

CS 302.1 - Automata Theory

Lecture 05

Shantanav Chakraborty

Center for Quantum Science and Technology (CQST)

Center for Security, Theory and Algorithms (CSTAR)

IIIT Hyderabad



Quick Recap

Grammars: A Grammar is a 4-tuple (V, Σ, P, S) , such that

- V is the set of **Variables**
 - Σ is the set of **Terminals**
 - P is the set of production **Rules**
 - S is the **Start Variable**
- [$(V \cup T)^* V (V \cup T)^* \rightarrow (V \cup T)^*$]
[The variable in the LHS of the first rule is generally the start variable]
- To show that a string $w \in L(G)$, we show that there exists a **derivation ending up in w** ($S \xrightarrow{*} w$).
 - The **language of the grammar**, $L(G)$ is $\{w \in \Sigma^* | S \xrightarrow{*} w\}$

Right Linear grammar: If the *rules* of the underlying grammar G are of the form

$$\begin{aligned}Var &\rightarrow Ter \; Var \\Var &\rightarrow Ter \\Var &\rightarrow \epsilon\end{aligned}$$

then it is **Right-linear grammar**.

Left linear grammar: If the *rules* of the underlying grammar G are of the form

$$\begin{aligned}Var &\rightarrow Var \; Ter \\Var &\rightarrow Ter \\Var &\rightarrow \epsilon\end{aligned}$$

then such a grammar is called **Left-linear grammar**.

Left-linear grammar \equiv Right-linear grammar \equiv DFA \equiv NFA \equiv Regular Expressions

Quick Recap

Grammars: A Grammar is a 4-tuple (V, Σ, P, S) , such that

- V is the set of **Variables**
- Σ is the set of **Terminals**
- P is the set of production **Rules** $[(V \cup T)^* V (V \cup T)^* \rightarrow (V \cup T)^*]$
[The variable in the LHS of the first rule is generally the start variable]
- S is the **Start Variable**
- To show that a string $w \in L(G)$, we show that there exists a **derivation ending up in w** ($S \xrightarrow{*} w$).
- The **language of the grammar**, $L(G)$ is $\{w \in \Sigma^* | S \xrightarrow{*} w\}$

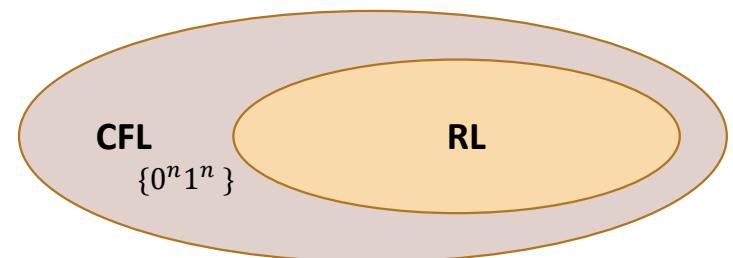
Context-Free Grammars: If the *rules* of the underlying grammar G are of the form

$$V \rightarrow (V \cup T)^*$$

then such a grammar is called **Context-Free**.

$$L(G) = \{\omega | \omega = 0^n 1^n, n \geq 0\}$$

So although $L(G)$ is not regular, it is context-free.



Context free Grammars

Context-Free Grammars: If the *rules* of the underlying grammar G are of the form

$$V \rightarrow (VUT)^*$$

then such a grammar is called **Context-Free**.

Regular languages \subset **Context Free Languages**.

Consider the Grammar G with the following rules:

$$S \rightarrow 0S1|SS|\epsilon$$

Strings that can be derived by G :

$$S \rightarrow \epsilon$$

$$\{\epsilon\}$$

Context free Grammars

Context-Free Grammars: If the *rules* of the underlying grammar G are of the form

$$V \rightarrow (VUT)^*$$

then such a grammar is called **Context-Free**.

Regular languages \subset **Context Free Languages**.

Consider the Grammar G with the following rules:

$$S \rightarrow 0S1|SS|\epsilon$$

Strings that can be derived by G :

$$S \rightarrow 0S1 \rightarrow 00S11 \dots$$

$$\{\epsilon, 01, 0011, \dots 0^n1^n\}$$

Context free Grammars

Context-Free Grammars: If the *rules* of the underlying grammar G are of the form

$$V \rightarrow (VUT)^*$$

then such a grammar is called **Context-Free**.

Regular languages \subset **Context Free Languages**.

Consider the Grammar G with the following rules:

$$S \rightarrow 0S1|SS|\epsilon$$

Strings that can be derived by G :

$$S \rightarrow \mathbf{0S1} \rightarrow 0\mathbf{SS1} \rightarrow 00\mathbf{S1S1} \rightarrow 001S1 \rightarrow 001\mathbf{0S11} \rightarrow 001011$$

$$\{\epsilon, 01, 0011, \dots 0^n1^n, 001011, \dots\}$$

Context free Grammars

Context-Free Grammars: If the *rules* of the underlying grammar G are of the form

$$V \rightarrow (VUT)^*$$

then such a grammar is called **Context-Free**.

Regular languages \subset **Context Free Languages**.

Consider the Grammar G with the following rules:

$$S \rightarrow 0S1|SS|\epsilon$$

Strings that can be derived by G :

$$S \rightarrow \mathbf{0S1} \rightarrow 0\mathbf{SS1} \rightarrow 00\mathbf{S1S1} \rightarrow 001S1 \rightarrow 001\mathbf{0S11} \rightarrow 001011$$

$$\{\epsilon, 01, 0011, \dots 0^n1^n, 001011, \dots\}$$

Show that the string $010101 \in L(G)$.

$$S \rightarrow SS \rightarrow SSS \rightarrow 0S1\ 0S1\ 0S1 \rightarrow 010101$$

Context free Grammars

Context-Free Grammars: If the *rules* of the underlying grammar G are of the form

$$V \rightarrow (VUT)^*$$

then such a grammar is called **Context-Free**.

Regular languages \subset **Context Free Languages**.

Consider the Grammar G with the following rules:

$$S \rightarrow 0S1|SS|\epsilon$$

Strings that can be derived by G :

$$S \rightarrow SS \rightarrow SSS \rightarrow 0S1SS \rightarrow 0S10S1S \rightarrow 0S10S10S1 \rightarrow 010101$$

$$\{\epsilon, 01, 0011, \dots 0^n1^n, 001011, 010101, \dots\}$$

Context free Grammars

Context-Free Grammars: If the *rules* of the underlying grammar G are of the form

$$V \rightarrow (VUT)^*$$

then such a grammar is called **Context-Free**.

Regular languages \subset **Context Free Languages**.

Consider the Grammar G with the following rules:

$$S \rightarrow 0S1|SS|\epsilon$$

Strings that can be derived by G :

$$S \rightarrow SS \rightarrow SSS \rightarrow 0S1SS \rightarrow 0S10S1S \rightarrow 0S10S10S1 \rightarrow 010101$$

$$\{\epsilon, 01, 0011, \dots 0^n1^n, 001011, 010101, \dots\}$$

What is $L(G)$?

Context free Grammars

Consider the Grammar G with the following rules:

$$S \rightarrow 0S1|SS|\epsilon$$

Strings that can be derived by G :

$$\{\epsilon, 01, 0011, \dots 0^n1^n, 001011, 010101, \dots\}$$

What is $L(G)$?

Context free Grammars

Consider the Grammar G with the following rules:

$$S \rightarrow 0S1|SS|\epsilon$$

Strings that can be derived by G :

$$\{\epsilon, 01, 0011, \dots 0^n1^n, 001011, 010101, \dots\}$$

What is $L(G)$?

You can see what the language is, if you replace **0** with **(** and **1** with **)**

Context free Grammars

Consider the Grammar G with the following rules:

$$S \rightarrow 0S1|SS|\epsilon$$

Strings that can be derived by G :

$$\{\epsilon, 01, 0011, \dots 0^n1^n, 001011, 010101, \dots\}$$

What is $L(G)$?

You can see what the language is, if you replace **0** with **(** and **1** with **)**

Strings that can be derived by G : $\{\epsilon, 01, 0011, \dots, 0^n1^n, 001011, 010101, \dots\}$

$$\{\epsilon, ()\}$$

Context free Grammars

Consider the Grammar G with the following rules:

$$S \rightarrow 0S1|SS|\epsilon$$

Strings that can be derived by G :

$$\{\epsilon, 01, 0011, \dots 0^n1^n, 001011, 010101, \dots\}$$

What is $L(G)$?

You can see what the language is, if you replace **0** with **(** and **1** with **)**

Strings that can be derived by G : $\{\epsilon, 01, 0011, \dots, 0^n1^n, 001011, 010101, \dots\}$

$$\{\epsilon, (), (()), \dots, \}$$

Context free Grammars

Consider the Grammar G with the following rules:

$$S \rightarrow 0S1|SS|\epsilon$$

Strings that can be derived by G :

$$\{\epsilon, 01, 0011, \dots, 0^n1^n, 001011, 010101, \dots\}$$

What is $L(G)$?

You can see what the language is, if you replace **0** with **(** and **1** with **)**

Strings that can be derived by G : $\{\epsilon, 01, 0011, \dots, 0^n1^n, 001011, 010101, \dots\}$

$$\{\epsilon, (), (()), \dots, (((\dots)))\}$$

Context free Grammars

Consider the Grammar G with the following rules:

$$S \rightarrow 0S1|SS|\epsilon$$

Strings that can be derived by G :

$$\{\epsilon, 01, 0011, \dots, 0^n1^n, 001011, 010101, \dots\}$$

What is $L(G)$?

You can see what the language is, if you replace **0** with **(** and **1** with **)**

Strings that can be derived by G : $\{\epsilon, 01, 0011, \dots, 0^n1^n, 001011, 010101, \dots\}$

$$\{\epsilon, (), (()), \dots, (((\dots)))\}, (())(()), \dots\}$$

Context free Grammars

Consider the Grammar G with the following rules:

$$S \rightarrow 0S1|SS|\epsilon$$

Strings that can be derived by G :

$$\{\epsilon, 01, 0011, \dots 0^n1^n, 001011, 010101, \dots\}$$

What is $L(G)$?

You can see what the language is, if you replace **0** with **(** and **1** with **)**

Strings that can be derived by G : $\{\epsilon, 01, 0011, \dots, 0^n1^n, 001011, 010101, \dots\}$

$$\{\epsilon, (), (()), \dots, (((\dots))) \}, (()()), 000, \dots\}$$

So, $L(G)$ is the language of all strings of properly nested parentheses.

$$L(G) = \{\omega \mid \omega \text{ is a correctly nested parenthesis}\}$$

Context free Grammars

Constructing CFG corresponding to a Language.

There is no fixed recipe for doing this. Requires some level of creativity.

Some tips might come in handy:

- Check if the CFL is a union of simpler languages. If $L(G) = L(G_1) \cup L(G_2)$ and G_1 and G_2 are known. If S_1 is the start variable for G_1 and S_2 is the start variable for G_2 then the rules of G :

$$\begin{aligned} S &\rightarrow S_1 | S_2 \\ S_1 &\rightarrow \dots \dots \\ S_2 &\rightarrow \dots \dots \end{aligned}$$

Context free Grammars

Constructing CFG corresponding to a Language.

