Syntactic analysis

Input: token sequence

Tasks:

Parsing: construct derivation according to **concrete syntax**, Tree construction according to **abstract syntax**, Error handling (detection, message, recovery)

Result: abstract program tree

Compiler module parser:

deterministic stack automaton, augmented by actions for tree construction **top-down parsers:** leftmost derivation; tree construction top-down or bottom-up **bottom-up parsers:** rightmost derivation backwards; tree construction bottom-up

Abstract program tree (condensed derivation tree):

represented by a data structure in memory for the translation phase to operate on, linear sequence of nodes on a file (costly in runtime), sequence of calls of functions of the translation phase.

Lecture Compiler I WS 2001/2002 / Slide 37

Objectives:

© 2001 bei Prof. Dr. Uwe Kastens

Relation between parsing and tree construction

In the lecture:

- Explain the tasks, use example on CI-16.
- Sources of prerequisites:
- context-free grammars: "Grundlagen der Programmiersprachen (2nd Semester), or "Berechenbarkeit und formale Sprachen" (3rd Semester),
- Tree representation in prefix form, postfix form: "Modellierung" (st Semester); see CI-5.

Suggested reading:

Kastens / Übersetzerbau, Section 4.1



Objectives:

Distinguish roles and properties of concrete and abstract syntax

In the lecture:

- Use the expression grammar of CI-39, CI-40 for comparison.
- Construct abstract syntax systematically.
- Context-free grammars specify trees not only strings! Is also used in software engineering to specify interfaces.

Suggested reading:

Kastens / Übersetzerbau, Section 4.1

Assignments:

- Generate abstract syntaxes from concrete syntaxes and symbol classes.
- Use Eli for that task. <u>Exercise 10</u>

Questions:

- Why is no information lost, when an expression is represented by an abstract program tree?
- Give examples for semantically irrelevant chain productions outside of expressions.
- Explain: XML-based languages are defined by context-free grammars. Their sentences are textual representations of trees.

© 2001 bei Prof. Dr. Uwe Kastens



Objectives:

Illustrate comparison of concrete and abstract syntax

In the lecture:

- Repeat concepts of "GdP" (slide GdP-2.5): Grammar expresses operator precedences and associativity.
- The derivation tree is constructed by the parser not necessarily stored as a data structure.
- · Chain productions have only one non-terminal symbol on their right-hand side.

Suggested reading:

Kastens / Übersetzerbau, Section 4.1

Suggested reading:

slide GdP-2.5

- How does a grammar express operator precedences and associativity?
- What is the purpose of the chain productions in this example.
- What other purposes can chain productions serve?



Objectives:

Illustrate comparison of concrete and abstract syntax

In the lecture:

- Repeat concepts of "GdP" (slide GdP-2.9):
- Compare grammars and trees.
- Actions create nodes of the abstract program tree.
- Symbol classes shrink node pairs that represent chain productions into one node

Suggested reading:

Kastens / Übersetzerbau, Section 4.1

Suggested reading: slide GdP-2.9

- Is this abstract grammar unambiguous?
- Why is that irrelevant?



Objectives:

Understand the construction schema

In the lecture:

Explanation of the method:

- Relate grammar constructs to function constructs.
- Each function plays the role of an acceptor for a symbol.
- accept function for reading and checking of the next token (scanner).
- Computation of decision sets on CI-42.
- Decision sets must be pairwise disjoint!

Suggested reading:

Kastens / Übersetzerbau, Section 4.2

- A parser algorithm is based on a stack automaton. Where is the stack of a recursive descent parser? What corresponds to the states of the stack automaton?
- Where can actions be inserted into the functions to output production sequences in postfix or in prefix form?

