

## Tutorial 2: Solutions

**Exercise 1:** The alphabet  $\Sigma$  is the set of all printable ASCII characters:

!"#\$%&'()\*+,-./0123456789:;<=>?@ABCDEFGHIJKLMNO  
PQRSTUVWXYZ[\]^\_`abcdefghijklmnopqrstuvwxyz{|}~

For each of the following languages, find an accepting grammar:

1. the decimal representations of odd numbers
2. the set of palindromes (words that are equal to their mirrored image)
3. the *arithmetic expressions*: every nonempty sequence of digits is an arithmetic expression, if  $e_1$  and  $e_2$  are two arithmetic expressions, then so are the words  $e_1+e_2$ ,  $e_1-e_2$ ,  $e_1*e_2$ ,  $e_1/e_2$ , and  $(e_1)$ .

**Solution:**

1. The grammar  $(\Sigma, \mathcal{N}, S, \mathcal{P})$  where  $\Sigma$  is the set of ASCII characters,  $\mathcal{N} = \{S, O, M\}$ , and  $\mathcal{P}$  contains the following rules:

$$\begin{aligned} S &\rightarrow MO & O &\rightarrow 1 \mid 3 \mid 5 \mid 7 \mid 9 \\ M &\rightarrow 0M \mid 1M \mid 2M \mid 3M \mid 4M \mid 5M \mid 6M \mid 7M \mid 8M \mid 9M \mid \varepsilon \end{aligned}$$

accepts the decimal representation of an odd number.

The grammar only produces words containing digits and the last digits of the word are odd numbers (since  $S \rightarrow MO$ , and the non-terminal  $O$  only produces odd numbers). Then, from the non-terminal  $M$  we can either derive  $\varepsilon$  or any digit.

**Note:** The notation  $S \rightarrow A \mid B$  is equivalent to the (union of) the rules  $S \rightarrow A$  and  $S \rightarrow B$ .

2. The grammar  $(\Sigma, \mathcal{N}, S, \mathcal{P})$  where  $\Sigma$  is the set of ASCII characters,  $\mathcal{N} = \{S\}$ , and  $\mathcal{P}$  is:

$$S \rightarrow \varepsilon \mid \alpha S \alpha \mid \alpha \quad \text{for all } \alpha \in \Sigma$$

accepts the language of palindrome letters.

The grammar accept  $\varepsilon$  ( $S \rightarrow \varepsilon$ ), it can derive a palindrome word with an even number of letters applying multiple times the rule  $S \rightarrow \alpha A \alpha$  (e.g.,  $S \rightarrow aAa \rightarrow abAba \rightarrow abba$ ), or it can derive a palindrome word with an odd number of letters when applying the rule  $S \rightarrow \alpha$  (e.g.,  $S \rightarrow a$ ,  $S \rightarrow aAa \rightarrow abAba \rightarrow abcba$ ).

3. The grammar  $(\Sigma, \mathcal{N}, S, \mathcal{P})$  where  $\Sigma$  is the set of ASCII characters,  $\mathcal{N} = \{S, D\}$ , and  $\mathcal{P}$  is:

$$\begin{aligned} S &\rightarrow D \mid S+S \mid S-S \mid S/S \mid S*S \mid (S) \\ D &\rightarrow 0D' \mid 1D' \mid 2D' \mid \dots \mid 9D' \\ D' &\rightarrow D \mid \varepsilon \end{aligned}$$

accepts the language of arithmetic expressions.

The production rules of the grammar follows the inductive definition of an arithmetic expression. An arithmetic expression is either a sequence of a digit (rule  $S \rightarrow D$ ,  $D'$  is used to end the construction of the Digit at some point), or is obtained as a binary operation on two other arithmetic expressions (rule  $S \rightarrow S$ ), or is an arithmetic expression enclosed in brackets (rule  $(S)$ ).

**Exercise 2:** Let  $L$  be the language accepted by the grammar  $(\Sigma, \mathcal{N}, S, \mathcal{P})$  where:

- $\Sigma = \{a, b, c\}$ ,  $\mathcal{N} = \{S, A, B, C, X, Y\}$ ,  $S$  is the start symbol, and
- the production rules (i.e. the elements of  $\mathcal{P}$ ) are:

$$\begin{aligned} S \rightarrow ABC \quad A \rightarrow AaX \quad Xa \rightarrow aX \quad XB \rightarrow BbX \quad Xb \rightarrow bX \quad XC \rightarrow Cc \\ A \rightarrow Y \quad Ya \rightarrow aY \quad Yb \rightarrow bY \quad YB \rightarrow Y \quad YC \rightarrow \epsilon \end{aligned}$$

Find two or three words that belong to  $L$  then guess what  $L$  is.

**Solution:** We can derive the following words applying the rule the grammar:

- $S \rightarrow ABC \rightarrow YBC \rightarrow YC \rightarrow \epsilon$   
Applying in order the rules  $S \rightarrow ABC$ ,  $A \rightarrow Y$ ,  $YB \rightarrow Y$ ,  $YC \rightarrow \epsilon$ .
- $S \rightarrow ABC \rightarrow AaXBC \rightarrow AaBbXC \rightarrow AaBbCc \rightarrow YaBbCc \rightarrow aYBbCc \rightarrow aYbCc \rightarrow abYCC \rightarrow abc$   
Applying in order the rules:  $S \rightarrow ABC$ ,  $A \rightarrow AaX$ ,  $XB \rightarrow BbX$ ,  $XC \rightarrow Cc$ ,  $A \rightarrow Y$ ,  $Ya \rightarrow aY$ ,  $YB \rightarrow Y$ ,  $Yb \rightarrow bY$ ,  $YC \rightarrow \epsilon$ .
- $S \rightarrow ABC \rightarrow AaXBC \rightarrow AaBbXC \rightarrow AaXaBbXC \rightarrow AaaXBbXC \rightarrow AaaBbXbXC \rightarrow$   
 $AaaBbbXXC \rightarrow AaaBbbXCc \rightarrow AaaBbbCcc \rightarrow YaaBbbCcc \rightarrow aYaBbbCcc \rightarrow$   
 $aaYBbbCcc \rightarrow aaYbbCcc \rightarrow aabYbCcc \rightarrow aabbYCcc \rightarrow aabbcc$

Applying in order the rules:

$$\begin{aligned} S \rightarrow ABC, A \rightarrow AaX, XB \rightarrow BbX, A \rightarrow AaX, Xa \rightarrow aX, XB \rightarrow BbX, Xb \rightarrow bX, \\ XC \rightarrow Cc, XC \rightarrow Cc, A \rightarrow Y, aY \rightarrow Ya, aY \rightarrow Ya, \\ YB \rightarrow Y, Yb \rightarrow bY, Yb \rightarrow bY, YC \rightarrow \epsilon. \end{aligned}$$

The language accepted by the grammar is the set  $\{a^n b^n c^n \mid n \in \mathbb{N}\}$ . The non-terminal  $A$  generates an instance of  $a$  and “stores” that the grammar generated one  $a$  in the non-terminal  $X$ . Note that, in order to derive a “new”  $b$  character (a  $b$  that was not already in the current derivation), the grammar requires to match the subword  $XB$ . Similarly, one need to have  $XC$  to generate a new  $c$ . The grammar has also rules to “move” the non-terminal  $X$  on the right of the word (i.e., switching  $X$  with the letters  $a$  and  $b$ ). The non-terminal  $Y$  is used to remove the instances of the non-terminal  $A, B, C$  and derive a word containing only terminal symbols.

Let us prove it formally<sup>1</sup>. We prove that the following properties are invariant when applying each of the rule of  $\mathcal{P}$ . If  $w$  is a word over  $\Sigma \cup \mathcal{N}$ , and  $l \in \Sigma \cup \mathcal{N}$ , we let  $|w|_l$  be the number of  $l$ s in  $w$ .

1.  $S$  will never appear after the first application of a production rule. If  $A$  disappear, it will not disappear. Indeed, the only way to produce a new  $A$  is to use  $S \rightarrow ABC$  that can not be used twice. The same applies for  $B$  and  $C$ .

---

<sup>1</sup>We do not expect you to do it in an exam.