There is no fixed recipe for doing this. Requires some level of creativity.

Some tips might come in handy:

- Check if the CFL is a union of simpler languages. If $L(G) = L(G_1) \cup L(G_2)$ and G_1 and G_2 are known. If S_1 is the start variable for G_1 and S_2 is the start variable for G_2 then the rules of G :

$$\begin{aligned} S &\rightarrow S_1 | S_2 \\ S_1 &\rightarrow \dots \dots \\ S_2 &\rightarrow \dots \dots \end{aligned}$$

- Grammars with rules such as $S \rightarrow aSb$ help generate strings where the corresponding machine would need unbounded memory to *remember* the number of a 's needed to verify that there are an equal number of b 's. This was not possible with regular expressions/linear grammars.

Context free Grammars

Constructing CFG corresponding to a Language.

- Check if the CFL is a union of simpler languages.
- Grammars with rules such as $S \rightarrow aSb$ help generate where the portions of a and b are equal.

Example: Construct the grammar G such that $L(G) = \{\omega \mid \omega \text{ has equal number of 0's and 1's}\}$

Context free Grammars

Constructing CFG corresponding to a Language.

- Check if the CFL is a union of simpler languages.
- Grammars with rules such as $S \rightarrow aSb$ help generate where the portions of a and b are equal.

Example: Construct the grammar G such that $L(G) = \{\omega \mid \omega \text{ has equal number of 0's and 1's}\}$

- The first thing to notice is that $L_1 = \{0^n 1^n, n \geq 0\} \subset L(G)$. We know the grammar for this language.
- Any string $\omega \in L_1$ has a series of 0's followed by an equal number of 1's.
- Again, consider L_2 to comprise all strings that start with a series of 1's followed by an equal number of 0's, i.e.

$$L_2 = \{1^n 0^n, n \geq 0\}$$

- The grammar for L_2 is similar to that of L_1 : replace the 0's with 1's and vice versa. Importantly, $L_2 = \{1^n 0^n, n \geq 0\} \subset L(G)$ also.
- Also, $L_1 \cup L_2 \subset L(G)$

Context free Grammars

Constructing CFG corresponding to a Language.

- Check if the CFL is a union of simpler languages.
- Grammars with rules such as $S \rightarrow aSb$ help generate where the portions of a and b are equal.

Example: Construct the grammar G such that $L(G) = \{\omega \mid \omega \text{ has equal number of 0's and 1's}\}$

- So $L'(G') = \{0^n 1^n \mid n \geq 0\} \cup \{1^n 0^n \mid n \geq 0\} \subset L(G)$
- Grammar for L_1 : $S \rightarrow 0S1|\epsilon$
- Grammar for L_2 : $S \rightarrow 1S0|\epsilon$
- Grammar for $L_1 \cup L_2$:

$$\begin{aligned} S &\rightarrow S_1 | S_2 \\ S_1 &\rightarrow 0S_11|\epsilon \\ S_2 &\rightarrow 1S_20|\epsilon \end{aligned}$$

Context free Grammars

Constructing CFG corresponding to a Language.

- Check if the CFL is a union of simpler languages.
- Grammars with rules such as $S \rightarrow aSb$ help generate where the portions of a and b are equal.

Example: Construct the grammar G such that $L(G) = \{\omega \mid \omega \text{ has equal number of 0's and 1's}\}$

- Grammar for $L_1 \cup L_2$:

$$\begin{aligned} S &\rightarrow S_1 | S_2 \\ S_1 &\rightarrow 0S_11|\epsilon \\ S_2 &\rightarrow 1S_20|\epsilon \end{aligned}$$

- Is that all? Is $L_1 \cup L_2 = L(G)$? $L_1 \cup L_2$ contains all strings that have equal number 0's followed by equal number of 1's or vice versa.

Context free Grammars

Constructing CFG corresponding to a Language.

- Check if the CFL is a union of simpler languages.
- Grammars with rules such as $S \rightarrow aSb$ help generate where the portions of a and b are equal.

Example: Construct the grammar G such that $L(G) = \{\omega \mid \omega \text{ has equal number of 0's and 1's}\}$

- Grammar for $L_1 \cup L_2$:

$$\begin{aligned} S &\rightarrow S_1 | S_2 \\ S_1 &\rightarrow 0S_11 | \epsilon \\ S_2 &\rightarrow 1S_20 | \epsilon \end{aligned}$$

- Is that all? Is $L_1 \cup L_2 = L(G)$? $L_1 \cup L_2$ contains all strings that have equal number 0's followed by equal number of 1's or vice versa.
- What about strings such as $s_1 = 0101 \dots$ and $s_2 = 1010 \dots$? For this we need to be able to go from

$$0S_11 \rightarrow 0S_21 \rightarrow 01S_201 \rightarrow \dots$$

Context free Grammars

Constructing CFG corresponding to a Language.

Example: Construct the grammar G such that $L(G) = \{\omega \mid \omega \text{ has equal number of 0's and 1's}\}$

- Grammar for $L_1 \cup L_2$:

$$\begin{aligned} S &\rightarrow S_1 | S_2 \\ S_1 &\rightarrow 0S_11|\epsilon \\ S_2 &\rightarrow 1S_20|\epsilon \end{aligned}$$

- What about strings such as $s_1 = 0101\cdots$ and $s_2 = 1010\cdots$? Add transitions $S_1 \rightarrow S_2$ and $S_2 \rightarrow S_1$.

$$\begin{aligned} S &\rightarrow S_1 | S_2 \\ S_1 &\rightarrow 0S_11|\epsilon \\ S_2 &\rightarrow 1S_20|\epsilon \\ S_1 &\rightarrow S_2 \\ S_2 &\rightarrow S_1. \end{aligned}$$

Context free Grammars

Constructing CFG corresponding to a Language.

Example: Construct the grammar G such that $L(G) = \{\omega \mid \omega \text{ has equal number of 0's and 1's}\}$

$$\begin{aligned} S &\rightarrow S_1 | S_2 \\ S_1 &\rightarrow 0S_11 | \epsilon \\ S_2 &\rightarrow 1S_20 | \epsilon \\ S_1 &\rightarrow S_2 \\ S_2 &\rightarrow S_1 \end{aligned}$$

- Can't we simplify this? We can replace S_1 and S_2 with a single Start variable as follows: $S \rightarrow 0S1|1S0|\epsilon$
- What kind of strings does the grammar generate? Well if we use Rule $S \rightarrow 0S1$, m times, we get to rules such as $0^m S 1^m$.
- Now applying the rule $S \rightarrow 1S0$, k times, we get $0^m 1^k S 0^k 1^m$.
- So the strings we obtain are of the form:

$$\{0^{m_1} 1^{n_1} 0^{m_2} 1^{n_2} \dots 0^{n_2} 1^{m_2} 0^{n_1} 1^{m_1}\} \cup \{1^{m_1} 0^{n_1} 1^{m_2} 0^{n_2} \dots 1^{n_2} 0^{m_2} 1^{n_1} 0^{m_1}\} \in L(G)$$

Context free Grammars

Constructing CFG corresponding to a Language.

Example: Construct the grammar G such that $L(G) = \{\omega \mid \omega \text{ has equal number of 0's and 1's}\}$

$$\begin{aligned} S &\rightarrow S_1 | S_2 \\ S_1 &\rightarrow 0S_11 | \epsilon \\ S_2 &\rightarrow 1S_20 | \epsilon \\ S_1 &\rightarrow S_2 \\ S_2 &\rightarrow S_1 \end{aligned}$$

- Simplified grammar:

$$S \rightarrow 0S1 | 1S0 | \epsilon$$

Context free Grammars

Constructing CFG corresponding to a Language.

Example: Construct the grammar G such that $L(G) = \{\omega \mid \omega \text{ has equal number of 0's and 1's}\}$

$$\begin{aligned} S &\rightarrow S_1 | S_2 \\ S_1 &\rightarrow 0S_11 | \epsilon \\ S_2 &\rightarrow 1S_20 | \epsilon \\ S_1 &\rightarrow S_2 \\ S_2 &\rightarrow S_1 \end{aligned}$$

- Simplified grammar:

$$S \rightarrow 0S1|1S0|\epsilon$$

- Is that all? What about strings such as $\{01\mathbf{10}, 0011\mathbf{1100}\}$?
- More generally, what about strings that are a concatenation of L_1 and L_2 ?

Context free Grammars

Constructing CFG corresponding to a Language.

Example: Construct the grammar G such that $L(G) = \{\omega \mid \omega \text{ has equal number of 0's and 1's}\}$

$$\begin{aligned} S &\rightarrow S_1 | S_2 \\ S_1 &\rightarrow 0S_11 | \epsilon \\ S_2 &\rightarrow 1S_20 | \epsilon \\ S_1 &\rightarrow S_2 \\ S_2 &\rightarrow S_1 \end{aligned}$$

- Simplified grammar:

$$S \rightarrow 0S1|1S0|\epsilon$$

- Is that all? What about strings such as $\{01\textcolor{red}{10}, 0011\textcolor{red}{1100}\}$?
- More generally, what about strings that are a concatenation of L_1 and L_2 ?
- Adding transitions like $S \rightarrow S_1S_2$ incorporates this.

Context free Grammars

Constructing CFG corresponding to a Language.

Example: Construct the grammar G such that $L(G) = \{\omega | \omega \text{ has equal number of 0's and 1's}\}$

$$\begin{aligned} S &\rightarrow S_1 | S_2 | S_1 S_2 \\ S_1 &\rightarrow 0 S_1 1 | \epsilon \\ S_2 &\rightarrow 1 S_2 0 | \epsilon \\ S_1 &\rightarrow S_2 \\ S_2 &\rightarrow S_1 \end{aligned}$$

- Simplify this further.

$$G: S \rightarrow SS | 0S1 | 1S0 | \epsilon$$

Parse trees for CFG

Consider the Grammar G with the following rules:

$$S \rightarrow 0S1|SS|\epsilon$$

Parse trees: These are ordered trees that provide alternative representations of the derivation of a grammar.

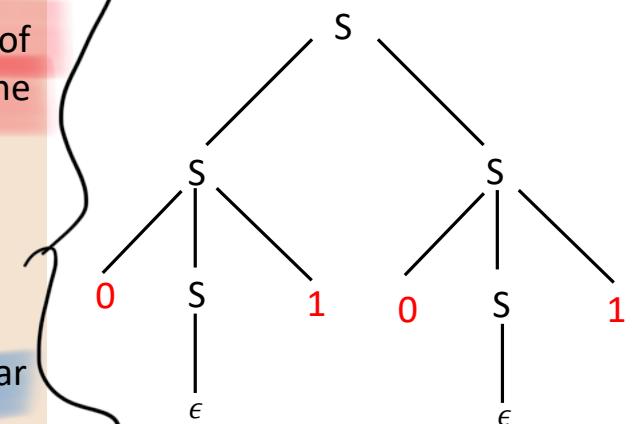
Parsing is a useful technique for compilers (Analysis of syntax eg: take sequence of tokens as input & output parse trees which provides structural representation of the input while checking for the correct syntax).

