

COMP4403 Crib Sheet

Parsing Theory

Context Free Grammars

- A context-free grammar consists of:
- A finite set, Σ , of terminal symbols.
 - A finite nonempty set of nonterminal symbols (disjoint from the terminal symbols).
 - A finite nonempty set of productions of the form $A \rightarrow \alpha$, where A is a nonterminal symbol, and α is a possibly empty sequence of symbols, each of which is either a terminal or nonterminal symbol.
 - A start symbol that must be a nonterminal symbol

Nullable

- A possibly empty sequence of symbols, α , is nullable if $\alpha \xRightarrow{*} \epsilon$ or $\alpha \xRightarrow{*} \text{Nullable rules}$
- ϵ is nullable
 - any terminal symbol is not nullable
 - a sequence of symbols is nullable if all of its constructs are nullable
 - a set of alternatives is nullable if any of its constructs are nullable
 - EBNF constructs for optionals and repetitions are nullable
 - a nonterminal is nullable if there is a production with a nullable right-hand side

Left and right associative operators

To remove the ambiguity and treat "-" as a left associate operator (as usual) we can rewrite the grammar to

$$\begin{aligned} E &\rightarrow E - T \\ E &\rightarrow T \\ T &\rightarrow N \end{aligned}$$

and to treat "-" as right associative we use

$$\begin{aligned} E &\rightarrow T - E \\ E &\rightarrow T \\ T &\rightarrow N \end{aligned}$$

Dangling Else

Dangling else is a problem that arises when we have an if-then-else statement and the else is ambiguous as to which if it belongs to.

Recursive-Descent Parsing

Translating EBNF to BNF

Replace optional construct $[S]$ by a new nonterminal $OptS$:

$$\begin{aligned} OptS &\rightarrow S|\epsilon \\ \text{For example, we can rewrite} \\ RelCondition &\rightarrow Exp [RelOp Exp] \\ \text{as} \\ RelCondition &\rightarrow Exp OptRelExp \\ OptRelExp &\rightarrow RelOp Exp|\epsilon \end{aligned}$$

Replace grouping construct (S) by a new nonterminal $GrpS$:

$$\begin{aligned} GrpS &\rightarrow S \\ \text{For example, we can rewrite} \\ RepF &\rightarrow (TIMES | DIVIDE) Factor RepF |\epsilon \\ \text{as} \\ RepF &\rightarrow TermOp Factor RepF |\epsilon \\ TermOp &\rightarrow TIMES | DIVIDE \end{aligned}$$

Recursive-descent parsing

A recursive-descent parser is a recursive program to recognise sentences in a language:

- Input is a stream of lexical tokens, represented by tokens of type `TokenStream`.
- Each nonterminal symbol N has a method `parseN` that:
 - recognises the longest string of lexical tokens in the input stream, starting from the current token, derivable from N ;
 - as it parses the input it moves the current location forward;
 - when it has finished parsing N , the current token is the token immediately following the last token matched as part of N (which may be at the end-of-file token).

parsing example

Write parsing example here

First & Follow Sets

First Sets

- The first set for a construct α records
- the set of terminal symbols α can start with, and

- if α is nullable, it contains the empty string ϵ to indicate that.

Calculating First Sets

Let

- α be a terminal symbol,
- $\alpha_1, \alpha_2, \dots, \alpha_n$ be strings of (terminal and nonterminal) symbols, and
- A be a nonterminal symbol defined by a single production

$$A \rightarrow \alpha_1|\alpha_2|\dots|\alpha_n,$$

then

$$Fst(\epsilon) = \{\epsilon\}$$

$$Fst(\alpha) = \{\alpha\}$$

$$Fst(\alpha_1|\alpha_2|\dots|\alpha_n) = Fst(\alpha_1) \cup Fst(\alpha_2) \cup \dots \cup Fst(\alpha_n)$$

$$Fst(A) = Fst(\alpha_1|\alpha_2|\dots|\alpha_n)$$

Let S_1, S_2, \dots, S_n be (terminal or nonterminal) symbols, then

$$\begin{aligned} Fst(S_1S_2\dots S_n) = \\ Fst(S_1) - \{\epsilon\} \\ \cup Fst(S_2) - \{\epsilon\} \quad \text{if } S_1 \text{ is nullable} \\ \cup \dots \\ \cup Fst(S_i) - \{\epsilon\} \quad \text{if } S_1-S_{i-1} \text{ is nullable} \\ \cup \dots \\ \cup Fst(S_n) - \{\epsilon\} \quad \text{if } S_1-S_{n-1} \text{ is nullable} \\ \cup \{\epsilon\} \quad \text{if } S_1-S_n \text{ is nullable} \end{aligned}$$

Calculating Algorithmically

Start with the first sets for all nonterminals being the empty set and note that the first set for every terminal symbol, α , is the singleton set $\{\alpha\}$. We then make a pass over all production in a grammar considering all alternatives and process as follows.

