

Context-Free Languages

- A **context-free grammar** is a phrase-structure grammar $G = (\mathcal{N}, \mathcal{T}, \mathcal{R}, s)$ in which each rule has only a single non-terminal on the left. **Context-free languages** are generated by context-free grammars.

Example: Let $\mathcal{N}_2 = \{s\}$, $\mathcal{T}_2 = \{\epsilon, a, b\}$, $\mathcal{R}_2 = \{s \rightarrow asb, s \rightarrow \epsilon\}$. Then, $G_2 = (\mathcal{N}_2, \mathcal{T}_2, \mathcal{R}_2, s)$ is context-free.

Chomsky Normal Form

- A CFG $G = (\mathcal{N}, \mathcal{T}, \mathcal{R}, s)$ is in **Chomsky normal form** if every rule is of the form $A \rightarrow BC$ or $A \rightarrow b$, $b \in \mathcal{T}$ except if $\epsilon \in L(G)$ in which case $S \rightarrow \epsilon$ is also a rule.

Theorem: Every context-free language L can be generated by a CFG in Chomsky normal form.

Example: $G_3 = (\mathcal{N}_3, \mathcal{T}_3, \mathcal{R}_3, S)$

(a)	S	→	cMnc
(b)	M	→	aMa
(c)	M	→	c

(d)	N	→	bNb
(e)	N	→	c

A Chomsky normal form grammar generating same language uses (c) & (e) and replaces others by:

- (a) $S \rightarrow CD, C \rightarrow c, D \rightarrow ME, E \rightarrow NC,$
- (b) $M \rightarrow AF, A \rightarrow a, F \rightarrow MA$
- (d) $N \rightarrow BG, B \rightarrow b, G \rightarrow NB$

Putting CFG in Chomsky Normal Form

Theorem: Every context-free language L can be generated by a CFG in Chomsky normal form.

Proof: If $\epsilon \in L$, add $s \rightarrow \epsilon$. Let L be generated by G . Convert G to G' in Chomsky normal form in stages.

a) eliminate from G ϵ -rules of the form $B \rightarrow \epsilon$ (except for $S \rightarrow \epsilon$) as follows: for each rule with ≥ 1 B in right-hand side, e.g. $A \rightarrow \alpha B \beta \gamma$ (α, β, γ are strings), add all possible rules formed by replacing B by ϵ in all possible ways e.g. $A \rightarrow \alpha \beta \gamma, A \rightarrow \alpha B \beta \gamma, A \rightarrow \alpha \beta \gamma$, giving four rules for one original rule.

b) For rules $A \rightarrow \alpha w_i \beta$ (α, β are strings) with $w_i \in \mathcal{T}$ replace it by $A \rightarrow \alpha Z_i \beta$ & add rule $Z_i \rightarrow w_i$, where Z_i is a new non-terminal. Continue until all rules have a single terminal on right or a string of non-terminals. This new grammar also generates L .

Chomsky Normal Form (cont.)

Proof (cont.) Rules are now of the form: a) $A \rightarrow b$ for $b \in \mathcal{T}$; b) $S \rightarrow \epsilon$; c) $A \rightarrow Z_1 Z_2 \dots Z_k$, for $Z_i \in \mathcal{N}$

Consider rules of type c) with $k = 1$. Cascading such rules gives derivations $A \xRightarrow{*} B$; delete all rules of type c) with $k = 1$ and replace them with $A \rightarrow B$ if $A \xRightarrow{*} B$. The same language is generated.

If $C \rightarrow D$ and $D \rightarrow b$, add $C \rightarrow b$, deleting all rules of the form $A \rightarrow B$. This generates same language; all remaining rules are of the form $S \rightarrow \epsilon, A \rightarrow b$ or $A \rightarrow Z_1 Z_2 \dots Z_k$ with $k \geq 2, Z_i \in \mathcal{N}$

Now replace all rules of the form $A \rightarrow Z_1 Z_2 \dots Z_k$ by the rules $A \rightarrow Z_1 N_1, N_1 \rightarrow Z_2 N_2, \dots, N_{k-3} \rightarrow Z_{k-2} N_{k-2}, N_{k-2} \rightarrow Z_{k-1} Z_k$ where each N_i is a new non-terminal.

This new grammar is in correct form and generates L .

Example

Let $G = (\mathcal{N}, \mathcal{T}, \mathcal{R}, E)$ be the grammar with $\mathcal{N} = \{E, T, F\}$, $\mathcal{T} = \{a, b, +, *, (,)\}$ and let \mathcal{R} have following rules:

(a)	E	→	E+T	(e)	F	→	(E)
(b)	E	→	T	(f)	F	→	a
(c)	T	→	T*F	(g)	F	→	b
(d)	T	→	F				

E, T, and F denote expressions, terms & factors. Easy to show that $E \xRightarrow{*} (a*b+a)*(a+b)$ and $E \xRightarrow{*} a*b+a$.