Features:

- The root node is the **Start variable**
- Branch out to nodes of the next level by following any of the rules of the grammar
- Stop when all the leaf nodes of the tree are terminals
- Read the terminals in the leaves from left to right.
- If w is the string obtained, then $S \xrightarrow{*} w$ and $w \in L(G)$

One derivation:

$$S \xrightarrow{*} SS \rightarrow 0S1S \rightarrow 0S10S1 \rightarrow 0101$$



Parse trees for CFG

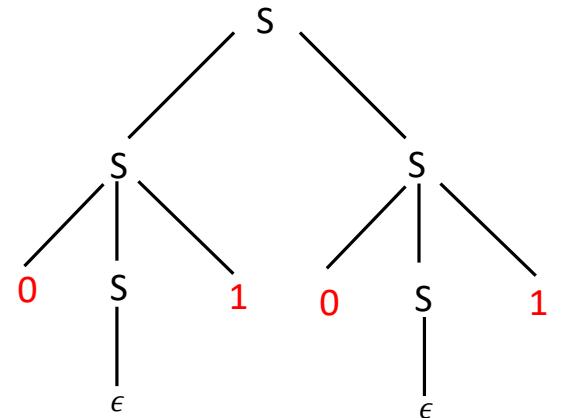
Consider the Grammar G with the following rules:

$$S \rightarrow 0S1|SS|\epsilon$$

Consider the following derivations for 0101:

1. $S \rightarrow SS \rightarrow 0S1S \rightarrow 0S10S1 \rightarrow 0101$
2. $S \rightarrow SS \rightarrow 0S1S \rightarrow 01S \rightarrow 010S1 \rightarrow 0101$
3. $S \rightarrow SS \rightarrow S0S1 \rightarrow S01 \rightarrow 0S101 \rightarrow 0101$

- The parse trees for all these derivations are the same.



Parse trees for CFG

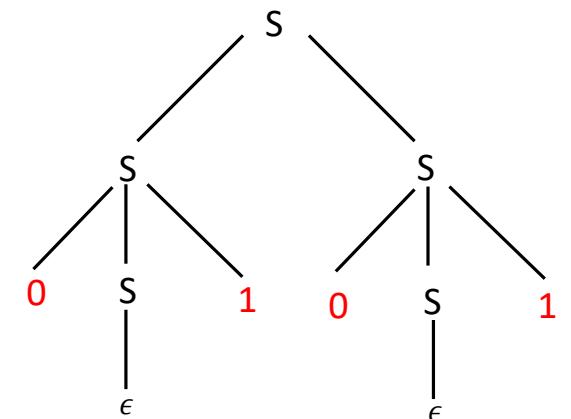
Consider the Grammar G with the following rules:

$$S \rightarrow 0S1|SS|\epsilon$$

Consider the following derivations for 0101:

1. $S \rightarrow SS \rightarrow 0S1S \rightarrow 0S10S1 \rightarrow 0101$
2. $S \rightarrow SS \rightarrow 0S1S \rightarrow 01S \rightarrow 010S1 \rightarrow 0101$
3. $S \rightarrow SS \rightarrow S0S1 \rightarrow S01 \rightarrow 0S101 \rightarrow 0101$

- The parse trees for all these derivations are the same.
- If a string is derived by replacing only the leftmost variable at every step, then the derivation is a **leftmost derivation**. (e.g. derivation 2.)
-rightmost variable = **rightmost derivation** (e.g. derivation 3.)
- Derivations may not always be **leftmost** or **rightmost** (e.g. derivation 1.)



Parse trees for CFG

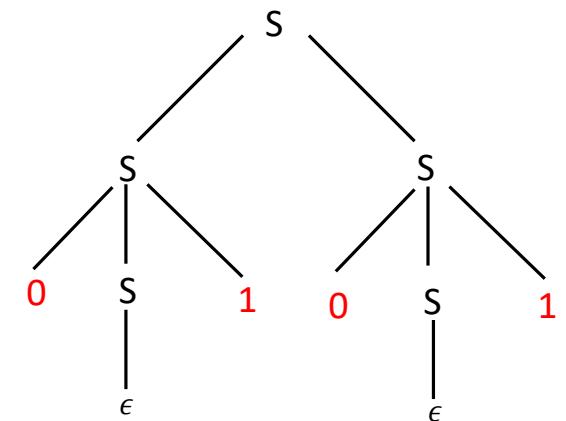
Consider the Grammar G with the following rules:

$$S \rightarrow 0S1|SS|\epsilon$$

Consider the following derivations for 0101:

1. $S \rightarrow SS \rightarrow 0S1S \rightarrow 0S10S1 \rightarrow 0101$
2. $S \rightarrow SS \rightarrow 0S1S \rightarrow 01S \rightarrow 010S1 \rightarrow 0101$
3. $S \rightarrow SS \rightarrow S0S1 \rightarrow S01 \rightarrow 0S101 \rightarrow 0101$

- The parse trees for all these derivations are the same.
- If a string is derived by replacing only the leftmost variable at every step, then the derivation is a **leftmost derivation**. (e.g. derivation 2.)
-rightmost variable = **rightmost derivation** (e.g. derivation 3.)
- Derivations may not always be **leftmost** or **rightmost** (e.g. derivation 1.)



Ambiguous grammars: A CFG G is said to be **ambiguous** if there exists $\omega \in L(G)$, such that there are **two or more leftmost derivations for ω** (or equivalently two or more rightmost derivations) or equivalently **two or more parse trees for ω** .

Parse trees for CFG

Ambiguous grammars: A CFG G is said to be **ambiguous** if there exists $\omega \in L(G)$, such that there are **two or more leftmost derivations for ω** (or equivalently two or more rightmost derivations) or equivalently **two or more parse trees for ω** .

Consider the Grammar G with the following rules: $S \rightarrow 0S1|SS|\epsilon$

Show that Grammar G is ambiguous, i.e. $\exists \omega \in L(G)$, such that there are two or more parse trees for ω .

Consider the string $\omega = 010101$:

- Show that there exist two different parse trees for **010101**.
- Show that there exist two leftmost derivations for **010101**.

Parse trees for CFG

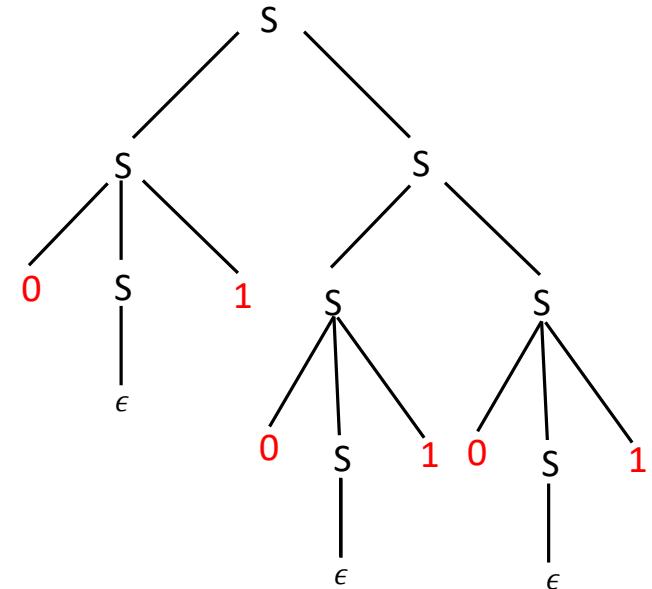
Ambiguous grammars: A CFG G is said to be **ambiguous** if there exists $\omega \in L(G)$, such that there are **two or more leftmost derivations for ω** (or equivalently two or more rightmost derivations) or equivalently **two or more parse trees for ω** .

Consider the Grammar G with the following rules: $S \rightarrow 0S1|SS|\epsilon$

Show that Grammar G is ambiguous, i.e. $\exists \omega \in L(G)$, such that there are two or more parse trees for ω .

Consider the string $\omega = 010101$:

- Show that there exist two different parse trees for 010101 .
- Show that there exist two leftmost derivations for 010101 .



Leftmost Derivation: $S \rightarrow SS$

Parse trees for CFG

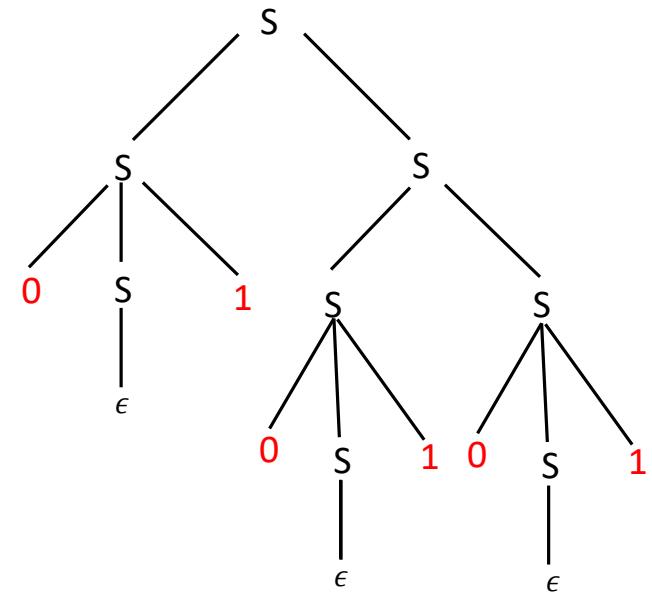
Ambiguous grammars: A CFG G is said to be **ambiguous** if there exists $\omega \in L(G)$, such that there are **two or more leftmost derivations for ω** (or equivalently two or more rightmost derivations) or equivalently **two or more parse trees for ω** .

Consider the Grammar G with the following rules: $S \rightarrow 0S1|SS|\epsilon$

Show that Grammar G is ambiguous, i.e. $\exists \omega \in L(G)$, such that there are two or more parse trees for ω .

Consider the string $\omega = 010101$:

- Show that there exist two different parse trees for 010101 .
- Show that there exist two leftmost derivations for 010101 .



Leftmost Derivation: $S \rightarrow SS$

Parse trees for CFG

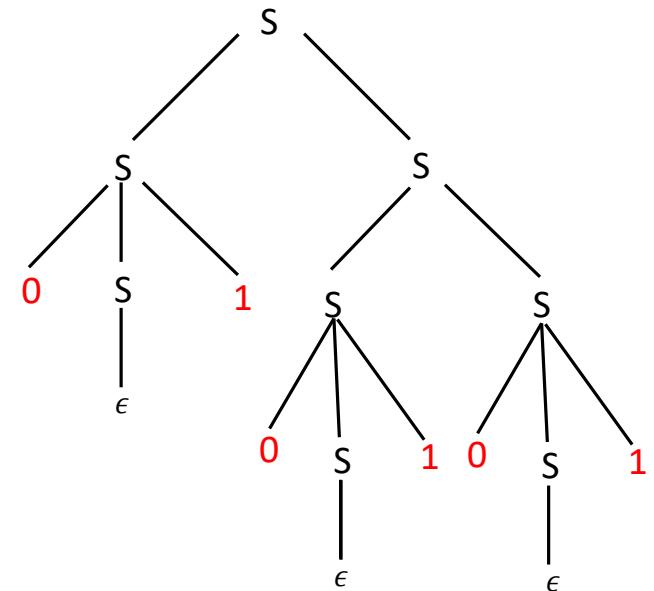
Ambiguous grammars: A CFG G is said to be **ambiguous** if there exists $\omega \in L(G)$, such that there are **two or more leftmost derivations for ω** (or equivalently two or more rightmost derivations) or equivalently **two or more parse trees for ω** .

