor Core

The core is organized into stages, with a total of 10 stages indicated on the right side of the diagram.

- Branch Prediction** (5 Stages): Includes Indirect Predictor w/path history, Global History Buffer, MicroBTB (64-entry), Branch Target Buffer (BTB) (2k-4k), and Return Stack.
- Instruction Fetch** (3 Stages): Includes **ITLB** (48 Entry) and **L1 Instruction Cache** (48KB, 3-way set-associative /64-Byte cache line/Parity). It feeds into a **128 bits (4 Instructions)** buffer.
- 3-way Instruction Decode**: Feeds into **Decode** blocks, which then feed into **Up to 3 macro-ops**.
- Register Rename** (Virtual to Physical Register Pool, 128 rename register): Feeds into **Dispatch States**.
- Dispatch States**: Feeds into **Up to 5 micro-ops Dispatch**.
- Commit**: Feeds into **Up to 5 micro-ops Dispatch**.
- Issue** (8-entry Queue per issue port): Feeds into **Up to 8 micro-ops Issue**.
- Up to 8 micro-ops Issue**: Feeds into various execution units:
  - S** (Store): Feeds into **Store** and **Load** units.
  - L** (Load): Feeds into **Store** and **Load** units.
  - F1** (Complex Cluster (NEON/FPU))
  - F0** (Complex Cluster (NEON/FPU))
  - M** (Multi-Cycle Cluster)
  - B** (Branch)
  - I1** (Simple Cluster 1)
  - I0** (Simple Cluster 0)
- Execution Units**:
  - Store** and **Load** units feed into the **Load-Store Unit** (Store Buffer).
  - F1** and **F0** feed into **AI/NEON & FPU ops Quad/MAC**.
  - M** and **B** feed into **ARM Integer shift-ALU, multiply & divide, CRC**.
  - I1** and **I0** feed into **Integer ALU & Shifter**.
- WriteBack** and **Retirement Buffer** (128 macro-ops): Receive data from the execution units.

### L2 Arbitration and Cache

- L2 Arbitration**: Connects to the **Shared L2 Cache** and the **Processor Arbitration (L1 Level)**.
- Processor Arbitration (L1 Level)**: Connects to the **L2 TLB** (1024-entry) and the **DTLB** (32-entry).
- L2 TLB** (1024-entry): Connects to the **DTLB** (32-entry).
- DTLB** (32-entry): Connects to the **Prefetch Engine**.
- DTLB** (32-entry) and **Prefetch Engine**: Feed into the **L1 Data Cache** (32KB, 2-way set-associative /64-Byte cache line, ECC).
- L1 Data Cache** (32KB): Feeds into the **WriteBack** and **Retirement Buffer**.

**48-bit Virtual Address** and **44-bit Physical Address** are shown at the bottom of the diagram.

Copyright (c) 2015 Hiroshige Goto All rights reserved.



# Branch prediction

- Deep pipelines have a problem with not taken conditional jumps
- Dynamic branch prediction
  - BTB
  - CPU uses history (patterns of some length) for branch prediction
  - Static prediction used when no history is available
- Static branch prediction
  - No hint
    - Forward jump is not taken, backward jump is taken
  - Hint
    - A compiler computes branch probability and generates appropriate branch hint
- Delay slot



# Speculative execution

- Execution of a code in advance which is not used later
- Significant disproportion between CPU and memory speed
- Usually used for read operations in advance
- CPU postpones write operations
- Memory barriers
- Code speculation
  - Check for successful execution
- Data speculation
  - Moved forward even with possible pointer-aliasing





# SIMD instructions

- Sometimes called multimedia instructions
- Integer and float data types
- Expression vectorization already in intermediate code or later in code generation
- Considerable (constant) execution speedup
- Sometimes difficult detection of vectorization opportunity



# VLIW

- Instruction encoding
  - Templates in IA-64
- ILP (Instruction Level Parallelism)
  - Concurrency encoded directly in instructions
  - Dependencies in concurrency group
    - RAW, WAW disallowed
    - WAR allowed

# SoC (System on Chip) + microcontrollers



- Small memory space
  - Optimization for size of code and data
    - One-bit variables
- Separated memory spaces
  - Code and several data address spaces
  - Different access
- Simple pipeline, no superscalarity



# Code generation

- Memory allocation
- Instruction selection
- Register allocation
- Instruction scheduling
- Output
  - Output file format
  - Objects
    - Code, data, relocations, external and public symbols, debug information



# Memory allocation

- Code
  - Jump resolving
- Static data
  - Global data
- Stack
  - Local variables and function parameters
- Size of data types
- Memory placement
- Data alignment
- Decide, what is in memory and what gets an allocated register