2. We have  $|w|_A \leq 1$ ,  $|w|_B \leq 1$  and  $|w|_C \leq 1$ . Indeed, it is the case at start ( $S$ ) and each of the rule preserve that property.
3. If  $|w|_A = 1$  then  $A$  is the first letter of  $w$ . Indeed, it is the case at start ( $|S|_A = 0$ ) and each rule put  $A$  at the beginning of the word if  $A$  is already present.
4. If  $|w|_B = 1$  then  $B$  is before any  $b$  in  $w$ .
5. If  $|w|_C = 1$  then  $C$  is before any  $c$  in  $w$ .
6. Any  $a$  is before  $B$  when  $|w|_B = 1$ . Assume it is not true. Let  $w$  be built from  $S$  with the rules  $\mathcal{P}$  and steps  $w_0 = S \rightarrow w_1 \rightarrow \dots \rightarrow w_n \rightarrow w$  such that  $w$  violates that property. Without loss of generality we assume  $w$  to be the first word in the sequence to violate the property. Since  $S$  does not violate the property, the sequence is well defined ( $n \geq 0$ ). We write  $w = u_1 B u_2 a u_3$  with  $u_i$  be (possibly empty) words over  $\Sigma \cup \mathcal{N}$ . We have the following cases
  - $w_n \rightarrow w$  implies  $B$ . If it is  $S \rightarrow ABC$ , using property 1. we have  $w = ABC$  which is not. If it is  $YB \rightarrow Y$ , then  $1 \geq |w|_B = |w_n|_B - 1 \leq 1 - 1 = 0$  using property 2., what is not. Then we have just applied  $XB \rightarrow BbX$ . Using 2.,  $w_n$  as only one  $B$ . We then know that  $w_n = u_1 X B u'_2 a u_3$  with  $u_2 = bX u'_2$ . Then  $w_n$  violates the property 6. and then contradicts the minimality of  $w$ .
  - $w_n \rightarrow w$  implies the  $a$  of  $w$  we are looking at. Then we do the same thing as the previous case to derive that  $w_n$  violates the minimality of  $w$ .
  - $w_n \rightarrow w$  does not implies the  $B$  of the  $a$ . Then  $w_n = u'_1 B u'_2 a u'_3$  with  $u'_i \rightarrow u_i$  for one  $i \in \{1, 2, 3\}$  and  $u'_j = u_j$  for  $j \neq i$ . Then  $w_n$  violates the minimality of  $w$ .

Then property 6. is an invariant.

7. Any  $b$  is before  $C$  when  $|w|_C = 1$ . We do the same proof as for property 6.
8.  $|w|_a = |w|_b + |w'|_X$  where  $w = w'w''$  and  $w'$  does not contains any  $b$  and is maximal for this property ( $w''$  is empty or starts with a  $b$ ). Indeed, the only rule that we need to pay attention is  $A \rightarrow AaX$ . But using property 3., we conclude that there is no problem.
9.  $|w|_b = |w|_c + |w''|_X$  where  $w = w'w''$  and  $w'$  does not contains any  $b$  and is maximal for this property ( $w''$  is empty or starts with a  $b$ ).
10. If  $w \neq S$  then every  $a$  is before any  $b$ . Indeed, if  $|w|_B = 1$  we conclude using 4. and 6. In the sequence to build  $w$ , when  $|w_i|_B = 0$  for the first time with  $i > 0$ , it is also the case (the rule applied is  $YB \rightarrow Y$ ) and the other rules let the property invariant when there is no more  $B$ .
11. If  $w \neq S$  then every  $b$  is before any  $c$ . We do the same as the previous proof using 5. and 7.

We are now ready to prove that the language is  $\{a^n b^n c^n \mid n \in \mathbb{N}\}$ .

- Assume  $w$  is accepted. Then it contains only terminal elements,  $a, b$ , or  $c$ . Using invariants 8. and 9. we conclude that  $|w|_a = |w|_b = |w|_c$ . Using 10. and 11. we conclude that  $w = a^n b^n c^n$ . To sum it up,  $w \in \{a^n b^n c^n \mid n \in \mathbb{N}\}$ .
- It remains to prove that if  $w \in \{a^n b^n c^n \mid n \in \mathbb{N}\}$ , then  $w$  is accepted. We just give a sequences of rules to build it :

$$\begin{array}{rcl}
S & \rightarrow & ABC \\
& \rightarrow & \dots \\
& \rightarrow & Y(aX)^n BC \\
& \xrightarrow{*} & Y a^n B (bX)^n C \quad \text{using } Xa \rightarrow aX \text{ } n \text{ times} \\
& \xrightarrow{*} & Y a^n B b^n C^n \quad \text{using } Xb \rightarrow bX \text{ } n \text{ times} \\
& \xrightarrow{*} & a^n Y b^n C^n \quad \text{using } Ya \rightarrow aY \text{ } n \text{ times} \\
& \rightarrow & a^n b^n C^n \quad \text{using } Yb \rightarrow bY \text{ } n \text{ times}
\end{array}$$

Then,  $w$  is accepted.

**Exercise 3:** Given a word  $w$  on  $\{a, b\}$  the grammar  $L_w$  is defined as follows:

$$\begin{aligned}
S &\rightarrow AwB & A &\rightarrow AX & Xba &\rightarrow ab & Xab &\rightarrow aXb & Xaa &\rightarrow aXa & Xbb &\rightarrow bXb & XbB &\rightarrow bXB \\
XB &\rightarrow Y & XaB &\rightarrow aY & aY &\rightarrow Ya & bY &\rightarrow Yb & XY &\rightarrow Y & AY &\rightarrow \varepsilon
\end{aligned}$$

1. Compute  $L_w$  for each word  $w$  in the set  $\{\varepsilon, ab, ba, aba\}$ , then guess what  $L_w$  is (compute  $L_w$  for more instances of  $w$  if needed, the pattern should appear clearly).
2. What if we remove the production rule  $XY \rightarrow Y$  from the grammar  $L_w$ ?

**Solution:**

1. We compute  $L_w$  for each word  $w$  in the set  $\{\varepsilon, ab, ba, aba\}$ :

- $w = \varepsilon$ :  $S \rightarrow AB \rightarrow AXB \rightarrow AY \rightarrow \varepsilon$ .  
The language  $L_\varepsilon = \emptyset$ .
- $w = ab$ :  $S \rightarrow AabB \rightarrow AXabB \rightarrow AaXbB \rightarrow AabXB \rightarrow AabY \rightarrow AaYb \rightarrow AYab \rightarrow ab$ .  
The language  $L_{ab} = \{ab\}$ .
- $w = ba$ :  $S \rightarrow AbaB \rightarrow AXbaB \rightarrow AabB \xrightarrow{*} ab$ .  
Note that we already saw the derivation from  $Aab$  in the case  $w = ab$ . The language  $L_{ba} = \{ab\}$ .
- $w = aba$ :  $S \rightarrow AabaB \rightarrow AXabaB \rightarrow AaXbaB \rightarrow AaabB \rightarrow AXaabB \rightarrow AaXabB \rightarrow AaaXbB \rightarrow AaabXB \rightarrow Aaaby \rightarrow AaaYb \rightarrow AaYab \rightarrow AYaab \rightarrow aab$ .  
The language  $L_{aba} = \{aab\}$ .

The language  $L_w$  is the language that contains the word  $a^n b^m$  where  $n$  is the number of  $a$  in  $w$  and  $m$  is the number of  $b$  in  $w$ .

2. If we remove  $XY \rightarrow Y$  any sequence  $XY$  cannot be removed anymore and one obtain an infinite derivation that cannot terminate:

$$S \rightarrow AB \rightarrow AXB \rightarrow AY \rightarrow AXY \rightarrow AXXY \rightarrow \dots$$

Instead, with the rule  $XY \rightarrow Y$ , the infinite derivation shown above still exists, but at any point in time we can decide to apply  $XY \rightarrow Y$  and obtain a finite derivation.

The language  $L_w$  does not change removing the rule  $XY \rightarrow Y$  from the grammar. Thus, even if the two grammars have different derivations they still accept the same language.

**Exercise 4:** In a grammar  $(\Sigma, \mathcal{N}, S, \mathcal{P})$ , a nonterminal symbol  $A$  is said to be *persistent* when for all production rules  $w \rightarrow w'$ , if  $A$  appears in  $w$  then  $A$  also appears in  $w'$ . A production rule is said to be *persistent* when it contains a persistent (nonterminal) symbol. What is the language recognized by the grammar if the start symbol  $S$  is persistent? Assuming that  $S$  is not persistent, prove that persistent symbols and persistent production rules can be removed from a grammar without changing its accepted language, i.e.

$$\llbracket \Sigma, \mathcal{N}, S, \mathcal{P} \rrbracket = \llbracket \Sigma, \mathcal{N} - \{\text{persistent}\}, S, \mathcal{P} - \{\text{persistent}\} \rrbracket.$$

**Solution:** If the initial symbol  $S$  is persistent then the language  $\llbracket G \rrbracket$  of  $G$  is empty since from  $S$  one can only apply persistent rules, hence always obtaining an  $S$  in every new derivation. The language is empty because there are no derivations starting from  $S$  that terminates with a word  $w \in \Sigma^*$ .