Consider the Grammar G with the following rules: $S \rightarrow 0S1|SS|\epsilon$

Show that Grammar G is ambiguous, i.e. $\exists \omega \in L(G)$, such that there are two or more parse trees for ω .

Consider the string $\omega = 010101$:

- Show that there exist two different parse trees for 010101 .
- Show that there exist two leftmost derivations for 010101 .



Leftmost Derivation: $S \rightarrow SS \rightarrow 0S1S$

Parse trees for CFG

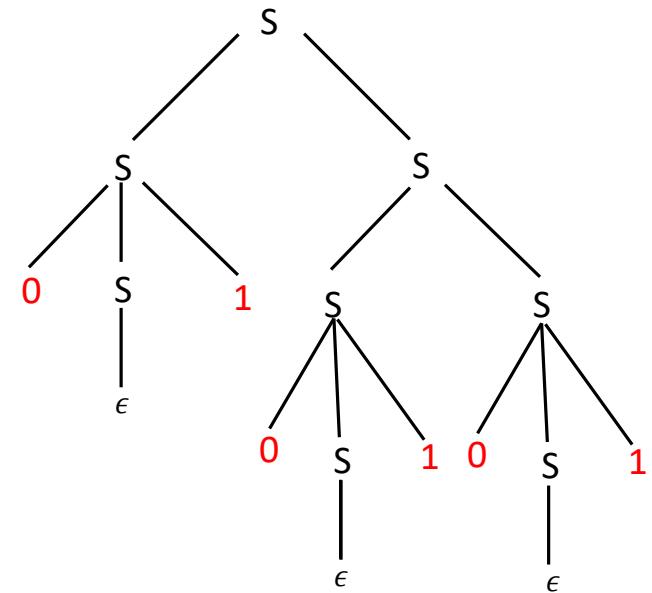
Ambiguous grammars: A CFG G is said to be **ambiguous** if there exists $\omega \in L(G)$, such that there are **two or more leftmost derivations for ω** (or equivalently two or more rightmost derivations) or equivalently **two or more parse trees for ω** .

Consider the Grammar G with the following rules: $S \rightarrow 0S1|SS|\epsilon$

Show that Grammar G is ambiguous, i.e. $\exists \omega \in L(G)$, such that there are two or more parse trees for ω .

Consider the string $\omega = 010101$:

- Show that there exist two different parse trees for 010101 .
- Show that there exist two leftmost derivations for 010101 .



Leftmost Derivation: $S \rightarrow SS \rightarrow 0S1S$

Parse trees for CFG

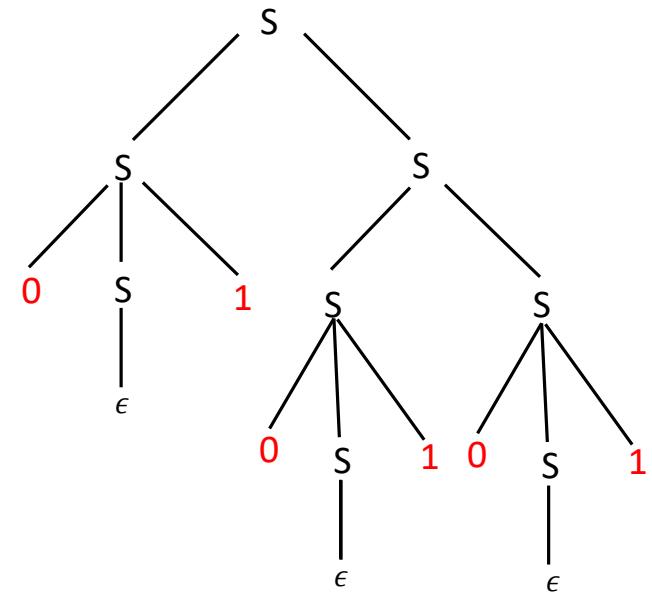
Ambiguous grammars: A CFG G is said to be **ambiguous** if there exists $\omega \in L(G)$, such that there are **two or more leftmost derivations for ω** (or equivalently two or more rightmost derivations) or equivalently **two or more parse trees for ω** .

Consider the Grammar G with the following rules: $S \rightarrow 0S1|SS|\epsilon$

Show that Grammar G is ambiguous, i.e. $\exists \omega \in L(G)$, such that there are two or more parse trees for ω .

Consider the string $\omega = 010101$:

- Show that there exist two different parse trees for 010101 .
- Show that there exist two leftmost derivations for 010101 .



Leftmost Derivation: $S \rightarrow SS \rightarrow 0S1S \rightarrow 01S$

Parse trees for CFG

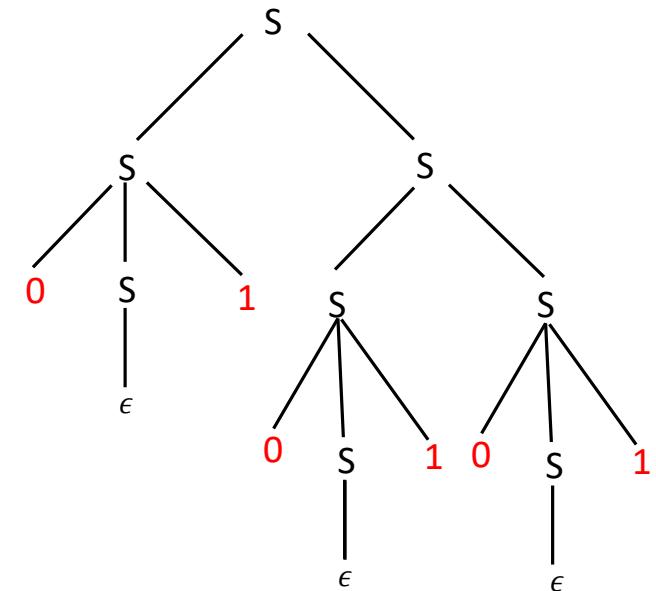
Ambiguous grammars: A CFG G is said to be **ambiguous** if there exists $\omega \in L(G)$, such that there are **two or more leftmost derivations for ω** (or equivalently two or more rightmost derivations) or equivalently **two or more parse trees for ω** .

Consider the Grammar G with the following rules: $S \rightarrow 0S1|SS|\epsilon$

Show that Grammar G is ambiguous, i.e. $\exists \omega \in L(G)$, such that there are two or more parse trees for ω .

Consider the string $\omega = 010101$:

- Show that there exist two different parse trees for 010101 .
- Show that there exist two leftmost derivations for 010101 .



Leftmost Derivation: $S \rightarrow SS \rightarrow 0S1S \rightarrow 01S$

Parse trees for CFG

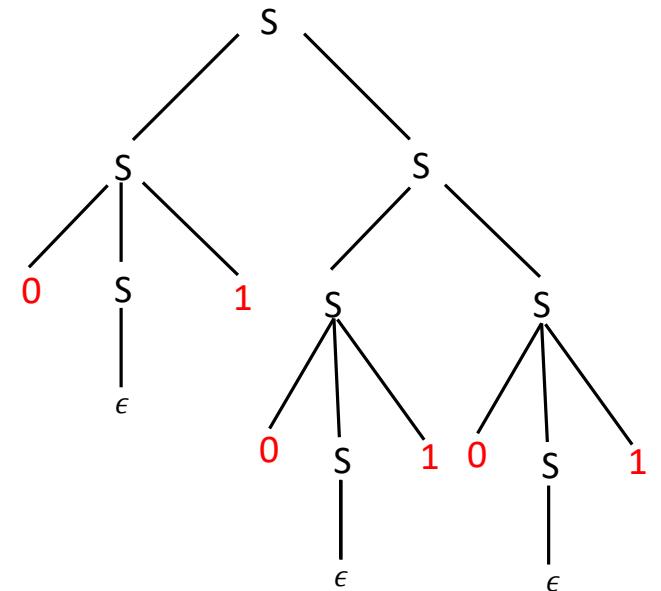
Ambiguous grammars: A CFG G is said to be **ambiguous** if there exists $\omega \in L(G)$, such that there are **two or more leftmost derivations for ω** (or equivalently two or more rightmost derivations) or equivalently **two or more parse trees for ω** .

Consider the Grammar G with the following rules: $S \rightarrow 0S1|SS|\epsilon$

Show that Grammar G is ambiguous, i.e. $\exists \omega \in L(G)$, such that there are two or more parse trees for ω .

Consider the string $\omega = 010101$:

- Show that there exist two different parse trees for 010101 .
- Show that there exist two leftmost derivations for 010101 .



Leftmost Derivation: $S \rightarrow SS \rightarrow 0S1S \rightarrow 01S \rightarrow 01SS$

Parse trees for CFG

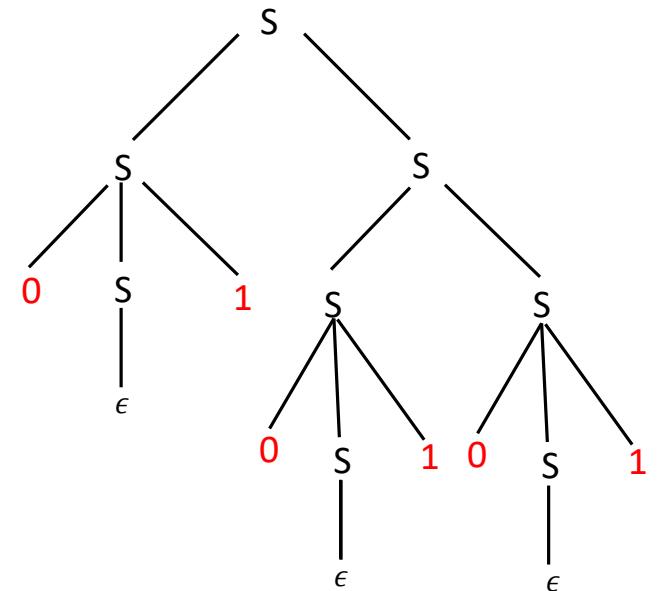
Ambiguous grammars: A CFG G is said to be **ambiguous** if there exists $\omega \in L(G)$, such that there are **two or more leftmost derivations for ω** (or equivalently two or more rightmost derivations) or equivalently **two or more parse trees for ω** .

Consider the Grammar G with the following rules: $S \rightarrow 0S1|SS|\epsilon$

Show that Grammar G is ambiguous, i.e. $\exists \omega \in L(G)$, such that there are two or more parse trees for ω .

Consider the string $\omega = 010101$:

- Show that there exist two different parse trees for 010101 .
- Show that there exist two leftmost derivations for 010101 .