# Instruction selection

- 1:N
  - One intermediate code instruction corresponds with N target instructions
  - Simple, the generated code is considerably suboptimal
- M:N
  - M intermediate code instructions correspond with N target instructions
  - NP-complete problem of selection
  - Heuristics
- More possibilities of code generation
  - CISC, SIMD, VLIW
- Non-orthogonal instruction set problem



# Register allocation

- What is placed to the memory and what is placed into a register of a required register type
- Non-orthogonal instruction set problem
- Coloring of a graph of variable liveness
  - The graph is constructed only for a selected variable type (e.g. integer variables)
  - The number of colors is equal to the number of available registers of the selected type
  - When there is no graph coloring, remove “less important” graph nodes and try again
  - We don't need to find an optimal graph coloring, we are looking for any graph coloring



# Instruction scheduling

- Rearrange instructions keeping semantic and avoiding pipeline stalls
  - RAW, WAW, WAR dependencies
  - Construct graph of data dependency
    - Oriented graph, where an edge represents instruction ordering for data availability
    - Find any topological order
  - Instruction latency, speculations, superscalarity, out-of-order execution
  - Usually performed upon one basic block



# Compiler principles

---

Run-time support

Jakub Yaghob





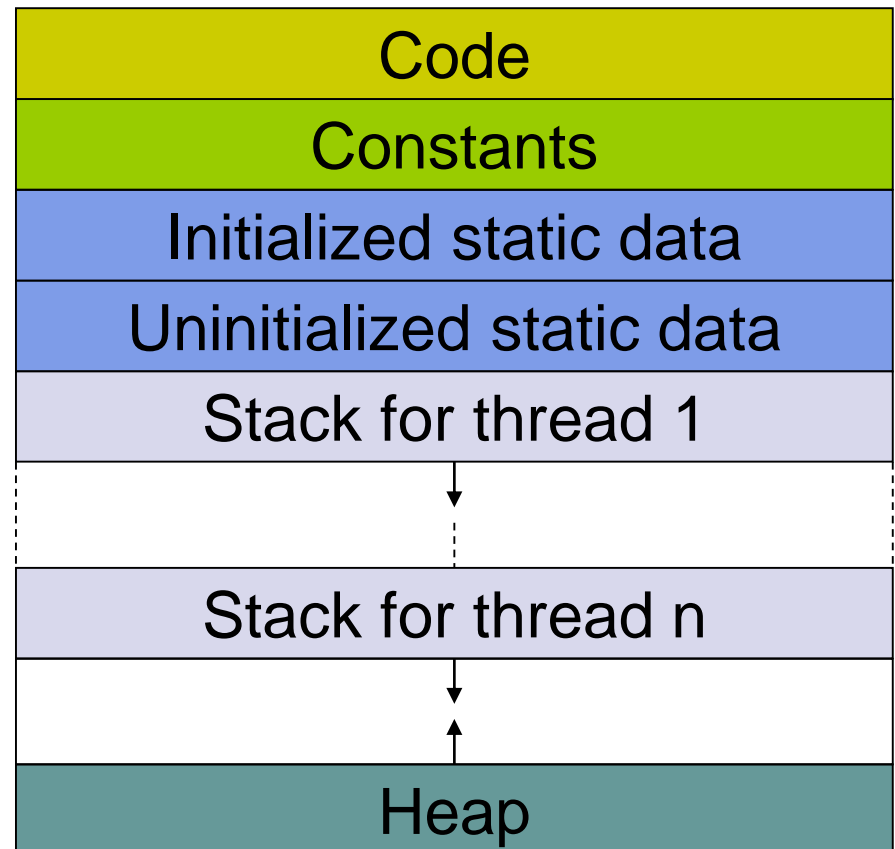
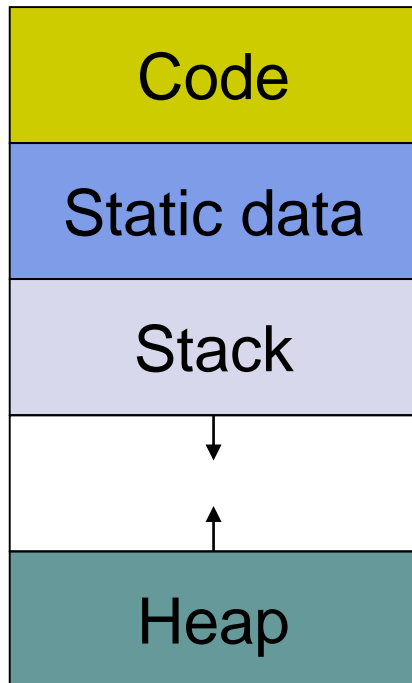
# Run-time support

- Static language support
  - Compiler
  - Library interface
    - Header files
- Dynamic language support
  - Run-time program environment
    - Storage organization
    - Memory content before execution
    - Constructors and destructors of global objects
  - Libraries
  - Calling convention