Let  $G = (\Sigma, \mathcal{N}, S, \mathcal{P})$  be a grammar,  $\mathcal{N}_P$  be the set of persistent non-terminal symbols of  $G$  (not containing  $S$ ),  $\mathcal{P}_P$  be the set of persistent production  $G$ , and  $G_P = (\Sigma, \mathcal{N} - \mathcal{N}_P, S, \mathcal{P} - \mathcal{P}_P)$  the grammar obtained removing the persistent non-terminals and productions from  $G$ .

We prove that  $\llbracket G \rrbracket = \llbracket G_P \rrbracket$ .

- $\llbracket G_P \rrbracket \subseteq \llbracket G \rrbracket$ .

If  $w \in \llbracket G_P \rrbracket$ , then there exists a derivation in  $G_P$  such that  $S \xrightarrow[G_P]{*} w$  (the notation  $\xrightarrow[G_P]{*}$  describes 0 or more derivations using any rule from the grammar  $G_P$ ). We show by induction on the length of the derivation that for any derivation  $S \xrightarrow[G_P]{*} w$  there is a derivation  $S \xrightarrow[G]{*} w$ .

- **Base case:** Suppose  $S \xrightarrow[G_P]{1} w$ . If  $S \xrightarrow[G_P]{1} w$ , then there exists a rule  $(S, w) \in \mathcal{P}_P$ . Since  $\mathcal{P}_P \subseteq \mathcal{P}$ , then  $(S, w) \in \mathcal{P}$  and hence  $S \xrightarrow[G]{1} w$ .
- **Inductive case:**  
Suppose  $S \xrightarrow[G_P]{n+1} w$ . By inductive hypothesis we have that  $S \xrightarrow[G_P]{n} v \xrightarrow[G_P]{1} w$  and  $S \xrightarrow[G]{n} v$ . Since  $v \xrightarrow[G_P]{1} w$ , there exists  $(\gamma, \delta) \in \mathcal{P}_P$  such that  $v = \alpha\gamma\beta$ , for some  $\alpha, \gamma \in (\Sigma \cup \mathcal{N})^*$ , and  $w = \alpha\delta\beta$  (this directly from the definition of a derivation for a grammar). As before,  $(\gamma, \delta) \in \mathcal{P}$  and then we can apply this rule to derive  $v \xrightarrow[G]{1} w$ . Thus we have that  $S \xrightarrow[G]{n} v \xrightarrow[G]{1} w$ .

Thus, if  $w \in \llbracket G_P \rrbracket$  then there exists a derivation  $S \xrightarrow[G_P]{*} w$  and hence a derivation  $S \xrightarrow[G]{*} w$ , and hence  $w \in \llbracket G \rrbracket$ .

- $\llbracket G \rrbracket \subseteq \llbracket G_P \rrbracket$

We first show that when we start from a word that includes a persistent non-terminal symbol  $A$ , all the possible derivations with the grammar  $G$  will always contain a  $A$ . In practice, we show that there are no derivations that can “remove” a non-terminal persistent symbol  $A$  from a word.

If  $A$  is a persistent non-terminal (i.e.,  $A \in \mathcal{N}_P$ ), then any word  $w = vAv' \in (\Sigma \cup \mathcal{N})^*$ , for some  $v, v' \in (\Sigma \cup \mathcal{N})^*$ , is such that  $w \xrightarrow[G]{*} w'$ , then  $w' = uAu'$  for some  $u, u' \in (\Sigma \cup \mathcal{N})^*$ .

We prove the lemma by induction on the length of the derivation.

- **Base case:** Suppose  $w \xrightarrow[G]{1} w'$  and  $w = vAv' \in (\Sigma \cup \mathcal{N})^*$ , for some  $v, v' \in (\Sigma \cup \mathcal{N})^*$ .  
If  $w \xrightarrow[G]{1} w'$ , then there exists a rule in  $(\delta, \gamma) \in \mathcal{P}$  that derives  $w'$  from  $vAv'$ . There are two possible cases:

1.  $(\delta, \gamma)$  is not persistent (i.e.,  $(\delta, \gamma) \in \mathcal{P} - \mathcal{P}_P$ ). In this case,  $A$  cannot appear in  $\delta$ , and hence the non-terminal  $A$  is not removed from  $w$ .
2.  $(\delta, \gamma)$  is persistent. In this case  $\delta = w_\delta A w'_\delta$  and  $\gamma = w_\gamma A w'_\gamma$  for some  $w_\delta, w'_\delta, w_\gamma, w'_\gamma \in (\Sigma \cup \mathcal{N})^*$ . Thus, the derivation with  $(\delta, \gamma)$  will contain at least one occurrence of the non-terminal  $A$  (i.e.,  $v A v' \xrightarrow[G]{1} u A u'$ ).

– **Inductive case:** Suppose  $w \xrightarrow[G]{n} w'' \xrightarrow[G]{1} w'$ . By inductive hypothesis,  $w'' = z A z', z, z' \in (\Sigma \cup \mathcal{N})^*$ .

We can prove that  $w'' = z A z' \xrightarrow[G]{1} w'' = u A u'$ , for some  $u, u' \in (\Sigma \cup \mathcal{N})^*$  as above, showing that if  $A$  is already in  $w''$  then  $A$  is also in  $w'$  after a single derivation step.

Now we prove  $\llbracket G \rrbracket \subseteq \llbracket G_P \rrbracket$  by contradiction. Suppose  $w \notin G_P$ . If  $w \notin G_P$ , then  $S \xrightarrow[G]{n} w$  is such that there is a derivation at a step  $0 \leq i \leq n$  that uses a persistent rule  $(w_\delta A w'_\delta, w_\gamma A w'_\gamma)$ , for some persistent non-terminal  $A$ .

Thus, we have  $S \xrightarrow[G]{i-1} v \xrightarrow[G]{i} v' \xrightarrow[G]{n-i} w$ . with  $v = \alpha A \beta$  and  $v' = \alpha' A \beta'$  for some  $\alpha, \beta, \alpha', \beta' \in (\Sigma \cup \mathcal{N})^*$ .

Since  $A$  is in the word  $v$ , then  $A$  is also in the word  $w$  (by the lemma we proved above. But this is a contradiction, since  $w \in \llbracket G \rrbracket$  and hence cannot contain non-terminals. We conclude that  $\llbracket G \rrbracket \subseteq \llbracket G_P \rrbracket$ .

**Exercise 5:** Prove that if we allow infinite sets of production rules and infinite sets of nonterminal symbols in the definition of a grammar, then for any language  $L$  there is a regular grammar such that  $L = \llbracket G \rrbracket$ .

**Solution:**

We define a regular grammar (with an infinite set of production rules and an infinite set of nonterminal symbols) that can accept an arbitrary language  $L$ .

The grammar is  $G = (\Sigma_L, \mathcal{N}, S, \mathcal{P})$  where:

- $\Sigma_L$  is the alphabet of  $L$ ;
- $\mathcal{N} = \{S\} \cup \{A_{w_i} \mid w \in L, 1 \leq i < |w|\}$ .

We define a non-terminal symbol for every letter  $w(i)$  in every (non-empty) word  $w$  of  $L$ . We will use a non-terminal  $A_{w_i}$  to keep track that the grammar generated the prefix  $w(1)w(2)\dots w(i)$  of  $w$ . We will further add a rule that generates the next character  $w(i+1)$  when matching the non-terminal  $A_{w_i}$ . We further have an initial symbol  $S$ . Note that  $|w|$  denotes the length of the word  $w$ .