Leftmost Derivation: $S \rightarrow SS \rightarrow 0S1S \rightarrow 01S \rightarrow 01SS$

Parse trees for CFG

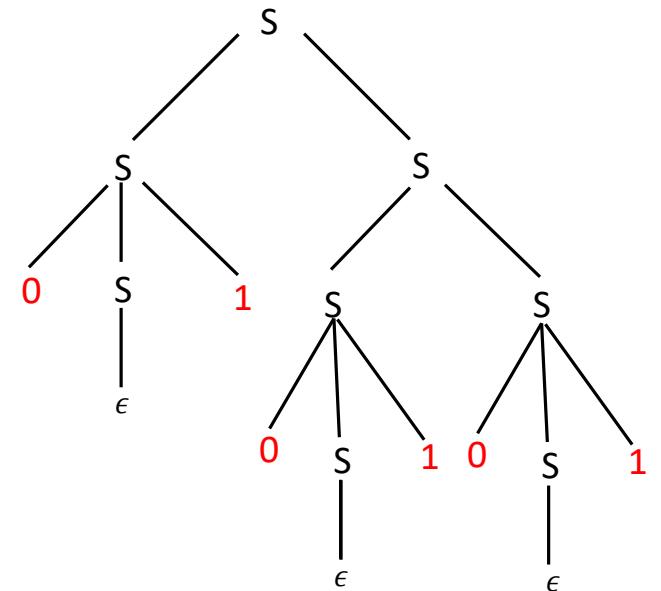
Ambiguous grammars: A CFG G is said to be **ambiguous** if there exists $\omega \in L(G)$, such that there are **two or more leftmost derivations for ω** (or equivalently two or more rightmost derivations) or equivalently **two or more parse trees for ω** .

Consider the Grammar G with the following rules: $S \rightarrow 0S1|SS|\epsilon$

Show that Grammar G is ambiguous, i.e. $\exists \omega \in L(G)$, such that there are two or more parse trees for ω .

Consider the string $\omega = 010101$:

- Show that there exist two different parse trees for 010101 .
- Show that there exist two leftmost derivations for 010101 .



Leftmost Derivation: $S \rightarrow SS \rightarrow 0S1S \rightarrow 01S \rightarrow 01SS \rightarrow 010S1S$

Parse trees for CFG

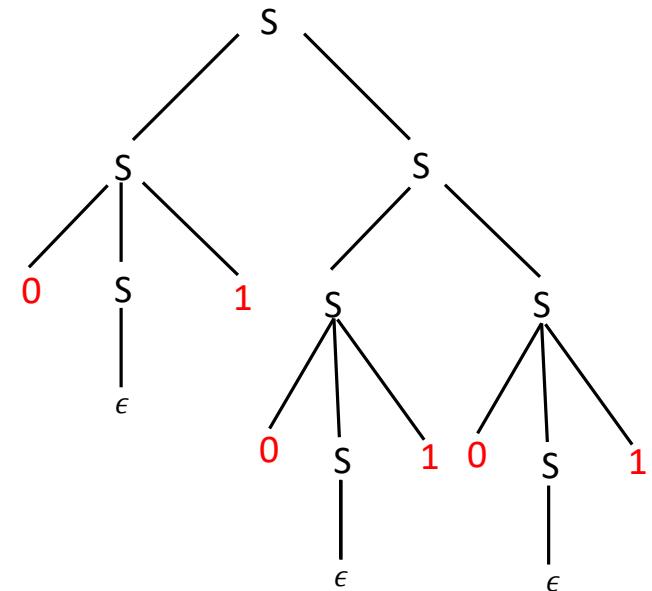
Ambiguous grammars: A CFG G is said to be **ambiguous** if there exists $\omega \in L(G)$, such that there are **two or more leftmost derivations for ω** (or equivalently two or more rightmost derivations) or equivalently **two or more parse trees for ω** .

Consider the Grammar G with the following rules: $S \rightarrow 0S1|SS|\epsilon$

Show that Grammar G is ambiguous, i.e. $\exists \omega \in L(G)$, such that there are two or more parse trees for ω .

Consider the string $\omega = 010101$:

- Show that there exist two different parse trees for 010101 .
- Show that there exist two leftmost derivations for 010101 .



Leftmost Derivation: $S \rightarrow SS \rightarrow 0S1S \rightarrow 01S \rightarrow 01SS \rightarrow 010\textcolor{red}{S}1S$

Parse trees for CFG

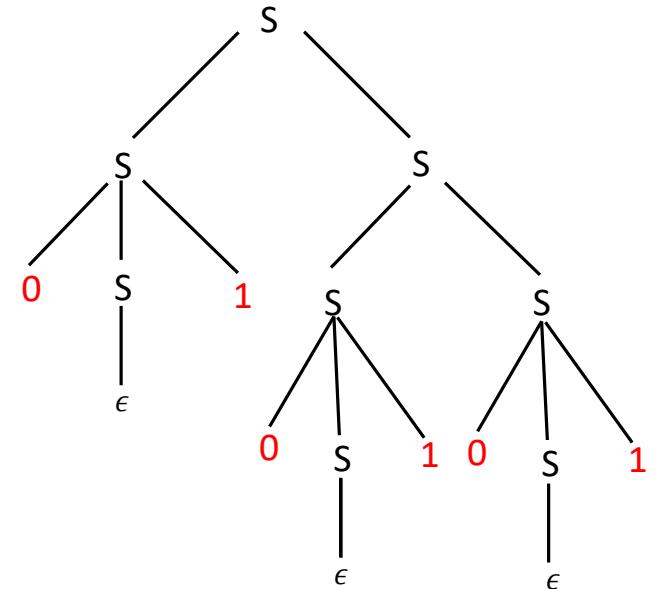
Ambiguous grammars: A CFG G is said to be **ambiguous** if there exists $\omega \in L(G)$, such that there are **two or more leftmost derivations for ω** (or equivalently two or more rightmost derivations) or equivalently **two or more parse trees for ω** .

Consider the Grammar G with the following rules: $S \rightarrow 0S1|SS|\epsilon$

Show that Grammar G is ambiguous, i.e. $\exists \omega \in L(G)$, such that there are two or more parse trees for ω .

Consider the string $\omega = 010101$:

- Show that there exist two different parse trees for 010101 .
- Show that there exist two leftmost derivations for 010101 .



Leftmost Derivation: $S \rightarrow SS \rightarrow 0S1S \rightarrow 01S \rightarrow 01SS \rightarrow 010S1S \rightarrow 0101S$

Parse trees for CFG

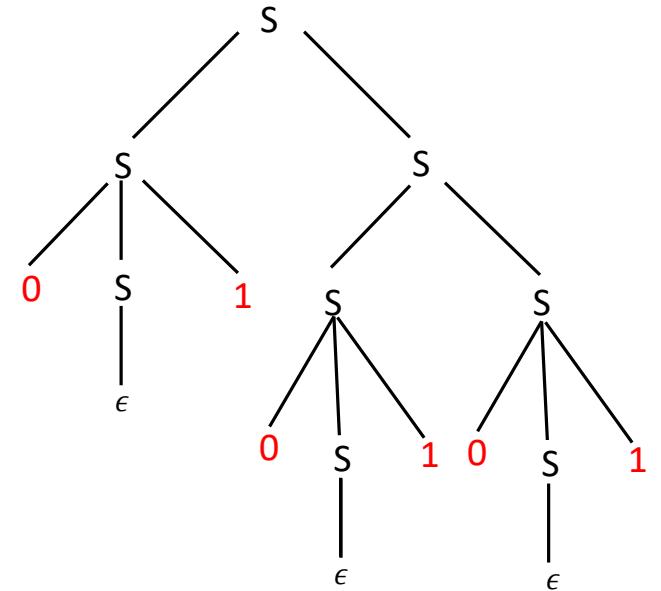
Ambiguous grammars: A CFG G is said to be **ambiguous** if there exists $\omega \in L(G)$, such that there are **two or more leftmost derivations for ω** (or equivalently two or more rightmost derivations) or equivalently **two or more parse trees for ω** .

Consider the Grammar G with the following rules: $S \rightarrow 0S1|SS|\epsilon$

Show that Grammar G is ambiguous, i.e. $\exists \omega \in L(G)$, such that there are two or more parse trees for ω .

Consider the string $\omega = 010101$:

- Show that there exist two different parse trees for 010101 .
- Show that there exist two leftmost derivations for 010101 .



Leftmost Derivation: $S \rightarrow SS \rightarrow 0S1S \rightarrow 01S \rightarrow 01SS \rightarrow 010S1S \rightarrow 0101\textcolor{red}{S}$

Parse trees for CFG

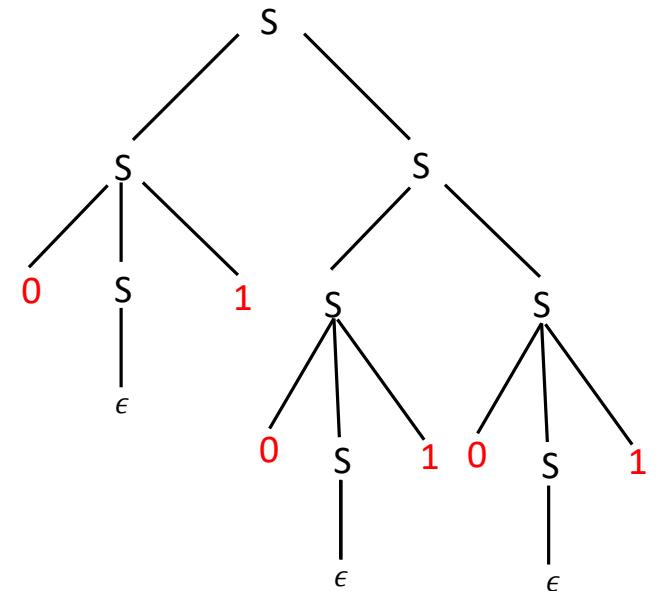
Ambiguous grammars: A CFG G is said to be **ambiguous** if there exists $\omega \in L(G)$, such that there are **two or more leftmost derivations for ω** (or equivalently two or more rightmost derivations) or equivalently **two or more parse trees for ω** .

Consider the Grammar G with the following rules: $S \rightarrow 0S1|SS|\epsilon$

Show that Grammar G is ambiguous, i.e. $\exists \omega \in L(G)$, such that there are two or more parse trees for ω .

Consider the string $\omega = 010101$:

- Show that there exist two different parse trees for 010101 .
- Show that there exist two leftmost derivations for 010101 .



Leftmost Derivation: $S \rightarrow SS \rightarrow 0S1S \rightarrow 01S \rightarrow 01SS \rightarrow 010S1S \rightarrow 0101S \rightarrow 01010S1$

Parse trees for CFG

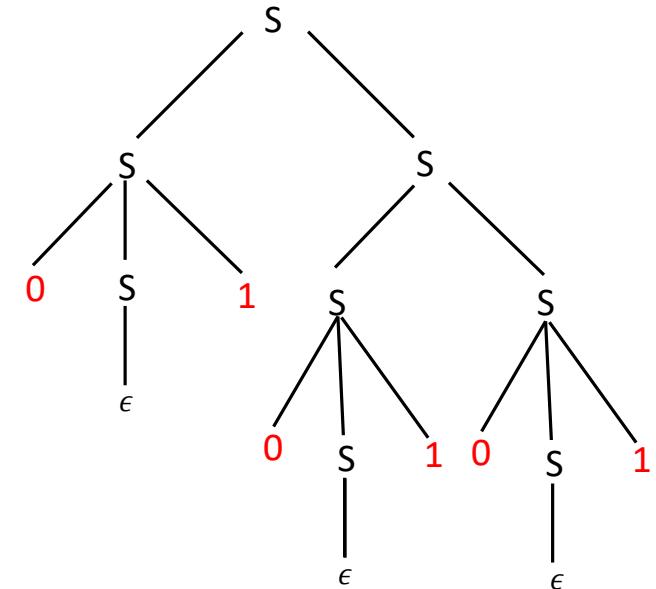
Ambiguous grammars: A CFG G is said to be **ambiguous** if there exists $\omega \in L(G)$, such that there are **two or more leftmost derivations for ω** (or equivalently two or more rightmost derivations) or equivalently **two or more parse trees for ω** .

