

LL(k): Deterministic top-down parsing

Parsing
ISCL-BA-06

Çağrı Çöltekin
ccolt@informatik.uni-tuebingen.de

University of Tübingen
Seminar für Sprachwissenschaft

Winter Semester 2020/21

version: 4793a39-20201010-01

So far ...

- Formal languages and automata
- General parsing techniques
 - Top-down – Bottom-up
 - Directional – non-directional
- Chart parsing
 - CKY
 - Early

Coming next:

- Deterministic context-free parsing
- Probabilistic context-free parsing
- Dependency parsing

Ç. Çöltekin, ISCL / University of Tübingen

Winter Semester 2020/21 1 / 18

Top-down parsing: recap Recursive descent Table driven parsing 13.1.1

Recap: top-down parsing

- General idea: try to generate the input using the grammar rules
 - Initialize with the start symbol
 - Rewrite each non terminal, replacing them with matching RHS in the grammar
 - When there are multiple options, follow one; backtrack and follow others when done
 - Repeat until input sentence is generated (or failed)
- If we always expand the left-most symbol first, the parser is directional, the resulting derivation is the left-most derivation
- Parsing proceeds with two actions:
predict expanding all RHS of the left-most non-terminal
match if the left-most item is a terminal, it has to match the next input symbol

Ç. Çöltekin, ISCL / University of Tübingen

Winter Semester 2020/21 2 / 18

Top-down parsing: recap Recursive descent Table driven parsing 13.1.2

Top-down parsing: an example

$S \rightarrow NP VP$ $NP \rightarrow d AN$ $NP \rightarrow aN$
 $VP \rightarrow v NP$ $AN \rightarrow a AN$ $AN \rightarrow n$

Matchd	Sent. Form	Input	Action
d n v a n		\$	match n
n	a AN \$	n	match \neq
	a AN VP \$	d n v a n	P: match \neq

Ç. Çöltekin, ISCL / University of Tübingen

Winter Semester 2020/21 3 / 18

Top-down parsing: recap Recursive descent Table driven parsing 13.1.3

Top-down parsing

- If we follow the predicted productions, we obtain a *leftmost* derivation
- Lots of unnecessary work, backtracking because of useless predictions
- Most of the unnecessary work is done in *predict*
- In this lecture we will look at ways to reduce this
- For some grammars, the unnecessary predictions can be completely avoided, resulting in a *deterministic* parser

Ç. Çöltekin, ISCL / University of Tübingen

Winter Semester 2020/21 4 / 18

Top-down parsing: recap Recursive descent Table driven parsing 13.1.3

Recursive descent parser

- Recursive descent parsers are top-down, recursive parsers where each non-terminal is implemented as a procedure
- For each symbol on a RHS, we either
 - call the sub-procedure (another nonterminal)
 - or match the input symbol

```
1: procedure A( )
2:   select a rule  $A \rightarrow X_1 \dots X_k$ 
3:   for  $i = 1$  to  $k$  do
4:     if  $X_i$  is a nonterminal then
5:       call  $X_i()$ 
6:     else if  $X_i =$  current input then
7:       advance the input pointer
8:     else
9:       return error
```

Ç. Çöltekin, ISCL / University of Tübingen

Winter Semester 2020/21 5 / 18

Top-down parsing: recap Recursive descent Table driven parsing 13.1.4

Recursive descent parser

some remarks

- The interesting idea is that now the parser is a program in a(ny) programming language
- In its general form a recursive descent parser is a backtracking parser
- If we can select a rule deterministically, then we can get a deterministic parser
- Deterministic parsing generally requires a *lookahead* mechanism:
 - Given the non-terminal to expand/rewrite, and the next input symbol(s), for some grammars, we can build a table that can deterministically guide a parser

Ç. Çöltekin, ISCL / University of Tübingen

Winter Semester 2020/21 6 / 18

Top-down parsing: recap Recursive descent Table driven parsing 13.1.5

Table driven parsing

$S \rightarrow NP VP$ $NP \rightarrow d AN$ $NP \rightarrow aN$
 $VP \rightarrow v NP$ $AN \rightarrow a AN$ $AN \rightarrow n$

non-term.	input (lookahead)				
	d	a	n	v	\$
S	$S \rightarrow NP VP$	$S \rightarrow NP VP$	$S \rightarrow NP VP$	$S \rightarrow NP VP$	$S \rightarrow NP VP$
NP	$NP \rightarrow d AN$	$NP \rightarrow aN$	$NP \rightarrow aN$	$NP \rightarrow aN$	$NP \rightarrow aN$
VP				$VP \rightarrow v NP$	
AN		$AN \rightarrow a AN$	$AN \rightarrow n$		