- We define a set of production rules that can generate any word in  $L$ :

$$\begin{aligned}
 \mathcal{P} = & \bigcup_{\varepsilon \in L} \{S \rightarrow \varepsilon\} \cup && \text{Accept the empty word} \\
 & \bigcup_{w \in L \setminus \{\varepsilon\}} \{S \rightarrow w(1)A_{w_1}\} \cup && \text{Derives the first character of } w \\
 & \bigcup_{w \in L \setminus \{\varepsilon\}} \bigcup_{1 \leq i < |w|} \{A_{w_i} \rightarrow w(i+1)A_{w_{i+1}}\} \cup && \text{Derives the prefix of length } i \text{ of } w \\
 & \bigcup_{w \in L \setminus \{\varepsilon\}} \{A_{w_{|w|}} \rightarrow \varepsilon\} && \text{Accepts the word } w
 \end{aligned}$$

In practice, the grammar can either generate the empty word  $\varepsilon$  (if it is part of  $L$ ), or non-deterministically choose to generate the first character of a word  $w \in L$ , picking one rule  $S \rightarrow w(1)A_{w_1}$ . At this point the grammar can only pick the rule that contains  $A_{w_1}$  on the left-hand side, hence starting to derive the word  $w$ , that is finally accepted with the rule  $A_{w_{|w|}}$ .

We prove that  $\llbracket G \rrbracket = L$ .

- $L \subseteq \llbracket G \rrbracket$  We prove that if  $w \in L$  then  $w \in \llbracket G \rrbracket$ .

If  $|w| = 0$ , then  $S \xrightarrow[G]{i} \varepsilon$ , and hence  $\varepsilon \in \llbracket G \rrbracket$ .

We then show that for any prefix of  $w$ ,  $|w| > 0$ , of length  $i \leq |w|$  there exists a derivation  $S \xrightarrow[G]{i} w(1)w(2)\dots w(i)A_{w_i}$ . We prove it by induction on the prefix length.

- **Base case.** For  $i = 1$  we have that  $S \xrightarrow[G]{1} w(1)A_{w_1}$ , since  $S \rightarrow w(1)A_{w_1} \in \mathcal{P}$ .
- **Inductive case.** For  $i = i + 1 \leq |w|$ , we have that  $S \xrightarrow[G]{i} w(1)\dots w(i)A_{w_i}$  by inductive hypothesis. By the definition of  $G$ , there exists a rule  $A_{w_i} \rightarrow w(i+1)A_{w_{i+1}} \in \mathcal{P}$ . Using the inductive hypothesis and the above rule we conclude that  $S \xrightarrow[G]{i} w(1)\dots w(i)A_{w_i} \xrightarrow[G]{1} w(1)\dots w(i+1)A_{w_{i+1}}$ .

Thus, there exist a derivation  $S \xrightarrow[G]{|w|} w(1)\dots w(|w|)A_{w_{|w|}}$ . We conclude that  $w \in \llbracket G \rrbracket$  since  $A_{w_{|w|}} \rightarrow \varepsilon \in \mathcal{P}$ .

- $\llbracket G \rrbracket \subseteq L$

Let  $w \in \llbracket G \rrbracket$ , and  $S \xrightarrow[G]{*} w$  its derivation using the grammar  $G$ .

We show that for any derivation  $S \xrightarrow[G]{i} w'$  either  $w'$  is part of the language  $L$  or  $w' = w''(1)\dots w''(i)A_{w_i}$ , for a  $w'' \in L$ .

We prove it by induction on the length of the derivation.

- **Base case.** Consider the case  $i = 1$  and the derivation  $S \xrightarrow[G]{1} w'$

After one derivation there are two sub-cases to consider: either a)  $S \xrightarrow[G]{1} \varepsilon$  or b)  $S \xrightarrow[G]{1} w(1)A_{w_1}$  for some word  $w$  from the definition of  $\mathcal{P}$  (i.e.,  $\mathcal{P}$  only has rules with  $S$  on the left-hand side for the rules  $S \rightarrow \varepsilon$  and  $S \rightarrow w(1)A_{w_1}$ ).

In case a),  $S \xrightarrow[G]{1} \varepsilon$  because of the production  $S \rightarrow \varepsilon$ . By construction we have that  $(S \rightarrow \varepsilon) \in \mathcal{P}$  iff  $\varepsilon \in L$ . Thus  $w' = \varepsilon$  is a word in  $L$ .

In case b),  $S \xrightarrow[G]{1} w(1)A_{w_1}$  because of the production  $S \rightarrow w(1)A_{w_1}$ . By construction,  $w(1)$  is a prefix of a word  $w$  in  $L$ .

- **Inductive case.** Consider the case  $i = i + 1$  and the derivation  $S \xrightarrow[G]{i+1} w'$ .

By inductive hypothesis we know that  $S \xrightarrow[G]{i} v$  is such that either  $v$  is a word in  $L$  or  $v = w(1)\dots w(i)A_{w_i}$  for a word  $w \in L$ . The first case is impossible, since  $v$  would only contain

non-terminal symbols ( $v \in \Sigma^*$ ), and the grammar is regular (i.e., all the left-hand sides of the rules are non-terminals). Thus, we know that  $S \xrightarrow{G} w(i)A_{w_i}$ .

We know that the only rules in the  $\mathcal{P}$  containing the non-terminal  $A_{w_i}$  are either  $A_{w_i} \rightarrow w(i+1)A_{w_{i+1}}$ , if  $i \leq |w|$ , or  $A_{|w|} \rightarrow \varepsilon$ , if  $i > |w|$ .

In the first case we have that  $w' = w(i) \dots w(i+1)A_{w_{i+1}}$  for a word  $w \in L$ , since  $S \xrightarrow{G} w(i) \dots w(i)A_{w_i} \xrightarrow{G} w(i) \dots w(i+1)A_{w_{i+1}}$ .

In the second case, we have that  $w' \in L$  since  $S \xrightarrow{G} w(1) \dots w(|w|)A_{|w|} \xrightarrow{G} w(1) \dots w(|w|)$ .

We conclude that if  $w \in \llbracket G \rrbracket$  then  $w \in L$ .

**Exercise 6:** Prove that if we allow *infinite sets of production rules* in the definition of a grammar, then for any language  $L$  there is a context-free grammar such that  $L = \llbracket G \rrbracket$ .

**Solution:**

We define a regular grammar (with an infinite set of production rules) that can accept an arbitrary language  $L$ . The grammar is  $G = (\Sigma_L, \mathcal{N}, S, \mathcal{P})$  where:

- $\Sigma_L$  is the alphabet of  $L$ ;
- $\mathcal{N} = \{S\}$ .
- $\mathcal{P} = \bigcup_{w \in L} \{S \rightarrow w\}$

The grammar just define a production rule for every word  $w \in L$ . Every rule  $S \rightarrow w$  can derive a word  $w \in L$  from the initial symbol  $S$ .

If  $w \in L$  then  $w \in \llbracket G \rrbracket$  because there exists a derivation  $S \xrightarrow{G} w$  with the rule  $S \rightarrow w \in \mathcal{P}$ .

If  $w \in \llbracket G \rrbracket$ , then there exists a derivation  $S \xrightarrow{*G} w$ . Since all the rules in  $\mathcal{P}$  have only  $S$  in the left-hand side, then  $S \xrightarrow{*G} w$  is  $S \xrightarrow{1G} w$ , for some rule  $S \rightarrow w \in \mathcal{P}$ . If  $S \rightarrow w \in \mathcal{P}$ , then  $w \in L$  by the definition of  $G$ .

**Exercise 7:** Given two grammars  $G_1 = (\Sigma_1, \mathcal{N}_1, S_1, \mathcal{P}_1)$  and  $G_2 = (\Sigma_2, \mathcal{N}_2, S_2, \mathcal{P}_2)$  find two grammars  $G'_1 = (\Sigma'_1, \mathcal{N}'_1, S'_1, \mathcal{P}'_1)$  and  $G'_2 = (\Sigma'_2, \mathcal{N}'_2, S'_2, \mathcal{P}'_2)$  such that:

1.  $\Sigma'_1 = \Sigma'_2$
2.  $\mathcal{N}'_1 \cap \mathcal{N}'_2 = \emptyset$
3.  $\llbracket G_1 \rrbracket = \llbracket G'_1 \rrbracket$  and  $\llbracket G_2 \rrbracket = \llbracket G'_2 \rrbracket$

Deduce that:

1. the union of two formal / context-free languages is formal / context-free.
2. the concatenation of two formal / context-free languages is formal / context-free.
3. the Kleene star of a formal / context-free language is formal / context-free.