Consider the Grammar G with the following rules: $S \rightarrow 0S1|SS|\epsilon$

Show that Grammar G is ambiguous, i.e. $\exists \omega \in L(G)$, such that there are two or more parse trees for ω .

Consider the string $\omega = 010101$:

- Show that there exist two different parse trees for 010101 .
- Show that there exist two leftmost derivations for 010101 .



Leftmost Derivation: $S \rightarrow SS \rightarrow 0S1S \rightarrow 01S \rightarrow 01SS \rightarrow 010S1S \rightarrow 0101S \rightarrow 01010S1 \rightarrow 010101$

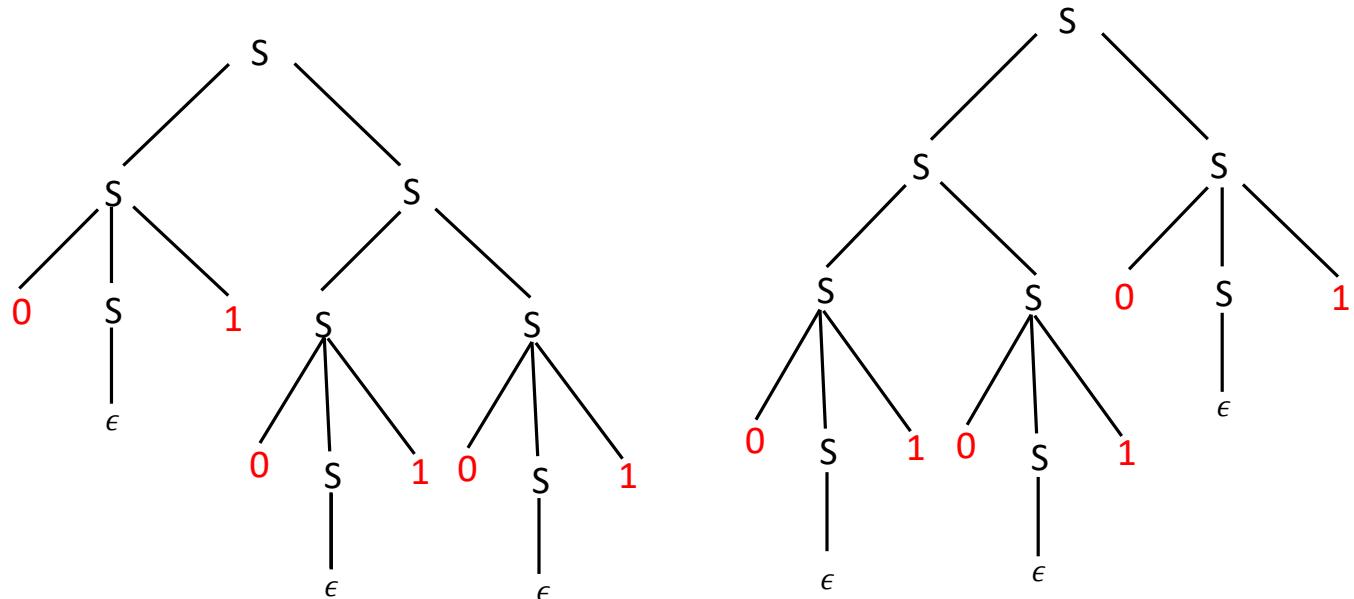
Parse trees for CFG

Ambiguous grammars: A CFG G is said to be **ambiguous** if there exists $\omega \in L(G)$, such that there are **two or more leftmost derivations for ω** (or equivalently two or more rightmost derivations) or equivalently **two or more parse trees for ω** .

Consider the Grammar G with the following rules: $S \rightarrow 0S1|SS|\epsilon$

Consider the string $\omega = 010101$:

- Show that there exist two different parse trees for 010101 .
- Show that there exist two leftmost derivations for 010101 .

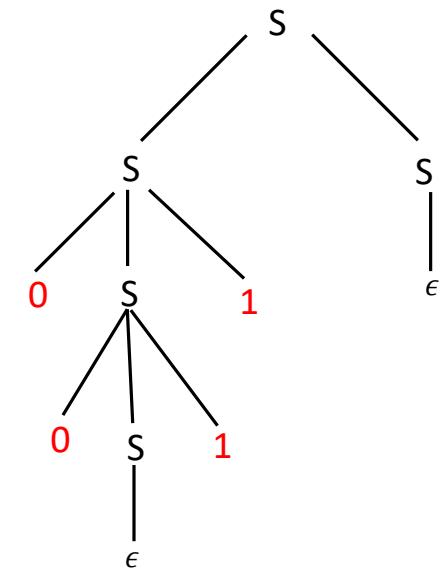
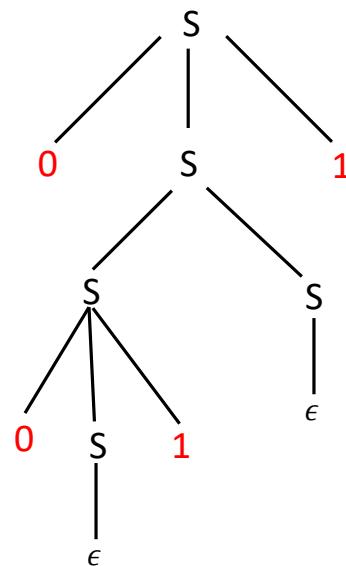
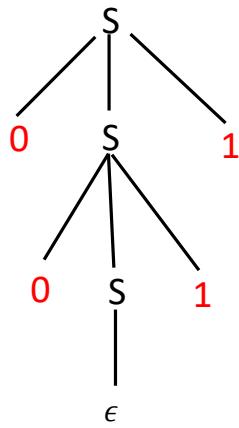


Leftmost Derivation: $S \rightarrow SS \rightarrow SSS \rightarrow 0S1SS \rightarrow 01SS \rightarrow 010S1S \rightarrow 0101S \rightarrow 01010S1 \rightarrow 010101$

Parse trees for CFG

Show that the Grammar G with the following rules: $S \rightarrow 0S1|SS|\epsilon$ is ambiguous.

Consider string $\omega = 0011$



LD: $S \rightarrow 0S1 \rightarrow 00S11 \rightarrow 0011$

LD: $S \rightarrow 0S1 \rightarrow 0SS1 \rightarrow 00S1S1 \rightarrow 001S1 \rightarrow 0011$

LD: $S \rightarrow SS \rightarrow 0S1S \rightarrow 00S11S \rightarrow 0011S \rightarrow 0011$

Ambiguity

Unique structures are important. For example:

- The syntax of a programming language can be represented by a CFG.
- A compiler
 - translates the code written in the programming language into a form that is suitable for execution.
 - checks if the underlying programming language is syntactically correct.
- Parse trees are data structures that represent such structures.
- Parse tree for the code helps analyze the syntax. So ambiguity might lead to different interpretations and hence, different outcomes for the same code.

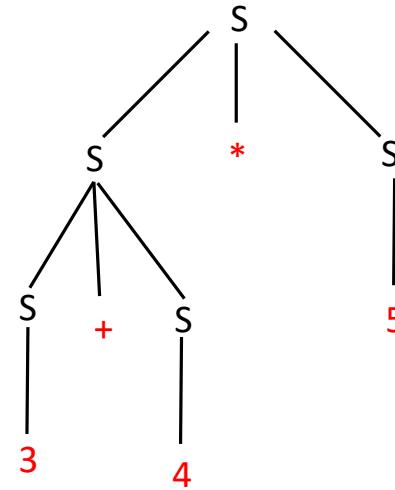
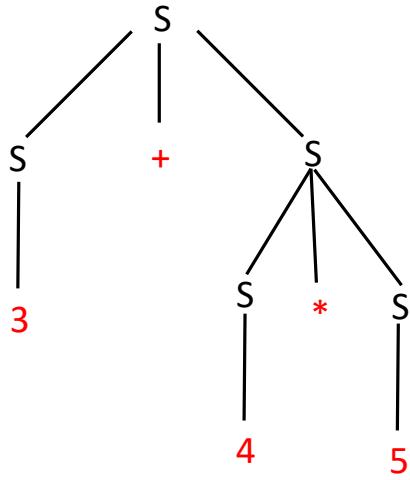
Ambiguity may not be desirable.

Ambiguity

Ambiguity may not be desirable.

Consider the grammar: $S \rightarrow S + S \mid S * S \mid 0 \mid 1 \mid 2 \mid \dots \mid 9$

and the derivation of the string $3 + 4 * 5$



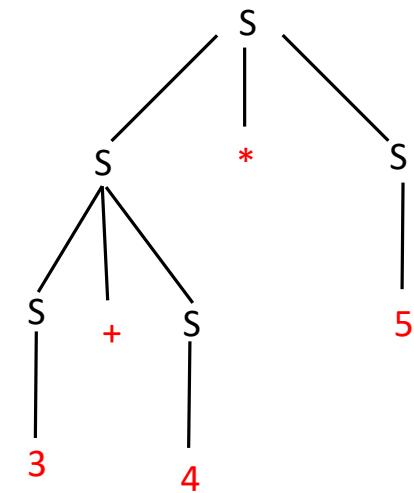
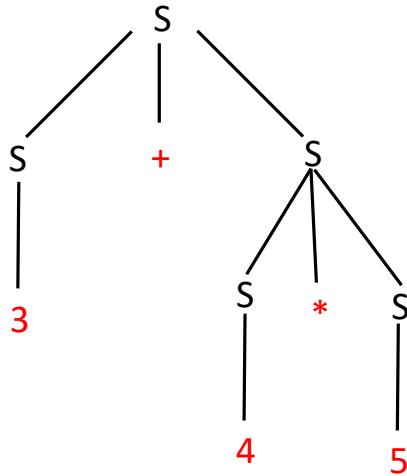
- The grammar contains no information on the precedence relations of the various arithmetic operations.
- The grammar may group $+$ before $*$

Ambiguity

Ambiguity may not be desirable.