- If there is a production of the form $N \rightarrow \epsilon$, we add ϵ to the first set of N .
- If there is a production of the form $N \rightarrow S_1S_2\dots S_n$, then for each $i \in 1..n$, if for all $j \in 1..i-1$, S_j is nullable, we add the current first set for S_i minus ϵ to the first set for N .
- If every construct S_1, \dots, S_n is nullable, we add ϵ to the first set for N .

We repeat the passes until no set is modified in the pass, in which case we are finished.

Example

$$A \rightarrow Bx|C$$

$$B \rightarrow Cy|D$$

$$C \rightarrow Dz|\epsilon$$

$$D \rightarrow Aw$$

A	{ }	{ ϵ }	{ ϵ, y }	{ ϵ, y, w }	{ ϵ, y, w }
B	{ }	{ y }	{ y }	{ y, w }	{ y, w }
C	{ ϵ }	{ ϵ }	{ ϵ, w }	{ ϵ, w, y }	{ ϵ, w, y }
D	{ }	{ w }	{ w, y }	{ w, y }	{ w, y }

Follow Sets

The follow set for a nonterminal, N , is the set of terminal symbols that may follow N in any context within the grammar. End-of-file is represented by the special terminal symbol $\$$ in which a follows N . A nonterminal, N , is followed by a terminal symbol, a , if there is a derivation from $S\$$

Calculating Follow Sets

We compute the Follow set for a nonterminal, N , using two rules.

- If there is a production of the form

$$A \rightarrow \alpha N \beta$$

then any symbols that can start β can follow N , and hence $Follow(N)$ must include all the terminal symbols in $First(\beta)$. Note that ϵ is not included even if it appears in $First(\beta)$.

$$First(\beta) - \{\epsilon\} \subseteq Follow(N)$$

- If there is a production of the form

$$A \rightarrow \alpha N \beta$$

and β is nullable, then any token that can follow A can also follow N . Hence,

$$Follow(A) \subseteq Follow(N)$$

The case where β is nullable includes the case when β is empty and the production is of the form

$$A \rightarrow \alpha N$$

Calculating Algorithmically

Start with all nonterminal symbols having an empty follow set, $\{\}$, except for the start symbol, S , which has the follow set $\{\$ \}$. We make a pass through the grammar examining the right side of every production. For each occurrence of a nonterminal within the right side of some production, we augment the Follow set for that nonterminal according to the following process. Assume we are processing an occurrence of N and the production is of the form

$$A \rightarrow \alpha N \beta$$

we add $First(\beta) - \{\epsilon\}$ to the Follow set computer for N so far, and if β is nullable, we also add the current Follow set for A to the Follow set for N . After making a complete pass, we repeat the process with the Follow sets computed so far until no Follow sets are modified, in which case we are done.

Example

$$\begin{array}{ll} S \rightarrow xAB & Fst(S) = \{x\} \\ A \rightarrow y|zB & Fst(A) = \{y, z\} \\ B \rightarrow \epsilon|Ax & Fst(B) = \{\epsilon, y, z\} \end{array}$$

S	$\{\$$	$\{\$$	$\{\$$	$\{\$$
A	$\{\}$	$\{y, z, \$, x\}$	$\{y, z, \$, x\}$	$\{y, z, \$, x\}$
B	$\{\}$	$\{\$, y, z\}$	$\{\$, y, z, x\}$	$\{\$, y, z, x\}$

LL(1) Grammar

A BNF grammar is LL(1) if for each nonterminal, N , where $N \rightarrow \alpha_1|\alpha_2|\dots|\alpha_n$,

- the First sets for each pair of alternatives for N are disjoint, and
- if N is nullable, $First(N)$ and $Follow(N)$ are disjoint.

Bottom Up Parsing

Shift/Reduce Parsing

Makes use of a parse stack and has three actions:

- shift** - push the next input symbol onto the stack,
- reduce** - if a sequence of symbols on the top of the stack, α , matches the right side of some production $N \rightarrow \alpha$, then the sequence α on top of the stack is replaced by N .
- accept** - if the stack contains just the start symbol and there is no input left, the input has been recognised and is accepted.

LR(0) Grammars

An LR(0) parsing item is of the form

$$N \rightarrow \alpha \bullet \beta$$

which indicates that, in matching N , α has been matched and β is yet to be matched where

- N is a nonterminal symbol,
- α and β are possibly empty sequences of (terminal and nonterminal) symbols such that $N \rightarrow \alpha \beta$ is a production of the grammar, and
- \bullet marks the current position in matching the right side.