**Solution:**

We define a function that renames the non-terminal symbols of  $\mathcal{N}_2$  ensuring that they are disjoint from  $\mathcal{N}_1$ :

$$\phi(A) = \begin{cases} A & \text{when } A \in \mathcal{N}_2 \setminus \mathcal{N}_1 \\ A' & \text{otherwise} \end{cases}$$

We assume that each  $A'$  is a new symbol (i.e., not contained in any of  $\Sigma_1, \Sigma_2, \mathcal{N}_1, \mathcal{N}_2$ ), and that  $\phi(A) = \phi(B)$  iff  $A = B$ .

We extend  $\phi$  to words in  $(\Sigma_2 \cup \mathcal{N}_2)^*$  where we rename only the symbols in a word that are from  $\mathcal{N}_2$ :

$$\phi(w) = \begin{cases} \varepsilon & \text{if } w = \varepsilon \\ \phi(a)\phi(w') & \text{if } w = aw' \text{ and } a \in \Sigma_2 \\ \phi(A)\phi(w') & \text{if } w = Aw' \text{ and } A \in \mathcal{N}_2 \end{cases}$$

Then, we define:

- $G'_1 = (\Sigma'_1 = \Sigma_1 \cup \Sigma_2, \mathcal{N}'_1 = \mathcal{N}_1, S'_1 = S_1, \mathcal{P}'_1 = \mathcal{P}_1)$
- $G'_2 = (\Sigma'_2 = \Sigma_1 \cup \Sigma_2, \mathcal{N}'_2 = \{\phi(A) \mid A \in \mathcal{N}_2\}, S'_2 = \phi(S_2), \mathcal{P}'_2 = \{\phi(w) \rightarrow \phi(w') \mid w \rightarrow w' \in \mathcal{P}_2\})$

Here we are also assuming that  $\mathcal{N}_1 \cap \Sigma_2 = \emptyset$  and  $\mathcal{N}_2 \cap \Sigma_1 = \emptyset$  (otherwise, we should further rename  $\mathcal{N}_1$  to force it to be disjoint from  $\Sigma_2$  and consider  $\Sigma_1$  when renaming  $\mathcal{N}_2$ ).

We can easily prove that  $\llbracket G_1 \rrbracket = \llbracket G'_1 \rrbracket$  (in fact, the two grammars have the same derivations). We can prove that  $\llbracket G_2 \rrbracket = \llbracket G'_2 \rrbracket$  showing that they have the same derivations but with the non-terminals renamed with  $\phi$  in  $G'_2$  (we also need an inverse of  $\phi$  when showing the direction  $\supseteq$ ).

**Closure properties of Context-free languages**

We first show that context-free languages are closed under union, concatenation, and Kleene closure.

If two languages  $L_1, L_2$  are context-free then there exists two context-free grammars  $G_1 = (\Sigma_1, \mathcal{N}_1, S_1, \mathcal{P}_1)$  and  $G_2 = (\Sigma_2, \mathcal{N}_2, S_2, \mathcal{P}_2)$  such that  $L_1 = \llbracket G_1 \rrbracket$  and  $L_2 = \llbracket G_2 \rrbracket$ . We assume that the two grammars  $G_1, G_2$  have two disjoint sets of non-terminal symbols (i.e.,  $\mathcal{N}_1 \cap \mathcal{N}_2 = \emptyset$ ). This is not a general restriction since the set of non-terminals can be renamed while still accepting the same language (with the construction shown above).

**• Union.**

We define the union of the two grammars as  $G_1 \cup G_2 = (\Sigma_{G_1 \cup G_2}, \mathcal{N}_{G_1 \cup G_2}, S, \mathcal{P}_{G_1 \cup G_2})$  where:

- $\Sigma_{G_1 \cup G_2} = \Sigma_1 \cup \Sigma_2$
- $\mathcal{N}_{G_1 \cup G_2} = \mathcal{N}_1 \cup \mathcal{N}_2 \cup \{S\}$
- $\mathcal{P}_{G_1 \cup G_2} = \mathcal{P}_1 \cup \mathcal{P}_2 \cup \{S \rightarrow S_1, S \rightarrow S_2\}$ .

The union of the grammars keeps the same productions of  $G_1$  and  $G_2$  and just adds two productions  $S \rightarrow S_1$  and  $S \rightarrow S_2$ . In this way, a derivation in  $G_1 \cup G_2$  will first non-deterministically “select” to follow either the derivations of  $G_1$  (and hence deriving a word in  $L_1$ ) or the derivations of  $G_2$  (and hence deriving a word in  $L_2$ ).

We first show that if  $L_1$  and  $L_2$  are context free, then  $G_1 \cup G_2$  is context free. Then we will show that the formal grammar  $G_1 \cup G_2$  recognizes the same language of  $L_1 \cup L_2$ .

We show that if  $G_1$  and  $G_2$  are context free grammars (they exists since  $L_1$  and  $L_2$  are context-free), then also  $G_1 \cup G_2$  is a context-free grammar. It is easy to see that all the production rules in  $\mathcal{P}_{G_1 \cup G_2}$  are of the form  $A \rightarrow w$ , where  $A \in \mathcal{N}_{G_1 \cup G_2}$  and  $w \in (\Sigma_{G_1 \cup G_2})^*$ . A rule in  $w \rightarrow w' \in \mathcal{P}_{G_1 \cup G_2}$  is either:

- A rule  $w \rightarrow w' \in \mathcal{P}_1$ . Since  $G_1$  is context-free, then  $w \rightarrow w'$  is such that  $w \in \mathcal{N}_1$  and  $w' \in (\Sigma_1 \cup \mathcal{N}_1)^*$ .
- A rule  $A \rightarrow w \in \mathcal{P}_2$ . This is a context-free rule for  $\mathcal{P}_{G_1 \cup G_2}$ , similarly to the rules from  $\mathcal{P}_1$ .
- $S \rightarrow S_1$ : this rule is context-free.
- $S \rightarrow S_2$ : this rule is context-free.

Now we show that  $\llbracket G_1 \cup G_2 \rrbracket = L_1 \cup L_2$ .

- $\llbracket G_1 \cup G_2 \rrbracket \supseteq L_1 \cup L_2$ .  
 If  $w \in L_1 \cup L_2$  then either  $w \in L_1$  or  $w \in L_2$ .  
 If  $w \in L_1$  then  $w \in \llbracket G_1 \rrbracket$ , and then there exists a derivation  $S_1 \xrightarrow{*}_{G_1} w$ . We have that  $S \xrightarrow{1}_{G_1 \cup G_2} S_1$  (since  $S \rightarrow S_1 \in \mathcal{P}_{G_1 \cup G_2}$ ) and that  $\mathcal{P}_1 \subseteq \mathcal{P}_{G_1 \cup G_2}$ . Thus, we have a derivation  $S \xrightarrow{1}_{G_1 \cup G_2} S_1 \xrightarrow{*}_{G_1 \cup G_2} w$  showing that  $w \in \llbracket G_1 \cup G_2 \rrbracket$ .  
 If  $w \in L_2$ , we have a symmetric proof showing that  $w \in \llbracket G_1 \cup G_2 \rrbracket$ .
- $\llbracket G_1 \cup G_2 \rrbracket \subseteq L_1 \cup L_2$   
 If  $w \in \llbracket G_1 \cup G_2 \rrbracket$  then there exists a derivation  $S \xrightarrow{*}_{G_1 \cup G_2} w$ .  
 Since the only derivation containing  $S$  on the left-hand side are  $S \rightarrow S_1$  and  $S \rightarrow S_2$ , then we have either that a)  $S \xrightarrow{1}_{G_1 \cup G_2} S_1 \xrightarrow{*}_{G_1 \cup G_2} w$  or b)  $S \xrightarrow{1}_{G_1 \cup G_2} S_2 \xrightarrow{*}_{G_1 \cup G_2} w$ .  
 In the case a), since  $\mathcal{N}_1$  and  $\mathcal{N}_2$  are disjoint and the grammar is context-free (i.e., every left-hand side of a rule is a single non-terminal) then the only possible derivations from the word  $S_1$  uses productions rules from  $G_1$ . Thus, there exists a derivation  $S_1 \xrightarrow{*}_{G_1} w$ , and hence  $w \in \llbracket G_1 \rrbracket$ ,  $w \in L_1$ , and  $w \in L_1 \cup L_2$ .  
 The case b) is symmetric to the case a).

Since  $G_1 \cup G_2$  is a context-free grammar and  $\llbracket G_1 \cup G_2 \rrbracket = L_1 \cup L_2$ , then the union  $L_1 \cup L_2$  is a context-free language.