Conversion: 1) no ϵ -rules, 2) in $A \rightarrow w_1 w_2 \dots w_k$ replace (by (,) by), + by +, and * by *, 3) cascade non-terminals.

(1)	E	→	E+T	(9)	T	→	b
(2)	E	→	T*F	(10)	F	→	(E)
(3)	E	→	(E)	(11)	F	→	a
(4)	E	→	a	(12)	F	→	b
(5)	E	→	b	(13)	(→	(
(6)	T	→	T*F	(14))	→)
(7)	T	→	(E)	(15)	+	→	+
(8)	T	→	a	(16)	*	→	*

The last step is to reduce to two the number of non-terminals on the right-hand side of a rule.

Example

$G = (\mathcal{N}, \mathcal{T}, \mathcal{R}, E)$, where $\mathcal{N} = \{A, B, C, D, E, F, G, T, (, +, *,)\}$, $\mathcal{T} = \{a, b, +, *, (,)\}$ and let \mathcal{R} have following rules:

(A)	E	→	EA	(L)	G	→	(E)
(B)	A	→	+T	(M)	T	→	a
(C)	E	→	TB	(N)	T	→	b
(D)	B	→	*F	(P)	F	→	(H
(E)	E	→	(C	(Q)	H	→	(E)
(F)	C	→	E)	(R)	F	→	a
(G)	E	→	a	(S)	F	→	b
(H)	E	→	b	(T)	(→	(
(I)	T	→	TD	(U))	→)
(J)	D	→	*F	(V)	+	→	+
(K)	T	→	(G	(W)	*	→	*

This grammar is in Chomsky normal form.

We illustrate bottom-up parsing.

Parsing Using Chomsky Normal Form

Let's parse $w = w_1 w_2 w_3 w_4 w_5 = a*b+a$ where $w_1=a, w_2=*, w_3=b, w_4=+, w_5=a$.

- w_i is derived by rules $A \rightarrow b$ where A is in $N_{i,i+1} = \{N \mid N \rightarrow w_i\}$.
- Thus, $N_{1,2} = \{E, T, F\}, N_{2,3} = \{*\}, N_{3,4} = \{E, T, F\}, N_{4,5} = \{+\}, N_{5,6} = \{E, T, F\}$.
- All other derivations are of form $A \rightarrow BC$.
- To derive $w_i w_{i+1}$, we use a nonterminal in $N_{i,i+2} = \{N \mid N \rightarrow MP, M \text{ in } N_{i,i+1}, P \text{ in } N_{i+1,i+2}\}$
- Not knowing which pair $w_i w_{i+1}$ is derived last we try all ways & compute $N_{i,i+2}$ for $i = 1, 2, 3$.
- We try all ways to derive $w_i w_{i+1} w_{i+2}$ computing $N_{i,i+3} = \{N \mid N \rightarrow MP, M \text{ in } N_{i,i+2}, P \text{ in } N_{i+2,i+3}\} \cup \{N \mid N \rightarrow MP, M \text{ in } N_{i,i+1}, P \text{ in } N_{i+1,i+3}\}$
- Continuing, we compute $N_{1,4}$. If it contains the start symbol, w is in our language and we have a parse for it.

More on Parsing

(A)	E	→	EA	(L)	G	→	(E)
(B)	A	→	+T	(M)	T	→	a
(C)	E	→	TB	(N)	T	→	b
(D)	B	→	*F	(P)	F	→	(H
(E)	E	→	(C	(Q)	H	→	(E)
(F)	C	→	E)	(R)	F	→	a
(G)	E	→	a	(S)	F	→	b
(H)	E	→	b	(T)	(→	(
(I)	T	→	TD	(U))	→)
(J)	D	→	*F	(V)	+	→	+
(K)	T	→	(G	(W)	*	→	*

- $N_{1,2} = \{E, T, F\}, N_{2,3} = \{*\}, N_{3,4} = \{E, T, F\}, N_{4,5} = \{+\}, N_{5,6} = \{E, T, F\}$
- $N_{1,3} = \emptyset, N_{2,4} = \{B, D\}, N_{3,5} = \emptyset, N_{4,6} = \{A\}$
- $N_{1,4} = \{E, T\}, N_{2,5} = \emptyset, N_{3,6} = \{E\}$
- $N_{1,5} = \emptyset, N_{2,6} = \emptyset$
- $N_{1,6} = \{E\}$

Summary

- Chomsky normal form
- Multiplication and closure of set matrices – basis for parsing of context-free languages.