Grammar conditions for recursive descent

A context-free grammar is **strong LL(1)**, if for any pair of productions that have the same symbol on their left-hand sides, the **decision sets are disjoint**:

productions:	A ::= u		A ::= v	
decision sets:	First (u Follow(A))	\cap	First (v Follow(A))	= Ø

First set and follow set:

First (u) := { $t \in T | v \in V^*$ exists and a derivation $u \Rightarrow^* t v$ } and $\varepsilon \in$ First (u) if $u \Rightarrow^* \varepsilon$ exists Follow (A) := { $t \in T | u, v \in V^*$ exist, $A \in N$ and a derivation $S \Rightarrow^* u A v$ such that $t \in$ First (v) }

Example: proc	luction	decision set			
p1: Prog p2: Bloc p3: Decl	::= Block # ::= begin Decls Stmts end ::= Decl : Decls	begin begin new	non-tei	rminal X First(X)	Follow(X)
p4: Decl p5: Decl p6: Stmt p7: Stmt p8: Stmt p9: Stmt	 3 ::= 5 ::= new Ident 5 ::= Stmts ; Stmt 5 ::= Stmt 5 ::= Block 1:= Ident := Ident 	Ident begin new begin Ident begin Ident begin Ident	Prog Block Decls Decl Stmts Stmt	begin begin ε new new begin Ident begin Ident	# ; end Ident begin ; ; end ; end

Lecture Compiler I WS 2001/2002 / Slide 42

Objectives:

Strong LL(1) can easily be checked

In the lecture:

- Explain the definitions using the example.
- First set: set of terminal symbols, which may begin some token sequence that is derivable from u.
- Follow set: set of terminal symbols, which may follow an A in some derivation.

Suggested reading:

Kastens / Übersetzerbau, Section 4.2, LL(k) conditions, computation of First sets and Follow sets

Questions:

The example grammar is not strong LL(1).

- Show where the condition is violated.
- Explain the reason for the violation.
- What would happen if we constructed a recursive descent parser although the condition is violated?

CI-42

Grammar trans	formations for LL(1)	CI-43
Consequences of strong LL(1) condition: A solution alternative productions that begin with • productions that are directly or indirect	strong LL(1) grammar can not ha the same symbols tly left-recursive.	ave
Simple grammar transformations that keep t	the defined language invariant:	
• left-factorization:	non-LL(1) productions	transformed
u, v, w ∈ V^ X ∈ N does not occur in the original grammar	A ::= v u A ::= v w	A ::= v X X ::= u X ::= w
• elimination of direct recursion :	A ::= A u A ::= v	A ::= v X X ::= u X X ::=
EBNF constructs can avoid violation of stro	ong LL(1) condition:	
for example repetition of u: A ::= additional condition: First branch in the function body: v	= v (u)* w :(u)	u} w
correspondingly for EBNF constructs u ⁺ ,	, [u]	

Objectives:

Understand transformations and their need

In the lecture:

- Argue why strong LL(1) grammars can not have such productions.
- Show why the transformations remove those problems.
- Replacing left-recursion by right recursion would usually distort the structure.
- There are more general rules for indirect recursion.
- Show EBNF productions in recursive descent parsers.

- Apply recursion elimination for expression grammars.
- Write a strong LL(1) expression grammar using EBNF.



Objectives:

Understand the decision basis of the automata

In the lecture:

Explain the meaning of the graphics:

- role of the stack: contains states of the automaton,
- accepted input: will not be considered again,
- lookahead: the next k symbols, not yet accepted
- leftmost derivation: leftmost non-terminal is derived next; rightmost correspondingly,
- consequences for the direction of tree construction,

Abbreviations

- LL: (L)eft-to-right, (L)eftmost derivation,
- LR: (L)eft-to-right, (R)ightmost derivation,
- LALR: (L)ook(A)head LR

Suggested reading:

Kastens / Übersetzerbau, Section Text zu Abb. 4.2-1, 4.3-1

Questions:

Use the graphics to explain why a bottom-up parser without lookahead (k=0) is reasonable, but a top-down parser is not.