#### • Concatenation.

We define the concatenation of two grammar as:  $G_1 G_2 = (\Sigma_{G_1 G_2} = \Sigma_1 \cup \Sigma_2, \mathcal{N}_{G_1 G_2} = \mathcal{N}_1 \cup \mathcal{N}_2 \cup \{S\}, S, \mathcal{P}_{G_1 G_2} = \mathcal{P}_1 \cup \mathcal{P}_2 \cup \{S \rightarrow S_1 S_2\})$ .

The production  $S \rightarrow S_1 S_2$  allows the grammar to accept the concatenation of a word derived following only productions in  $G_1$  (and hence in  $L_1$ ) to a word derived following only productions in  $G_2$  (and hence in  $L_2$ ).

All the production rules in  $\mathcal{P}_{G_1 G_2}$  are of the form  $A \rightarrow w$ , where  $A \in \mathcal{N}_{G_1 G_2}$  and  $w \in (\Sigma_{G_1 G_2} \cup \mathcal{N}_{G_1 G_2})^*$ . The only “new” rule is  $S \rightarrow S_1 S_2$ , which is context-free. Thus the grammar is context-free.

Now we show that  $\llbracket G_1 G_2 \rrbracket = L_1 L_2$ .

- $\llbracket G_1 G_2 \rrbracket \supseteq L_1 L_2$ .  
 If  $w \in L_1 \cap L_2$  then there exists two words  $w_1 \in L_1$  and  $w_2 \in L_2$ .  
 Then there exist two derivations  $S_1 \xrightarrow{n}_{G_1} w_1$  and  $S_2 \xrightarrow{m}_{G_2} w_2$ .  
 We can construct a derivation in  $G_1 G_2$  from  $S$  to  $w = w_1 w_2$  as:  $S \xrightarrow{1}_{G_1 G_2} S_1 S_2 \xrightarrow{n}_{G_1} w_1 S_2 \xrightarrow{m}_{G_2} w_1 w_2$

–  $\llbracket G_1 G_2 \rrbracket \subseteq L_1 L_2$ .

If  $w \in \llbracket G_1 G_2 \rrbracket$  then there exists a derivation  $S \xrightarrow[G_1 G_2]{*} w$ .

We show that a derivation in  $G_1 G_2$  always derives a word  $w = w_1 w_2$  such that  $w_1$  can be derived from  $G_1$  and  $w_2$  can be derived from  $G_2$  (i.e.,  $S \xrightarrow[G_1]{*} w_1$  and  $S \xrightarrow[G_2]{*} w_2$ ). This is sufficient to show that a derivation in  $G_1 G_2$  deriving a word  $w \in (\Sigma_1 \cup \Sigma_2)^*$  is the concatenation of two words derived from  $G_1$  and  $G_2$ .

\* **Base case.**

In the first step the only possible derivation is  $S \xrightarrow[G_1 G_2]{1} S_1 S_2$ .  $S_1$  and  $S_2$  are two possible derivations of  $G_1$  and  $G_2$  (with 0 steps).

\* **Inductive case.** Suppose we have  $S \xrightarrow[G_1 G_2]{n} w'_1 w'_2 \xrightarrow[G_1 G_2]{1} w$ .

By inductive hypothesis we have that  $S \xrightarrow[G_1]{n} w'_1$  and  $S \xrightarrow[G_2]{n} w'_2$ , and that  $w'_1 \in (\Sigma_1 \cup \mathcal{N}_1)^*$  and  $w'_2 \in (\Sigma_2 \cup \mathcal{N}_2)^*$ .

Thus, there is a rule  $(\gamma, \delta) \in \mathcal{P}_{G_1 G_2}$  such that  $w'_1 w'_2 \xrightarrow[G_1 G_2]{1} w$ . We know that both  $G_1$  and  $G_2$  are context-free, and their non-terminal symbols are disjoint. Thus, it is the case that the derivation is done with a rule in  $\mathcal{P}_{G_1 G_2}$  that is either from  $\mathcal{P}_1$  or  $\mathcal{P}_2$ . In practice, we have either  $w'_1 w'_2 \xrightarrow[G_1 G_2]{1} w_1 w'_2$  or  $w'_1 w'_2 \xrightarrow[G_1 G_2]{1} w'_1 w_2$ .

To see this, notice that  $w'_1 \in (\Sigma_1 \cup \mathcal{N}_1)^*$  and a rule from  $\mathcal{P}_{G_1 G_2}$  is of the form  $(A, \delta)$ , with either a)  $A = S$ , b)  $A \in \mathcal{N}_1$ , or c)  $A \in \mathcal{N}_2$ . We cannot be in case a), since the derivation is at step  $n > 1$ . We cannot be in case c) either, since  $w'_1 \in (\Sigma_1 \cup \mathcal{N}_1)^*$  and  $(\Sigma_1 \cup \mathcal{N}_1) \cap \mathcal{N}_2 = \emptyset$ . Thus, the only possible case to apply a rule matching a substring in  $w_1$  is b). The same reasoning works for  $w'_2$ .

Thus, if  $S \xrightarrow[G_1 G_2]{*} w = w_1 w_2$ , then  $S \xrightarrow[G_1]{*} w_1$  and  $S \xrightarrow[G_2]{*} w_2$ . Thus, if  $w$  is a word on  $\llbracket G_1 G_2 \rrbracket$ , then  $w = w_1 w_2$  and  $w_1 \in \llbracket G_1 \rrbracket = L_1$  and  $w_2 \in \llbracket G_2 \rrbracket = L_2$ , and hence  $w \in L_1 L_2$ .

• **Kleene closure.** Let  $G_1 = (\Sigma_1, \mathcal{N}_1, S_1, \mathcal{P}_1)$  be the grammar that accepts the formal language  $L_1$ .

We define the Kleene closure of  $G_1$  as:  $G_1^* = (\Sigma_{G_1^*} = \Sigma_1, \mathcal{N}_{G_1^*} = \mathcal{N}_1 \cup \{S\}, S, \mathcal{P}_{G_1^*} = \mathcal{P}_1 \cup \{S \rightarrow \varepsilon, S \rightarrow S_1 S\})$ .

The rule  $S \rightarrow \varepsilon$  allows the grammar to recognize  $\varepsilon \in L_1^0$  and terminate a derivation, while  $S \rightarrow S_1 S$  allows the grammar to derive a word  $w \in L_1$  concatenated to a word that can be derived again from  $S$ .

$G_1^*$  is a context-free grammar, since  $G_1$  is a context-free grammar and the additional rules to  $\mathcal{P}_1$  are  $S \rightarrow \varepsilon$  and  $S \rightarrow S_1 S$ .

Now we show that  $\llbracket G_1^* \rrbracket = L_1^*$ .

–  $\llbracket G_1^* \rrbracket \supseteq L_1^*$ .

We show that if  $w \in L_1^*$  then there exists a derivation  $S \xrightarrow[G_1^*]{*} w$ , and thus we conclude that  $w \in \llbracket G_1^* \rrbracket$ .

If  $w \in L_1^*$ , then we have that  $w \in \bigcup_{i \geq 0} L_1^i$ , from the definition of the Kleene closure.

We have that if  $w = \varepsilon$  (i.e.,  $w \in L_1^0$ ), then  $S \xrightarrow[G_1^*]{1} \varepsilon$  (because  $S \rightarrow \varepsilon \in \mathcal{P}_{G_1^*}$ ).

We then show that for any  $w' \in L_1^n$ ,  $n > 1$ , the grammar  $G_1^*$  can derive the word  $w'S$ . That is, the grammar can derive a word  $w' \in L_1^n$  and always concatenate a new word (including a word from  $w'' \in L_1$ ) with derivation from the non-terminal  $S$ . This allows us to conclude the first language inclusion. We prove this by induction on the number of words from  $L_1$  concatenated in  $w$ .

\* **Base case.** Suppose  $w \in L_1$ .

Then there exists a derivation  $S_1 \xrightarrow[G_1]{*} w$  of  $G_1$  that we can use to show that there is the following derivation in  $G_1^*$ , since all the productions of  $G_1$  are in  $G_1^*$ :  $S \xrightarrow[G_1^*]{1} S_1 S \xrightarrow[G_1^*]{*} w S_1$ .