Parsing Automaton

An LR(0) parsing automaton consists of

- a finite set of states, each of which consists of a set of LR(0) parsing items, and
- transitions between states, each of which is labelled by a transition symbol.

Each state in an LR(0) parsing automaton must have only one associated parsing action.

Automaton States

If a state has an LR(0) item of the form

$$N \rightarrow \alpha \bullet M \beta$$

where the nonterminal M has productions

$$\begin{array}{l} M \rightarrow \alpha_1 \\ M \rightarrow \alpha_2 \\ \vdots \\ M \rightarrow \alpha_m \end{array}$$

then the state also includes the derived items

$$\begin{array}{l} M \rightarrow \bullet \alpha_1 \\ M \rightarrow \bullet \alpha_2 \\ \vdots \\ M \rightarrow \bullet \alpha_m \end{array}$$

Goto States

If a state s_i has an LR(0) item of the form

$$N \rightarrow \alpha \bullet x \beta$$

where x is either a terminal symbol or a nonterminal symbol, then there is a goto state, s_j , from the state s_1 on x , and s_j includes a kernel item of the form

$$N \rightarrow \alpha x \bullet \beta$$

If there are multiple items in s_i with the same x immediately to the right of the \bullet then the goto state s_j includes all those items but with the \bullet after that occurrence of x rather than before it.

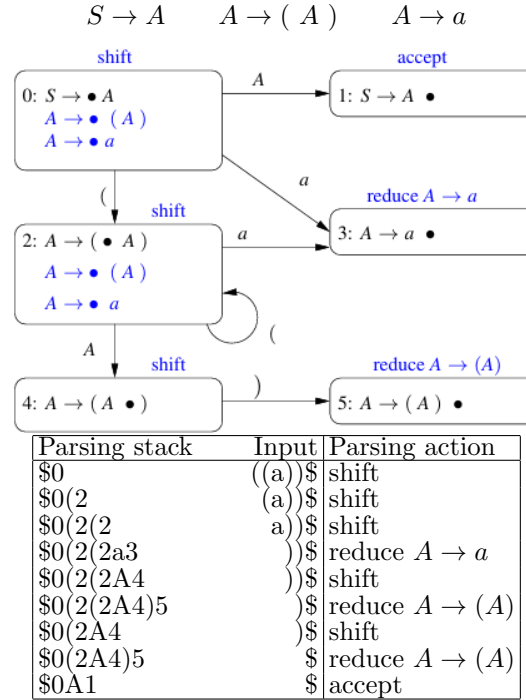
Actions

An LR(0) item of the form

- $N \rightarrow \alpha \bullet a \beta$ where a is a terminal symbol, indicates the state containing the item has a shift **parsing** action
- $S' \rightarrow S \bullet$ where S' is the (introduced) start symbol for the grammar, indicates the state containing the item has an **accept** action
- $N \rightarrow \alpha \bullet$ where N is not the (introduced) start symbol for the grammar, indicates the state containing the item has a parsing action **reduce** $N \rightarrow \alpha$

A shift action at end-of-file is an error, as is an accept action when the input is not at end-of-file.

Example



Parsing Conflicts

If a state in an LR(0) parsing automaton has more than one action, there is a parsing action conflict.

A grammar is LR(0) if none of the states in its LR(0) parsing automaton contains a parsing action conflict.

LR(1) Grammars

An LR(1) parsing item is a pair

$$[N \rightarrow \alpha \bullet \beta, T]$$

consisting of

- an LR(0) parsing item $N \rightarrow \alpha \bullet \beta$, and
- a set T of terminal symbols called a look-ahead set.

The above item indicates that in matching N , α has been matched and β is yet to be matched, in a context in which N can be followed by a terminal symbol in the set T .

Parsing Automaton

The kernel item of the initial state is

$$[S' \rightarrow \bullet S, \$]$$

where S is the start symbol of the grammar and we introduce a fresh replacement start

symbol S' and production $S' \rightarrow S$. This fresh production is used to determine when parsing has completed.