Ç. Çöltekin, ISCL / University of Tübingen

Winter Semester 2020/21 7 / 18

Top-down parsing: recap Recursive descent Table driven parsing 13.1.6

Table driven parsing: example

non-term.	input (lookahead)				
	d	a	n	v	\$
S	$S \rightarrow NP VP$	$S \rightarrow NP VP$	$S \rightarrow NP VP$	$S \rightarrow NP VP$	$S \rightarrow NP VP$
NP	$NP \rightarrow d AN$	$NP \rightarrow aN$	$NP \rightarrow aN$	$NP \rightarrow aN$	$NP \rightarrow aN$
VP				$VP \rightarrow v NP$	
AN		$AN \rightarrow a AN$	$AN \rightarrow n$		
Matchd	Sent. Form		Input	Action	
d n v a n			\$	match n	

Ç. Çöltekin, ISCL / University of Tübingen

Winter Semester 2020/21 8 / 18

Top-down parsing: recap Recursive descent Table driven parsing 13.1.7

FIRST and FOLLOW sets

- FIRST and FOLLOW sets are useful for both top-down and bottom-up table driven parsers
- FIRST set of a non-terminal A, FIRST(A), is the set of initial terminal symbols of all strings generated by A
- FOLLOW set of a non-terminal A, FOLLOW(A), is the set of initial terminals that may follow any A according to the grammar
- Both sets generalize to any sentential form
- FIRST and FOLLOW sets are also useful for error recovery during parsing



Ç. Çöltekin, ISCL / University of Tübingen

Winter Semester 2020/21 9 / 18

Top-down parsing: recap Recursive descent Table driven parsing 13.1.8

Computing the FIRST set

- The FIRST set of a terminal symbols contains only itself
- To compute the FIRST sets of nonterminals, repeat the following until no new symbols are added to any of the sets
 - For each rule $X \rightarrow Y_1 Y_2 \dots Y_k$ in the grammar,
 - place all terminals in FIRST(Y_i) if $Y_1 Y_2 \dots Y_{i-1} \Rightarrow \epsilon$
 - if ϵ is in all FIRST(Y_i) for all $i = 1, \dots, k$, add ϵ to FIRST(X)
 - if the rule processed is $X \rightarrow \epsilon$, add ϵ to FIRST(X)
- Then, FIRST set of any sentential form, FIRST($X_1 X_2 \dots X_k$) can be computed:
 - For $i = 1, \dots, k$
 - 1. Add all non- ϵ symbols from X_i to FIRST($X_1 X_2 \dots X_k$)
 - 2. If $\epsilon \in$ FIRST(X_i), stop
 - if $\epsilon \in$ FIRST(X_i) for all $i = 1, \dots, k$, add ϵ to FIRST($X_1 X_2 \dots X_k$)

Ç. Çöltekin, ISCL / University of Tübingen

Winter Semester 2020/21 10 / 18

Top-down parsing: recap Recursive descent Table driven parsing 13.1.9

Computing the FOLLOW set

- Calculate the FIRST sets
- Place ϵ in the FOLLOW(S)
- For a production $A \rightarrow \alpha \beta$, add everything in FIRST(β) except ϵ to FOLLOW(A)
- For a production $A \rightarrow \alpha \beta$ or $A \rightarrow \alpha \beta \gamma$ where FIRST(β) contains ϵ , add all items in FOLLOW(A) to FOLLOW(B)
- Repeat 3 until no more items are added to any of the FOLLOW sets

Ç. Çöltekin, ISCL / University of Tübingen

Winter Semester 2020/21 11 / 18

LL(1) grammars

- A grammar is called and LL(1) grammar, if we can find a table similar to our example:
 - If there is only a single prediction for each (non-terminal, lookahead) pair, then the grammar is an LL(1) grammar
- L's stand for *Left-to-right* and *Leftmost derivation*, (1) indicates the number of lookahead symbols needed
- If we increase the number of lookahead symbols, we get LL(k) grammars
- LL(k) grammar can be parsed with a top-down parser without backtracking
- Not every context free grammar is LL(k)
- But, programming language grammars are mostly LL(1)

LL(1) grammars

formal definition

- If a grammar is LL(1) then whenever $A \rightarrow \alpha$ and $A \rightarrow \beta$ are two rules in the grammar, then
 - The sets of non-terminals of strings derived from α and β are disjoint
 - Only one (or none) of α and β can derive the empty string
 - If $\beta \rightarrow \epsilon$, α cannot start with a terminal that may follow A
- In other words:
 - FIRST(α) and FIRST(β) are disjoint
 - if ϵ is in FIRST(α), then FIRST(β) and FOLLOW(A) are disjoint sets