\* **Inductive case.** Suppose  $w \in L_1^n$ . Thus, we know that  $w = w''w'$ , for some  $w'' \in L_1^{n-1}$  and  $w' \in L_1$ . So, we also have a derivation in  $G_1$  such that  $S_1 \xrightarrow[G_1]{*} w'$

By inductive hypothesis, we have that there exists a derivation  $S \xrightarrow[G_1^*]{*} w''S$  and hence we have the following derivation:

$$S \xrightarrow[G_1^*]{*} w''S \xrightarrow[G_1^*]{1} w''S_1S \xrightarrow[G_1^*]{*} w''w'S$$

Applying the above result, we have that if  $w \in L_1^*$  then there exists a derivation  $S \xrightarrow[G_1^*]{*} wS \xrightarrow[G_1^*]{1} w$ , and hence we conclude that  $w \in \llbracket G_1^* \rrbracket$ .

–  $\llbracket G_1^* \rrbracket \subseteq L_1$ .

Given a word  $w \in \llbracket G_1^* \rrbracket$  we want to show that  $w \in L_1$ . As usual, we know that there exists a derivation  $S \xrightarrow[G_1^*]{*} w$ .

Showing that the derivation  $S \xrightarrow[G_1^*]{*} w$  produces a word in  $L_1^*$  is more involved, because the grammar can “start” several derivations from several non-terminals  $S_1$ . For example:

$$S \xrightarrow[G_1^*]{1} S_1 S \xrightarrow[G_1^*]{1} S_1 S_1 S \xrightarrow[G_1^*]{1} S_1 S_1 S_1 S \dots$$

We can show that, in any possible derivations of  $G_1^*$ , we obtain a word  $w = u_1 \dots u_n$  that is the concatenation of either: a) a word from  $L_1$  (i.e.,  $u_i \in L_1$ ), b) a word (including terminal and non-terminals) that can be derived from  $G_1$ , c) the non-terminal  $S$ .

Thus, when the derivation accepts a word  $w \in \Sigma_1^*$ , we know that  $w \in L_1^*$ .

We prove this by induction on the number of derivations in  $S \xrightarrow[G_1^*]{*} w$

\* **Base case.** If  $S \xrightarrow[G_1^*]{1} w$  we either have that  $S \xrightarrow[G_1^*]{1} \varepsilon$  or  $S \xrightarrow[G_1^*]{1} S_1 S$ .

\* **Inductive case.** If  $S \xrightarrow[G_1^*]{n} w$ , we know by inductive hypothesis that  $S \xrightarrow[G_1^*]{n-1} u'_1 \dots u'_k$  where a  $u_i$  is either a)  $u'_i \in L_1$ , b)  $u'_i \in (\Sigma_1 \cup \mathcal{N}_1)^*$  and  $S \xrightarrow[G_1^*]{*} u'_i$ , or c)  $u'_i = S$ .

Since the grammar is context-free, then the derivation rule at the step  $n$  is only applied to one of the substrings  $u_i$  if they are either in the cases b) or in c).

In the case b), then one must apply a rule from  $G_1$  (since  $u_i$  only contains non-terminals from  $G_1$ ). Since we have that  $S \xrightarrow[G_1^*]{*} u'_i$  and we apply a rule from  $G_1$ , then we have that

$$S \xrightarrow[G_1^*]{*} u'_i \xrightarrow[G_1^*]{*} u_i, \text{ with } u_i \text{ again in case b) or a) (if the word } u_1 \in \Sigma_1^*).$$

In the case c) the only derivations are with the rules  $S \rightarrow \varepsilon$  or  $S \rightarrow S_1 S$ . In the first case,

we have the derivation:  $S \xrightarrow[n-1]{G_1^*} u'_1 \dots u'_n \xrightarrow{1}{G_1^*} u'_1 \dots u'_{i-1} u'_{i+1} \dots u'_k$ . In the second case, we

have the derivation:  $S \xrightarrow[n-1]{G_1^*} u'_1 \dots u'_n \xrightarrow{1}{G_1^*} u'_1 \dots u'_{i-1} S_1 S u'_{i+1} \dots u'_k$ .

Since  $S \xrightarrow{n}{G_1^*} w$  with  $w \in \Sigma^*$ , then  $w \in L_1$ .

**Formal languages.** Formal languages are closed under union, concatenation, and Kleene closure.

However, the grammars provided before for the union/concatenation/Kleene closure of context-free may not recognize the union/concatenation/Kleene closure if the languages are not context-free.

To see this, consider the grammars with the following productions ( $S_1$  and  $S_2$  are the starting symbols):

$$S_1 \rightarrow b$$

and:

$$S_2 \rightarrow AabA \quad Aa \rightarrow a \quad bA \rightarrow bb$$

We have that  $\llbracket G_1 \rrbracket = \{b\}$ ,  $\llbracket G_2 \rrbracket = \{abb\}$ , and  $\llbracket G_1 \rrbracket \llbracket G_2 \rrbracket = \{babbb\}$ . However, the language obtained with the grammar  $G_1 G_2$  as defined above would be  $\llbracket G_1 \cup G_2 \rrbracket = \{babbb, bbabb\}$ .

The word  $bbabb$  is not part of the concatenation  $\llbracket G_1 \rrbracket \llbracket G_2 \rrbracket$  and is obtained with the derivation:

$$S \rightarrow S_1 S_2 \rightarrow b S_2 \rightarrow b A a b A \rightarrow b b a b A \rightarrow b b a b b$$

The issue here is that a rule from  $G_2$  is applied also on a terminal from  $G_2$ . This example shows that the above definition of concatenation of grammars is not correct to show that unconstrained grammars are closed under concatenation. We can find similar examples for the union and the Kleene closure.

**Exercise 8:** Let  $\Sigma_1$  and  $\Sigma_2$  be two alphabets. A language morphism is a mapping  $h : \Sigma_1^* \rightarrow \Sigma_2^*$  such that for all words  $w, w' \in \Sigma_1^*$  we have  $h(ww') = h(w)h(w')$ .

1. Prove that  $h$  preserves the empty word (i.e.  $h(\varepsilon) = \varepsilon$ ).
2. Prove that  $h$  is entirely determined by its restriction to  $\Sigma_1$  (i.e. the words of length 1).
3. Prove that any map  $f : \Sigma_1 \rightarrow \Sigma_2^*$  extends to a language morphism in a unique way.

Let  $\mathcal{N}$  be a set of nonterminal symbols and  $\tilde{h}$  be the unique extension of  $h$  to  $(\Sigma_1 \cup \mathcal{N})^*$  such that  $\tilde{h}(A) = A$  for all  $A \in \mathcal{N}$ . Given a grammar  $G$  on  $\Sigma_1 \cup \mathcal{N}$ , we define  $\tilde{h}(G)$  as the grammar on  $\Sigma_2 \cup \mathcal{N}$  whose production rules are  $(\tilde{h}(w), \tilde{h}(w'))$  where  $(w, w')$  is a production rule of  $\mathcal{P}_1$ .

Prove the following statements:

4.  $h(\llbracket G \rrbracket) \subseteq \llbracket \tilde{h}(G) \rrbracket$ .
5. If the grammar  $G$  is context-free, then so is  $\tilde{h}(G)$  and  $h(\llbracket G \rrbracket) = \llbracket \tilde{h}(G) \rrbracket$ .
6. If the language  $\llbracket G \rrbracket$  is regular then so is  $\llbracket \tilde{h}(G) \rrbracket$ .  $\triangleleft$ : In general, the grammar  $\tilde{h}(G)$  is not regular.

**Solution:**

1. We prove that  $h(\varepsilon) = \varepsilon$  proving that the length of  $|h(\varepsilon)|$  is 0.

$$\begin{aligned} |h(\varepsilon)| &= |h(\varepsilon\varepsilon)| & \varepsilon &= \varepsilon\varepsilon \\ |h(\varepsilon)h(\varepsilon)| & & \text{Definition of morphism } h(ww') &= h(w)h(w') \\ |h(\varepsilon)| + |h(\varepsilon)| & & |ww'| &= |w| + |w'| \end{aligned}$$

Now we can subtract  $|h(\varepsilon)|$  on both sides:

$$\begin{aligned} |h(\varepsilon)| - |h(\varepsilon)| &= |h(\varepsilon)| + |h(\varepsilon)| - |h(\varepsilon)| \\ 0 &= |h(\varepsilon)| \end{aligned}$$

Since  $|h(\varepsilon)| = 0$ , then  $h(\varepsilon) = \varepsilon$ .