Derived Items

If a state has an LR(1) item of the form

$$[N \rightarrow \alpha \bullet M \beta, T]$$

where the nonterminal M has productions

$$M \rightarrow \alpha_1|\alpha_2|\dots|\alpha_m$$

and $T = \{a_1, a_2, \dots, a_n\}$, then the state also includes the derived items

$$[N \rightarrow \bullet \alpha_1, T] \dots [N \rightarrow \bullet \alpha_m, T]$$

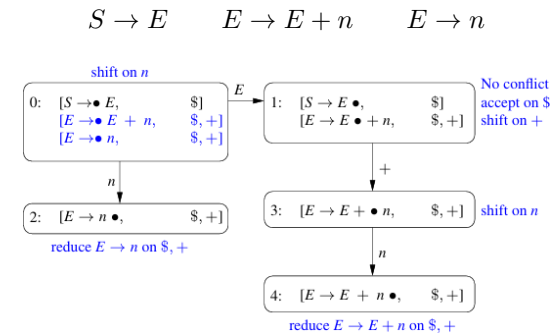
where if β is not nullable $T' = First(\beta)$ and if β is nullable, $T' = First(\beta) - \{\epsilon\} \cup T$.

Parsing Actions

An LR(1) item of the form

- $[N \rightarrow \alpha \bullet a \beta, T]$, where a is a terminal symbol, indicates the state containing the item has a **shift** parsing action if the next input x is a
- $[S' \rightarrow \bullet \$, \$]$, where S' is the added start symbol for the grammar, indicates the state containing the item has an **accept** action if there is no more input.
- $[N \rightarrow \alpha \bullet, T]$, where N is not S' , indicates the state containing the item has a parsing action **reduce** $N \rightarrow \alpha$ if the next input x is in T .

Example



Parsing Action Conflicts

If a state in an LR(1) parsing automaton has more than one action for a look-ahead terminal symbol, there is a parsing action conflict.

A grammar is LR(1) if none of the states in its LR(1) parsing automaton contains a parsing action conflict.

LALR(1) Parsing

- Look-Ahead merged LR(1) parsing where each state of an LALR(1) parsing automaton consists of a set of LR(1) items.
- An LALR(1) parsing automaton can be formed from an LR(1) parsing automaton by merging states that have identical sets of LR(0) items but possibly different look-ahead sets.
- Each item in the merged state consists of one of the LR(0) items in the states being merged with a look-ahead set that is the union of the look-ahead sets of that item in all the states being merged.
- The parsing actions for LALR(1) parsing are defined in the same way as for LR(1) parsing.
- A grammar is LALR(1) if none of the states in its LALR(1) parsing automaton contains a parsing action conflict.

Left Factoring & Left Recursion

Left Factoring Productions

Not all EBNF grammars are suitable for Recursive-Descent Parsing, however, sometimes we can rewrite them into a form that is suitable.

Left Factor Rewriting Rule

To remove the left factor from

$$A \rightarrow \alpha \beta \mid \alpha \gamma,$$

we can rewrite the production using an additional nonterminal A' as

$$A \rightarrow \alpha A'$$

$$A' \rightarrow \beta \mid \gamma$$

Left Recursive Productions

A production of the form

$$E \rightarrow E + T \mid T$$

is not suitable for RDP because the left recursion in the grammar leads to an infinite recursion.

Immediate Left Recursion Rewriting Rule (Simple)

To remove the left recursion from

$$A \rightarrow A \alpha \mid \beta,$$

we can rewrite the production as

$$A \rightarrow \beta A'$$

$$A' \rightarrow \epsilon \mid \alpha A'$$

Immediate Left Recursion Rewriting Rule (General)

To remove the left recursion from the general case

$$A \rightarrow A\alpha_1|A\alpha_2|\dots|A\alpha_n|\beta_1|\beta_2|\dots|\beta_m,$$

we can use grouping to see the structure in the same form as the simple case

$$A \rightarrow A(\alpha_1|\alpha_2|\dots|\alpha_n)|(\beta_1|\beta_2|\dots|\beta_m)$$

which allows us to rewrite the production as follows,

$$A \rightarrow (\beta_1|\beta_2|\dots|\beta_m)A'$$

$$A' \rightarrow \epsilon|(\alpha_1|\alpha_2|\dots|\alpha_n)A'$$

Indirect Left Recursion Rewriting Rule

The following productions have indirect recursion from $A \rightarrow B \rightarrow C \rightarrow A$.

$$A \rightarrow B \alpha$$

$$B \rightarrow C \beta$$

$$C \rightarrow A \gamma_1 \mid \gamma_2$$

To remove the recursion we can first collapse B into A as follows.

$$A \rightarrow C \beta \alpha$$

$$C \rightarrow A \gamma_1 \mid \gamma_2$$

Then we can collapse C into A .

$$A \rightarrow (A \gamma_1 \mid \gamma_2) \beta \alpha$$

↓

$$A \rightarrow A \gamma_1 \beta \alpha \mid \gamma_2 \beta \alpha$$

This now leaves a simple direct left recursion which we can remove as follows.

$$A \rightarrow \gamma_2 \beta \alpha A'$$

$$A' \rightarrow \gamma_1 \beta \alpha A'$$