Consider the grammar: $S \rightarrow S + S \mid S * S \mid 0 \mid 1 \mid 2 \mid \dots \mid 9$

and the derivation of the string $3 + 4 * 5$



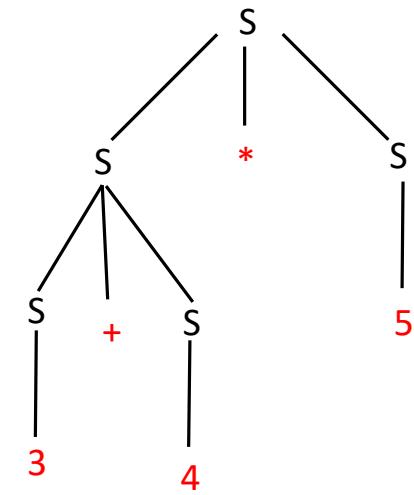
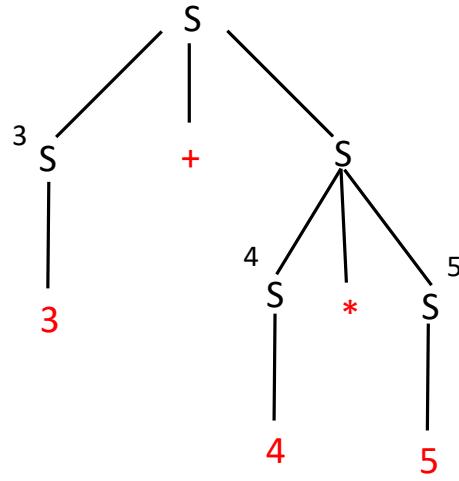
- What will be the result obtained from each of these *parsings*?

Ambiguity

Ambiguity may not be desirable.

Consider the grammar: $S \rightarrow S + S \mid S * S \mid 0 \mid 1 \mid 2 \mid \dots \mid 9$

and the derivation of the string $3 + 4 * 5$



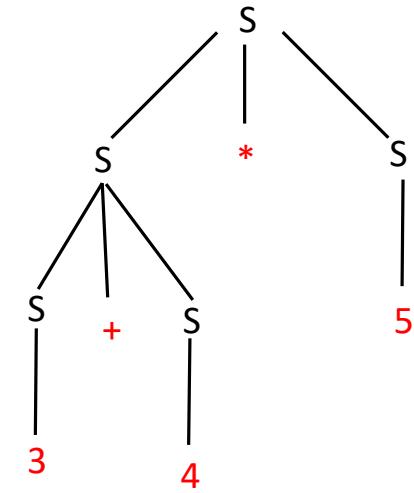
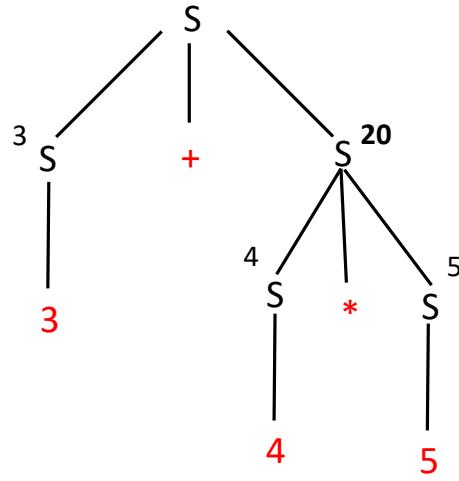
- If the compiler compiles the left parse tree

Ambiguity

Ambiguity may not be desirable.

Consider the grammar: $S \rightarrow S + S \mid S * S \mid 0 \mid 1 \mid 2 \mid \dots \mid 9$

and the derivation of the string $3 + 4 * 5$



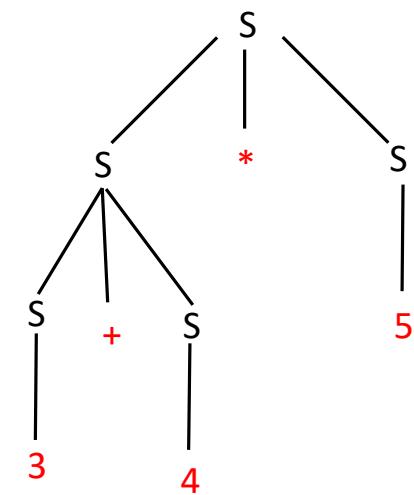
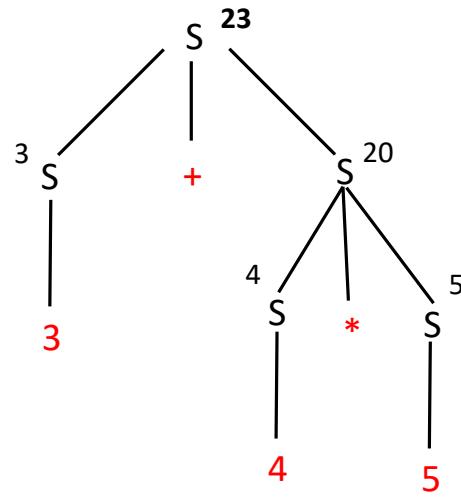
- If the compiler compiles the left parse tree

Ambiguity

Ambiguity may not be desirable.

Consider the grammar: $S \rightarrow S + S \mid S * S \mid 0 \mid 1 \mid 2 \mid \dots \mid 9$

and the derivation of the string $3 + 4 * 5$



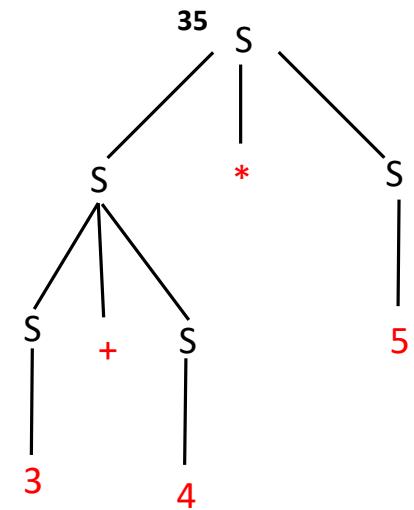
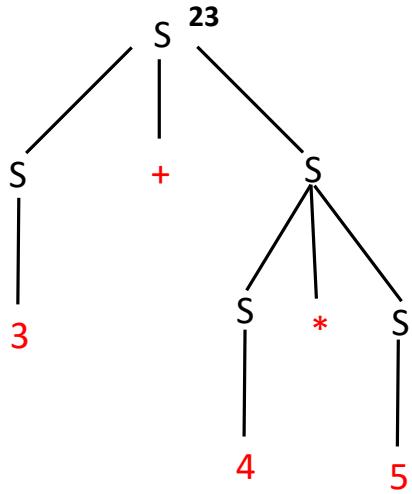
- If the compiler compiles the left parse tree. Outcome = 23

Ambiguity

Ambiguity may not be desirable.

Consider the grammar: $S \rightarrow S + S \mid S * S \mid 0 \mid 1 \mid 2 \mid \dots \mid 9$

and the derivation of the string $3 + 4 * 5$



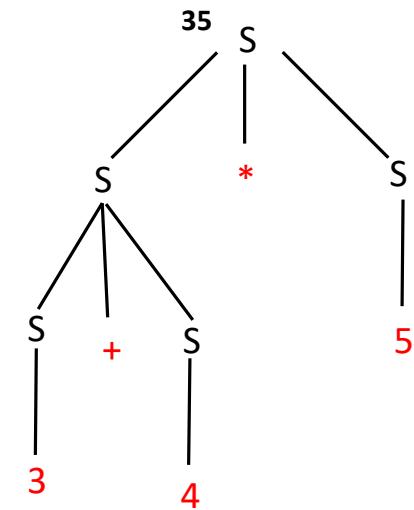
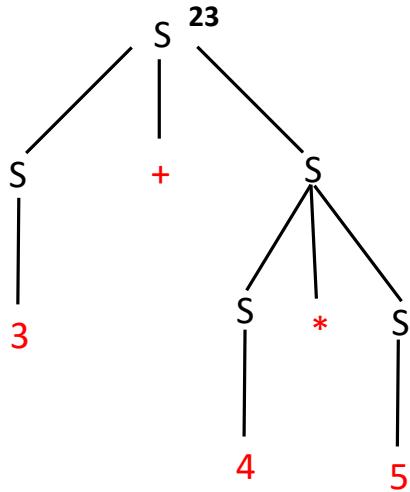
- If the compiler compiles the **right** parse tree. Outcome = **35**

Ambiguity

Ambiguity may not be desirable.

Consider the grammar: $S \rightarrow S + S \mid S * S \mid 0 \mid 1 \mid 2 \mid \dots \mid 9$

and the derivation of the string $3 + 4 * 5$



- How can we get rid of this ambiguity?

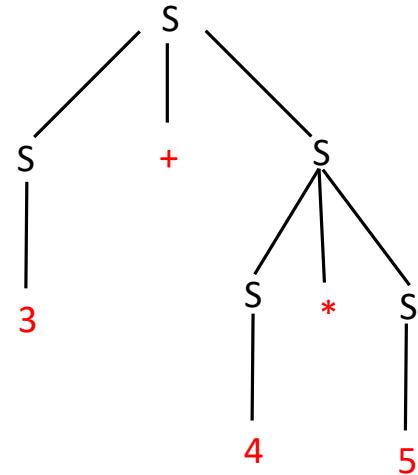
Ambiguity

Consider the grammar: $S \rightarrow S + S \mid S * S \mid 0 \mid 1 \mid 2 \mid \dots \mid 9$

How can we get rid of this ambiguity? Change the production rules

1) Add parenthesis

New Grammar: $S \rightarrow (S + S) \mid (S * S) \mid 0 \mid 1 \mid 2 \mid \dots \mid 9$



Old Parse tree (before adding parenthesis)

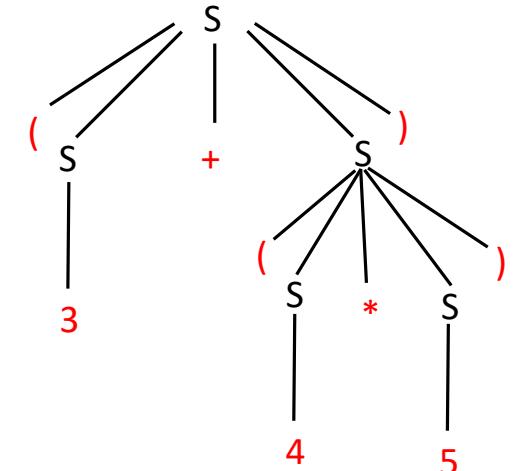
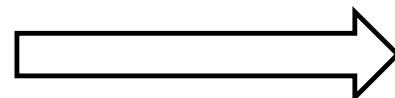
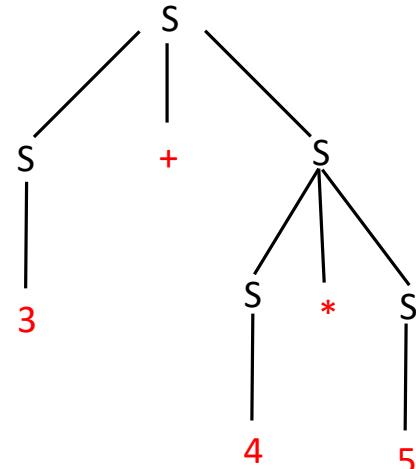
Ambiguity

Consider the grammar: $S \rightarrow S + S \mid S * S \mid 0 \mid 1 \mid 2 \mid \dots \mid 9$