# Memory organization

- Memory organization during procedural program execution





# Activation record

Return value
Actual parameters
Control link
Access link
Saved machine status
Local data
Temporaries

- Control link
  - Activation record of the caller
- Access link
  - Pointer to nonlocal data held in other activation records
- Saved machine status
  - Return address to the code
  - Registers



# Calling convention

- Calling convention
  - Public name mangling
  - Call/return sequence for functions and procedures
    - Housekeeping responsibility
  - Parameter passing
    - Registers, stack
    - Order of passed parameters
  - Return value
    - Registers, stacks
  - Registers role
    - Parameter passing, scratch, preserved



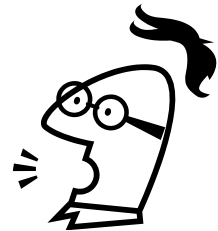
# Public name mangling

- Real meaning
  - mangle
    - mandlovat
    - rozsekat, roztrhat, rozbít, rozdrtit, těžce poškodit, potlouci, pohmoždit
    - *přen.* pokazit, **z**netvořit, **k** nepoznání **z**měnit, překroutit, zkomolit
- Examples:

```
long f1(int i, const char *m, struct s *p)
```

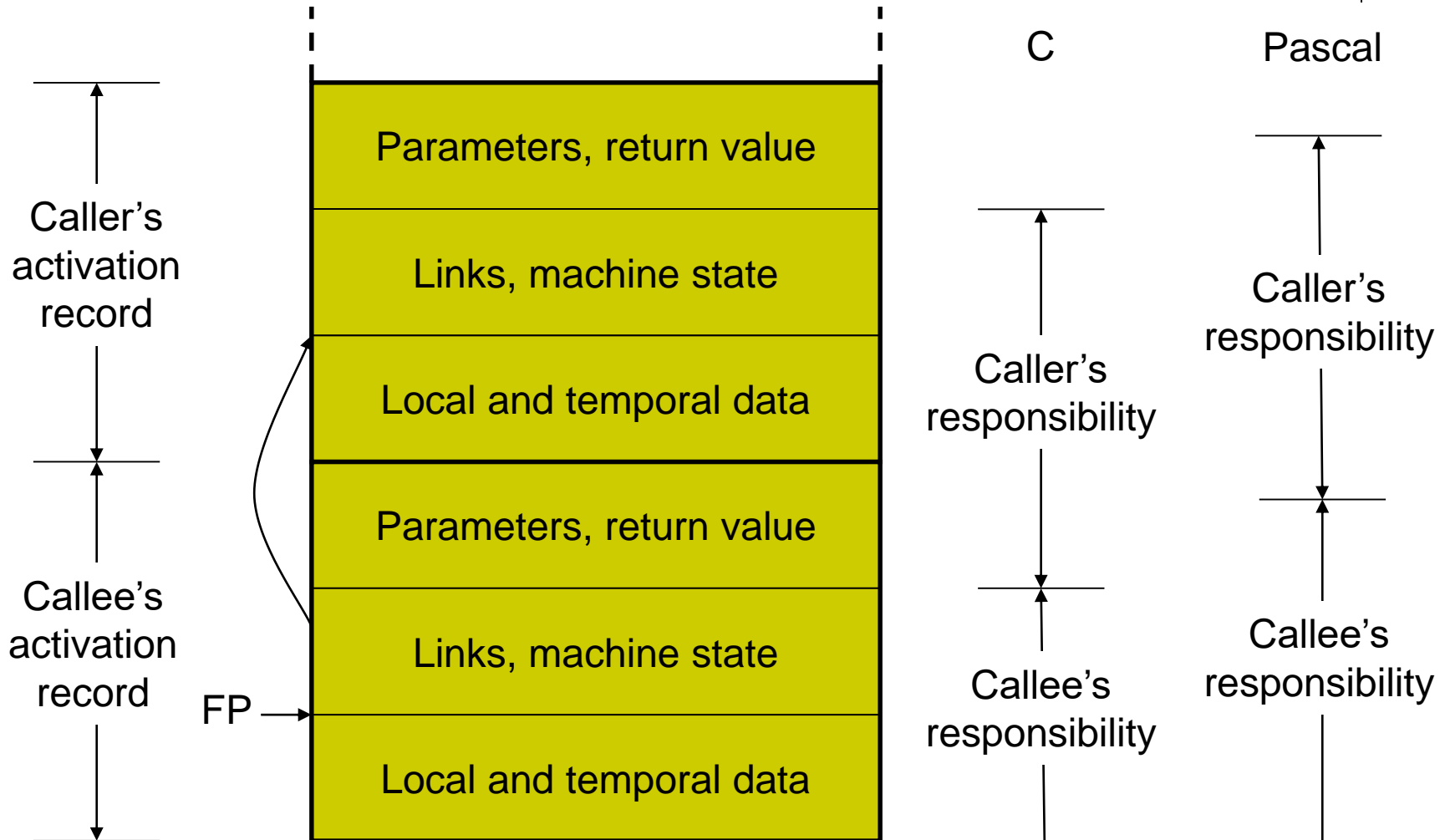
```
__f1  
@f1@12  
__f1@12  
?f1@@YAJHPBDPAUs@@@Z  
__f1  
__Z2f1iPKcPls  
f1  
?f1@@YAJHPEBDPEAUs@@@Z
```

```
MSVC IA-32 C __cdecl  
MSVC IA-32 C __fastcall  
MSVC IA-32 C __stdcall  
MSVC IA-32 C++  
GCC IA-32 C  
GCC IA-32 C++  
MSVC IA-64 C  
MSVC IA-64 C++
```