Objectives:

Fundamental notions of LR automata

In the lecture:

Explain

- meaning of an item,
- lookahead in the input and right context in the automaton.
- There is no right context set in case of an LR(0) automaton.

Suggested reading:

Kastens / Übersetzerbau, Section 4.3

Questions:

• What contains the right context set in case of a LR(3) automaton?



Objectives:

Understand LR(1) states and operations

In the lecture:

Explain

- Sets of items,
- shift transitions,
- reductions.

Suggested reading:

Kastens / Übersetzerbau, Section 4.3

Questions:

• Explain: A state is encoded by a number. A state represents complex information which is important for construction of the automaton.



Objectives:

Example for states, transitions, and automaton construction

In the lecture:

Use the example to explain

- the start state,
- the creation of new states,
- transitions into successor states,
- transitive closure of item set,
- push and pop of states,
- consequences of left-recursive and right-recursive productions,
- use of right context to decide upon a reduction,
- erläutern.

Suggested reading:

Kastens / Übersetzerbau, Section 4.3

- Describe the subgraphs for left-recursive and right-recursive productions. How do they differ?
- How does a LR(0) automaton decide upon reductions?



Objectives:

Understand the method

In the lecture:

Explain using the example on CI-47:

- transitive closure,
- computation of the right context sets,
- relation between the items of a state and those of one of its successor

Suggested reading:

Kastens / Übersetzerbau, Section 4.3

- Explain the role of the right context.
- Explain its computation.

Operations of the LR(1) automaton

shift x (terminal or non-terminal):	Example:		
from current state q under x into the successor state q' ,	stack	input	reduction
push qʻ	1 (a;a;b;b)#	
reduce p:	12	a;a;b;b)#	
apply production p B ::= u.	123	;a;b;b)#	р3
pop as many states.	12	;a;b;b)#	
as there are symbols in u , from the	124	;a;b;b)#	
new current state make a shift with B	1245	a;b;b)#	
	12456	;b;b)#	p2
error:	12	;b;b)#	
the current state has no transition	124	;b;b)#	
under the next input token,	1245	b;b)#	
issue a message and recover	12457	;b)#	
	124578	b)#	
stop:	1245787)#	р5
recuce start production,	124578)#	
see # in the input	1245789)#	p4
	1245)#	
	1 2 4 5 10) #	
	1 2 3 5 10 11	· #	p1
	1	#	•

Lecture Compiler I WS 2001/2002 / Slide 49

Objectives:

© 2001 bei Prof. Dr. Uwe Kastens

Understand how the automaton works

In the lecture:

Explain operations

- Why does the automaton behave differently on a-sequences and b-sequences?
- Which behaviour is better?



Objectives:

Understand LR conflicts

In the lecture:

Explain: In certain situations the given input token t can not determine

- Reduce-reduce: which reduction is to be taken;
- Shift-reduce: whether the next token is to be shifted, a reduction is to be made.

Suggested reading:

Kastens / Übersetzerbau, Section 4.3

- Why can a shift-shift conflict not exist?
- In LR(0) automata items do not have a right-context set. Explain why a state with a reduce item may not contain any other item.



Objectives:

See a conflict in an automaton

In the lecture:

Explain

- the construction
- a solution of the conflict: The automaton can be modified such that in state 6, if an else is the next input token, it is shifted rather than a reduction is made. In that case the ambiguity is solved such that the else part is bound to the inner if. That is the structure required in Pascal and C. Some parser generators can be instructed to resolve conflicts in this way.

Suggested reading:

Kastens / Übersetzerbau, Section 4.3



Objectives:

Understand relations between LR classes

In the lecture:

Explain:

- LALR(1), SLR(1), LR(0) automata have the same number of states,
- compare their states,
- discuss the grammar classes for the example on slide CI-47.

Suggested reading:

Kastens / Übersetzerbau, Section 4.3

Questions:

• Assume that the LALR(1) contruction for a given grammar yields conflicts. Classify the potential reasons using the LR hierarchy.