How can we get rid of this ambiguity? Change the production rules

1) Add parenthesis

New Grammar: $S \rightarrow (S + S) \mid (S * S) \mid 0 \mid 1 \mid 2 \mid \dots \mid 9$



$$(3 + (4 * 5)) = 23$$

Ambiguity

Consider the grammar: $S \rightarrow S + S \mid S * S \mid 0 \mid 1 \mid 2 \mid \dots \mid 9$

How can we get rid of this ambiguity? Change the production rules

- 1) Add parentheses
- 2) Add new variables

Ambiguity

Consider the grammar: $S \rightarrow S + S \mid S * S \mid 0 \mid 1 \mid 2 \mid \dots \mid 9$

How can we get rid of this ambiguity? Change the production rules

- 1) Add parentheses
- 2) Add new variables

New Grammar:

$$\begin{aligned} E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow 0 \mid 1 \mid 2 \mid \dots \mid 9 \mid E \end{aligned}$$

Ambiguity

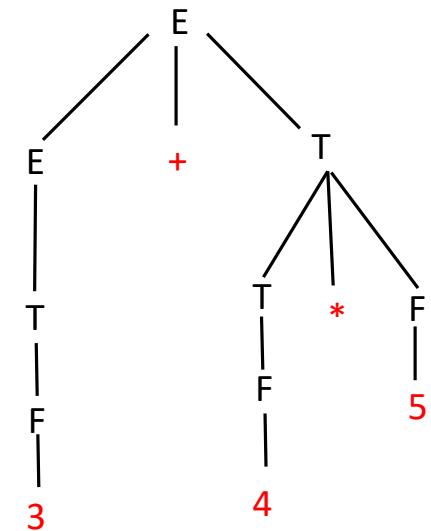
How can we get rid of this ambiguity? Change the production rules

- 1) Add parentheses
- 2) Add new variables

New Grammar:

$$\begin{aligned} E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow 0 \mid 1 \mid 2 \mid \dots \mid 9 \mid E \end{aligned}$$

Parse tree to derive: $3 + (4 * 5)$



Ambiguity

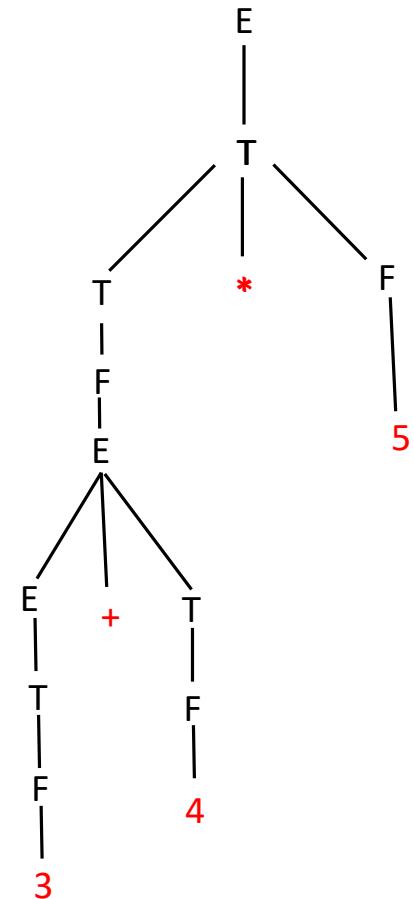
How can we get rid of this ambiguity? Change the production rules

- 1) Add parentheses
- 2) Add new variables

New Grammar:

$$\begin{aligned} E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow 0 \mid 1 \mid 2 \mid \dots \mid 9 \mid E \end{aligned}$$

Parse tree to derive: $(3 + 4) * 5$



Ambiguity

How can we get rid of this ambiguity? Change the production rules

- 1) Add parentheses
- 2) Add new variables

- In general, it is not possible to write an algorithm that takes as input a grammar G and outputs, YES if G is ambiguous and NO, otherwise. (**Undecidable**)
- So removing ambiguity is impossible in general.

Chomsky Normal Form

Often it is easier to work with CFG in a simple standardized form - the Chomsky Normal Form (CNF) is one of them.

Chomsky Normal Form

A CFG G is in CNF if every rule of G is of the form

$$\begin{aligned}Var &\rightarrow Var\ Var \\Var &\rightarrow ter \\Start\ Var &\rightarrow \epsilon\end{aligned}$$

where Var can be any variable, including the Start Variable, $Start\ Var$.

Chomsky Normal Form

Often it is easier to work with CFG in a simple standardized form - the Chomsky Normal Form (CNF) is one of them.

Chomsky Normal Form

A CFG G is in CNF if every rule of G is of the form

$$\begin{aligned}Var &\rightarrow Var \; Var \\Var &\rightarrow ter \\Start \; Var &\rightarrow \epsilon\end{aligned}$$

where Var can be any variable, including the Start Variable, $Start \; Var$.

Why are CNFs useful?

- Suppose you are given a CFG G and a string w as input and you have to write an algorithm that decides whether G generates w .
- Your algorithm outputs YES if G generates w and NO, otherwise.

Chomsky Normal Form

A CFG G is in **CNF** if every rule of G is of the form

$$\begin{aligned}Var &\rightarrow Var \; Var \\Var &\rightarrow ter \\Start \; Var &\rightarrow \epsilon\end{aligned}$$

where Var can be any variable, including the Start Variable, $Start \; Var$.

Why are CNFs useful?

- Suppose you are given a CFG G as and a string w as input and you have to **write an algorithm that decides whether G generates w** .
 - The algorithm outputs YES if G generates w and NO, otherwise.
 - ❖ One idea is to go through ALL derivations one by one and output YES if any of them generates w .
 - ❖ However, infinitely many derivations may have to tried.
 - ❖ So if G does not generate w , the algorithm will never stop.
 - ❖ So this problem appears to be **undecidable**.

Chomsky Normal Form

A CFG G is in **CNF** if every rule of G is of the form

$$\begin{aligned}Var &\rightarrow Var \; Var \\Var &\rightarrow ter \\Start \; Var &\rightarrow \epsilon\end{aligned}$$

where Var can be any variable, including the Start Variable, $Start \; Var$.

Why are CNFs useful?

Suppose you are given a CFG G and a string w as input and you have to **write an algorithm that decides whether G generates w** . This problem appears to be **undecidable**.

Chomsky Normal Form

A CFG G is in **CNF** if every rule of G is of the form

$$\begin{aligned}Var &\rightarrow Var \; Var \\Var &\rightarrow ter \\Start \; Var &\rightarrow \epsilon\end{aligned}$$

where Var can be any variable, including the Start Variable, $Start \; Var$.

Why are CNFs useful?

Suppose you are given a CFG G as and a string w as input and you have to **write an algorithm that decides whether G generates w** .

- Converting G first to a CNF alleviates this and **makes the problem decidable**.
- It limits the number of steps in derivations required to generate any $w \in L(G)$.
- If $w \in L(G)$, then a CFG in Chomsky Normal Form has **derivations of $2n - 1$ steps** for input strings w of length n (We will prove this shortly).

Chomsky Normal Form

A CFG G is in **CNF** if every rule of G is of the form

$$\begin{aligned}Var &\rightarrow Var \; Var \\Var &\rightarrow ter \\Start \; Var &\rightarrow \epsilon\end{aligned}$$

where Var can be any variable, including the Start Variable, $Start \; Var$.

A CFG in Chomsky Normal Form has derivations of $2n - 1$ steps for generating strings $w \in L(G)$ of length n .

Why are CNFs useful?

Suppose you are given a CFG G as and a string w as input and you have to **write an algorithm that decides whether G generates w** .

1. Convert G to CNF.
2. List all derivations of $2n - 1$ steps, where $|w| = n$. (There are a finite number of these)
3. If ANY of these derivations generate w , output YES, otherwise output NO.

Chomsky Normal Form

A CFG G is in **CNF** if every rule of G is of the form

$$\begin{aligned}Var &\rightarrow Var \; Var \\Var &\rightarrow ter \\Start \; Var &\rightarrow \epsilon\end{aligned}$$

where Var can be any variable, including the Start Variable, $Start \; Var$.

- 1) A CFG in Chomsky Normal Form has derivations of $2n - 1$ steps for generating strings $w \in L(G)$ of length n .
- 2) Any CFL can be generated by a CFG written in Chomsky Normal Form.

To prove 1) use induction!

Chomsky Normal Form

Prove that a CFG in Chomsky Normal Form has derivations of $2n - 1$ steps for generating strings $w \in L(G)$ of length n .

Proof: Note that any CFG in CNF can be written as:

$$\begin{array}{ll} A \rightarrow BC & [B, C \text{ are not start variables}] \\ A \rightarrow a & [a \text{ is a terminal}] \\ S \rightarrow \epsilon & [S \text{ is the Start Variable}] \end{array}$$

We will prove this by **induction**.

(Basic step) Let $|w| = 1$. Then **one** application of the second rule would suffice. So any derivation of w would need $2|w| - 1 = 1$ step.

(Inductive hypothesis) Assume the statement of the theorem to be true for any string of length at most k where $k \geq 1$. Now we shall show that it holds for any $w \in L(G)$ such that $|w| = k + 1$.

Chomsky Normal Form

Prove that a CFG in Chomsky Normal Form has derivations of $2n - 1$ steps for generating strings $w \in L(G)$ of length n .

Proof: Note that any CFG in CNF can be written as:

$$\begin{array}{ll} A \rightarrow BC & [B, C \text{ are not start variables}] \\ A \rightarrow a & [a \text{ is a terminal}] \\ S \rightarrow \epsilon & [S \text{ is the Start Variable}] \end{array}$$

We will prove this by **induction**.

(Basic step) Let $|w| = 1$. Then **one** application of the second rule would suffice. So any derivation of w would need $2|w| - 1 = 1$ step.

(Inductive hypothesis) Assume the statement of the theorem to be true for any string of length at most k where $k \geq 1$. Now we shall show that it holds for any $w \in L(G)$ such that $|w| = k + 1$.

Since $|w| > 1$, any derivation will start from the rule $A \rightarrow BC$. So $w = xy$, where $B \xrightarrow{*} x$, $|x| > 0$ and $C \xrightarrow{*} y$, $|y| > 0$. But since $|x|, |y| \leq k$, and we have that by the inductive hypothesis: (i) number of steps in the derivation $B \xrightarrow{*} x$ is $2|x| - 1$ and (ii) number of steps in the derivation $C \xrightarrow{*} y$ is $2|y| - 1$. So the number of steps in the derivation of w is

$$1 + (2|x| - 1) + (2|y| - 1) = 2(|x| + |y|) - 1 = 2|w| - 1 = 2(k + 1) - 1.$$

Chomsky Normal Form

A CFG G is in **CNF** if every rule of G is of the form

$$\begin{aligned}Var &\rightarrow Var\;Var \\Var &\rightarrow ter \\Start\;Var &\rightarrow \epsilon\end{aligned}$$

where Var can be any variable, including the Start Variable, $Start\;Var$.

- 1) A CFG in Chomsky Normal Form has derivations of $2n - 1$ steps for generating strings $w \in L(G)$ of length n .
- 2) Any CFL can be generated by a CFG written in Chomsky Normal Form.

Thank You!