# Call/return sequence

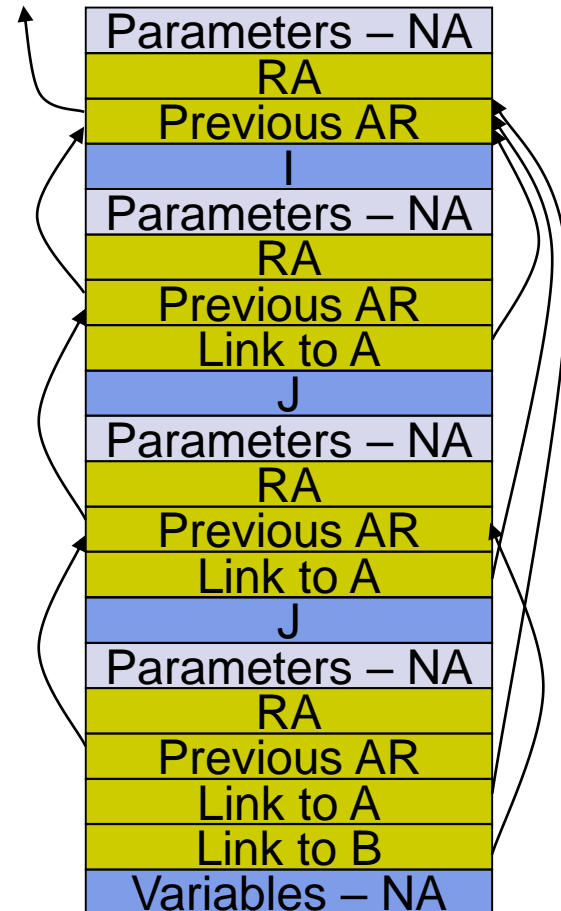


# Access to semantically superior variables



- Pascal
  - Nested functions

```
procedure A;  
var I:integer;  
  procedure B;  
    var J:integer;  
      procedure C;  
        begin ...I+J... end;  
      begin  
        if I>0 B else C;  
      end;  
    begin B end;  
  end;  
begin A end;
```







# Parameter passing

- Call by value
  - Actual parameter is evaluated and the value is passed
  - Input parameters, the parameter is like a local variable
  - C, non-VAR parameters in Pascal
- Call by reference
  - The caller passes a pointer to the storage
  - Input/output parameters
  - & in C++, VAR parameters in Pascal
- Call by name
  - Like a macro – actual expression substituted at the point of use



# Dynamic memory

- Allocation algorithms
  - Continuous blocks of variable length
- Garbage collector
  - Implicit deallocation
  - Reference counting
  - Markers
    - Stop the program at some point of execution
    - All pointers must be known including types
    - All blocks marked as unused
    - Go recursively through pointers and mark used blocks
    - Unused blocks are removed
    - Possible memory consolidation

# Compiler principles

---

Interpreted languages

Jakub Yaghob





# What are we talking about?

- The source language is not translated to a CPU code. Instead, it is translated to an abstract machine code
- Native compiled interpreter simulates the abstract machine



# Why?

- No space for a compiler
  - 8-bit computers and BASIC
- Portability
  - The same abstract machine can run on different OSs and different CPU architectures
  - AS/400, Java
- Security
  - Better control of executed instructions



# Problems

- Speed
  - It can be partially solved by JIT (Just-In-Time compilation)
    - Whenever the interpreter should execute abstract machine code, it is immediately compiled to native CPU instructions and stored in a cache
  - AoT (Ahead-of-Time)
    - Compile during installation
- Portability
  - Changing abstract machine behavior causes problems with portability
    - Java
- How to design the abstract machine/code
  - It should cover all source languages
    - .NET



# Dynamic memory

- When supported, always with garbage collector
  - Pointers under control
  - Easy programming
  - Faster dynamic memory
    - Program usually does not need too much memory, therefore GC is not invoked at all and all memory is just quickly allocated
    - The abstract machine simulator usually takes more memory than would have