© 2001 bei Prof. Dr. Uwe Kastens



Objectives:

Implementation of LR tables

In the lecture:

Explanation of

- pair of tables and their entries,
- unreachable entries,
- compression techniques, derived from general table compression,
- Singleton reduction states yield an effective optimization.

- Why are there no error entries in the nonterminal part?
- Why are there unreachable entries?
- Why does a parser need a shift-reduce operation if the optimization of LR(0)-reduction states is applied?



Objectives:

Accept strong requirements

In the lecture:

- The reasons for and the consequences of the requirements are discussed.
- Some of the requirements hold for error handling in general not only that of the syntactic analysis.



Objectives:

Error position from the view of the parser

In the lecture:

Explain the notions with respect to parser actions using the examples.

Questions:

Assume the programmer omitted an opening parenthesis.

- Where is the error position?
- What is the symptom the parser recognizes?



Objectives:

Understand error recovery

In the lecture:

Explain the notions with respect to parser actions using the examples.

Questions:

Assume the programmer omitted an opening parenthesis.

• What could be a suitable repair?

Recovery method: simulated continuation

CI-57

Problem: Determine a continuation point close to the error position and reach it.

Idea: Use parse stack to determine a set of tokens as potential continuation points.

Steps of the method:

- 1. Save the contents of the parse stack when an error is recognized. Skip the error token.
- Compute a set D ⊆ T of tokens that may be used as continuation point (anchor set) Let a modified parser run to completion: Instead of reading a token from input it is inserted into D; (modification given below)
- 3. Find a continuation point d: Skip input tokens until a token of D is found.
- Reach the continuation point d: Restore the saved parser stack as the current stack. Perform dedicated transitions until d is acceptable. Instead of reading tokens (conceptually) insert tokens. Thus a correct prefix is constructed.
- 5. Continue normal parsing.

Augment parser construction for steps 2 and 4:

For each parser state select a transition and its token,

such that the parser empties its stack and terminates as fast as possible.

This selection can be generated automatically.

The quality of the recovery can be improved by influence on the computation of D.

Lecture Compiler I WS 2001/2002 / Slide 57

Objectives:

© 2001 bei Prof. Dr. Uwe Kastens

Error recovery can be generated

In the lecture:

- Explain the idea and the steps of the method.
- The method yields a correct parse for any input!
- Other, less powerful methods determine sets D statically at parser construction time, e. g. semicolon and curly bracket for errors in statements.

Questions:

• How does this method fit to the general requirements for error handling?

			CI-58	
Parser generators				
PGS Cola Lalr Yacc Bison Llgen Deer	Univ. Karlsruhe; in Eli Univ. Paderborn; in Eli Univ. / GMD Karlsruhe Unix tool Gnu Amsterdam Compiler Ki Univ. Colorado, Bouder	EliLALR(1), table-drivenn EliLALR(1), optional: table-driven or directly programmedruheLALR(1), table-drivenLALR(1), table-drivenLALR(1), table-driveniler Kit LL(1), recursive descentouderLL(1), recursive descent		
Form of grammar specification: EBNF: Cola, PGS, Lalr; BNF: Yacc, Bison				
Error recovery: simulated continuation, automatically generated: Cola, PGS, Lalr error productions, hand-specified: Yacc, Bison				
Actions: staten at the anywh	nents in the implementation end of productions: here in productions:	on language	Yacc, Bison Cola, PGS, Lalr	
Conflict r modifi order rules f	r esolution: cation of states (reduce if of productions: for precedence and assoc	:) ciativity:	Cola, PGS, Lalr Yacc, Bison Yacc, Bison	
Implementation languages: C: Cola, Yacc, BisonC, Pascal, Modula-2, Ada: PGS, Lalr				

Objectives:

© 2001 bei Prof. Dr. Uwe Kastens

Overview over parser generators

In the lecture:

• Explain the significance of properties

Suggested reading:

Kastens / Übersetzerbau, Section 4.5