Construction of LL(1) table

- If there are no ϵ productions, the table can be easily constructed from the FIRST sets
- Otherwise, after computing FIRST and FOLLOW sets, the following procedure fills the LL(1) table
 - For each rule $A \rightarrow \alpha$ in the grammar
 - For each terminal a in FIRST(α), add $A \rightarrow \alpha$ to table cell $[A, a]$
 - If ϵ is in FIRST(α), then for each terminal b in FOLLOW(A) add $A \rightarrow \alpha$ to table cell $[A, b]$

Example

calculating FIRST sets

$$S \rightarrow BA \quad A \rightarrow aBA \mid \epsilon \quad B \rightarrow CD \quad D \rightarrow bCD \mid \epsilon \quad C \rightarrow cSc \mid d$$

- Repeat until no additions
 - For each $X \rightarrow Y_1Y_2 \dots Y_k$
 - place all terminals in FIRST(Y_i)
 - if $Y_1Y_2 \dots Y_{i-1} \Rightarrow \epsilon$
 - if ϵ is in all FIRST(Y_i) for all $i = 1, \dots, k$, add ϵ to FIRST(X)
 - if the rule processed is $X \rightarrow \epsilon$, add ϵ to FIRST(X)
 - $$\begin{aligned} \text{FIRST}(S) &= \{c, d\} \\ \text{FIRST}(A) &= \{a, c\} \\ \text{FIRST}(B) &= \{c, d\} \\ \text{FIRST}(C) &= \{c, d\} \\ \text{FIRST}(D) &= \{b, c\} \end{aligned}$$

Example

calculating FOLLOW sets

$$S \rightarrow BA \quad A \rightarrow aBA \mid \epsilon \quad B \rightarrow CD \quad D \rightarrow bCD \mid \epsilon \quad C \rightarrow cSc \mid d$$

- Place $\$$ in the FOLLOW(S)
- For a production $A \rightarrow \alpha\beta$, add everything in FIRST(β) except ϵ to FOLLOW(B)
- For $A \rightarrow \alpha\beta$, or $A \rightarrow \alpha\beta\gamma$ where FIRST(β) contains ϵ , add items in FOLLOW(A) to FOLLOW(B)

	S	A	B	C	D
FIRST	{c,d}	{a,c}	{c,d}	{c,d}	{b,c}

$$\begin{aligned} \text{FIRST}(S) &= \{c, \$\} \\ \text{FIRST}(A) &= \{c, \$\} \\ \text{FIRST}(B) &= \{a, c, \$\} \\ \text{FIRST}(C) &= \{a, b, c, \$\} \\ \text{FIRST}(D) &= \{a, c, \$\} \end{aligned}$$

Example

constructing the LL(1) table

$$S \rightarrow BA \quad A \rightarrow aBA \mid \epsilon \quad B \rightarrow CD \quad D \rightarrow bCD \mid \epsilon \quad C \rightarrow cSc \mid d$$

- For each rule $A \rightarrow \alpha$ in the grammar
 - For each terminal a in FIRST(α), add $A \rightarrow \alpha$ to table cell $[A, a]$
 - If ϵ is in FIRST(α), then for each terminal b in FOLLOW(A) add $A \rightarrow \alpha$ to table cell $[A, b]$

	a	b	c	d
S			BA	BA
A	aBA		ϵ	
B			CD	CD
C			cSc	d
D	ϵ	bCD	ϵ	

Summary

- LL(1) grammars can be parsed deterministically (without backtracking) using top-down parsers
- Like any top-down parser, left-recursion needs additional care
- Not every context free grammar is LL(k), but programming language grammars are mostly LL(1)
- LL(k) parsing is intuitive and relatively easy to construct by hand, but LR(k) grammars (bottom-up, deterministic) are more powerful (next lecture)
- Suggested reading: Grune and Jacobs (2007, ch.8), Aho et al. (2007, Section 4.4)

Next:

- Deterministic bottom-up parsing
- Suggested reading: Grune and Jacobs (2007, ch.9), Aho et al. (2007, Section 4.5–4.7)

Acknowledgments, references, additional reading material

- Aho, Alfred V., Monica S. Lam, Ravi Sethi, and Jeffrey D. Ullman (2007). *Compilers: Principles, Techniques, & Tools*. Pearson/Addison-Wesley, now Pearson Education.
- Grune, Dirk, and Cordelia Jacobs (2007). *Parsing Techniques: A Practical Guide*, second. *Monographs in Computer Science*. The book online is available at https://algorithms.wtf/Book/PTAPG_for_Acknowledgments.pdf. Springer New York, now 9780307095800.

Exercise

compute the FIRST and FOLLOW sets, and LL(1) table for $S \rightarrow iESQ \mid a \quad Q \rightarrow eS \mid \epsilon \quad E \rightarrow b$