2. We need to prove that  $h(w) = h(w(1))h(w(2)) \dots h(w(n))$  (i.e., it is completely determined by the restriction of  $h$  to a single letter of the alphabet).

We can prove it by induction on the length of the word  $w$ .

- If  $w = \varepsilon$  then  $h(\varepsilon) = \varepsilon$  (from the result in 8.1).
- Consider a word  $w$  such that  $|w| = n > 0$ . We have  $h(w(1) \dots w(n-1)) = h(w(1)) \dots h(w(n-1))$  by inductive hypothesis.

$$\begin{aligned} h(w) &= h(w(1) \dots w(n-1)w(n)) \\ &= h(w(1) \dots w(n-1))h(w(n)) && \text{Since } h(ww') = h(w)h(w') \\ &= h(w(1)) \dots h(w(n-1))h(w(n)) && \text{By inductive hypothesis} \end{aligned}$$

3. We want to show that any map  $f : \Sigma_1 \rightarrow \Sigma_2^*$  extends to a morphism  $h : \Sigma_1^* \rightarrow \Sigma_2^*$ .

We can define a morphism  $h$  as the extension of  $f$  as follows:

$$h(w) = \begin{cases} \varepsilon & \text{if } w = \varepsilon \\ h(w')f(a) & \text{if } w = w'a, \text{ for some } w' \in \Sigma^* \end{cases}$$

We know that (if  $h$  is a morphism), it is uniquely determined by its restrictions to  $\Sigma_1$ . The map  $f$  is such restriction:  $h(a) = h(\varepsilon)f(a) = \varepsilon f(a) = f(a)$ .

We just need to prove that  $h$  is a morphism showing that for every  $w, w' \in \Sigma^*$   $h(ww') = h(w)h(w')$ . We prove it by induction on the length of strings  $w' \in \Sigma^*$ :

- When  $w' = \varepsilon$  we have  $h(w\varepsilon) = h(w)$
- Let  $|w'| = m > 0$  and  $|w| = n$ .

$$\begin{aligned} h(ww') &= h(w(1) \dots w(n)w'(1) \dots w'(m)) \\ &= h(w(1) \dots w(n)w'(1) \dots w'(m-1))f(w'(m)) && \text{Definition of } h \\ &= h(w(1) \dots w(n))h(w'(1) \dots w'(m-1))f(w'(m)) && \text{By inductive hypothesis} \\ &= h(w(1) \dots w(n))h(w'(1) \dots w'(m-1))h(w'(m)) && \forall a \in \Sigma. h(a) = f(a) \\ &= h(w(1) \dots w(n))h(w'(1) \dots w'(m-1)w'(m)) && \text{Definition of } h \end{aligned}$$

4. Let  $G = (\Sigma_1, \mathcal{N}, S, \mathcal{P}_1)$  and  $\tilde{h}(G) = (\Sigma_2, \mathcal{N}, S, \tilde{\mathcal{P}})$  where  $\tilde{h}(\mathcal{P}) = \{(\tilde{h}(w), \tilde{h}(w')) \mid (w, w') \in \mathcal{P}_1\}$ .

We prove that  $h(\llbracket G \rrbracket) \subseteq \llbracket \tilde{h}(G) \rrbracket$ .

$w \in h(\llbracket G \rrbracket)$  means that there exists a word  $v \in \Sigma_1^*$  and a derivation in  $G$  such that  $h(v) = w$  and  $S \xrightarrow[G]{*} v$ .

We show by induction on the length of the derivation that for any derivation  $S \xrightarrow[G]{*} v$  there exists a derivation  $S \xrightarrow[\tilde{h}(G)]{*} w$  and  $h(v) = w$ .

- **Base case:** If  $S \xrightarrow{1}_G v$  then there exists  $S \rightarrow v \in \mathcal{P}$ , and hence.  $S \rightarrow \tilde{h}(v) \in \tilde{h}(\mathcal{P}_1)$ . Thus there exists a derivation  $S \xrightarrow{\tilde{h}(G)} \tilde{h}(v)$ .

- **Inductive case:** Suppose  $S \xrightarrow{n}_G v' \xrightarrow{1}_G v$ .

Then there exists a rule  $(w, w') \in \mathcal{P}_1$  such that for two  $\alpha, \beta \in \Sigma_1^*$ ,  $v' = \alpha w \beta \xrightarrow{1}_G \alpha w' \beta = v$ .

By inductive hypothesis we also have that  $S \xrightarrow{n}_{\tilde{h}(G)} \tilde{h}(v')$ .

Also, we have that  $(\tilde{h}(w), \tilde{h}(w')) \in \tilde{h}(\mathcal{P})$  and that, since  $\tilde{h}$  is a morphism,  $\tilde{h}(v') = \tilde{h}(\alpha)\tilde{h}(w)\tilde{h}(\beta)$ .

We thus have the following derivation in  $\tilde{h}(G)$ :  $v' = \tilde{h}(\alpha)\tilde{h}(w)\tilde{h}(\beta) \xrightarrow{\tilde{h}(G)} \tilde{h}(\alpha)\tilde{h}(w')\tilde{h}(\beta)$ .

Then, since  $\tilde{h}(h)$  is a morphism, we have that  $\tilde{h}(\alpha)\tilde{h}(w')\tilde{h}(\beta) = \tilde{h}(\alpha w' \beta) = \tilde{h}(v)$

We conclude that there is a derivation:  $S \xrightarrow{n}_{\tilde{h}(G)} \tilde{h}(v') \xrightarrow{\tilde{h}(G)} \tilde{h}(v)$ .

Since for any derivation  $S \xrightarrow{*}_G v$  there exists a derivation  $S \xrightarrow{*}_{\tilde{h}(G)} w$  and  $h(v) = w$ , we conclude that if  $w \in \llbracket G \rrbracket$  then  $w \in \llbracket \tilde{h}(G) \rrbracket$

5. We start showing that if  $G$  is context-free then  $\tilde{h}(G)$  is context free.

We show that  $\forall (w, w') \in \tilde{h}(\mathcal{P})$ ,  $(w, w')$  is context-free. If  $(w, w') \in \tilde{h}(\mathcal{P})$  then there is a rule  $(v, v') \in \mathcal{P}_1$  such that  $\tilde{h}(v) = w$  and  $\tilde{h}(v') = w'$ . Since  $G$  is context-free, then  $(v, v')$  is such that  $v = A$ , for some non-terminal symbol  $A \in \mathcal{N}$ , and  $v' \in (\Sigma_1 \cap \mathcal{N})^*$ . Then, from the definition of  $\tilde{h}(G)$  and  $\tilde{h}$ , we know that  $w = \tilde{h}(v) = \tilde{h}(A) = A$ ,  $w' = \tilde{h}(v')$ , and  $\tilde{h}(v') \in (\Sigma_2 \cap \mathcal{N})^*$ . Thus, any rule  $(w, w') \in \tilde{h}(\mathcal{P})$  is also context-free. We conclude that  $\tilde{h}(G)$  is context-free.

We shown above in 8.4 that  $h(\llbracket G \rrbracket) \subseteq \llbracket \tilde{h}(G) \rrbracket$  for an unrestricted grammar. Here we prove that if  $G$  is context-free, then  $h(\llbracket G \rrbracket) \supseteq \llbracket \tilde{h}(G) \rrbracket$ . We show that if  $w \in \llbracket \tilde{h}(G) \rrbracket$  then  $w \in h(\llbracket G \rrbracket)$ .

If  $w \in \llbracket \tilde{h}(G) \rrbracket$  then there exists a derivation  $S \xrightarrow{*}_{\tilde{h}(G)} w$ . We show that there exists a derivation  $S \xrightarrow{*}_G v$  such that  $h(v) = w$ , and hence that if  $w \in \llbracket \tilde{h}(G) \rrbracket$ , then there exists a  $v \in \Sigma_1^*$  such that  $v \in \llbracket G \rrbracket$  and  $h(v) = w$ , and hence that  $w \in h(\llbracket G \rrbracket)$ .

We prove the existence of the derivation by induction on the derivation length.

- if  $S \xrightarrow{1}_{\tilde{h}(G)} w$ , then  $(S, w) \in \tilde{h}(\mathcal{P})$ . But there exists  $(S, v) \in \mathcal{P}_1$  such that  $w = \tilde{h}(v)$ . Hence  $S \xrightarrow{1}_G v$ , and  $h(v) = w$ .
- Suppose  $S \xrightarrow{n}_{\tilde{h}(G)} w' \xrightarrow{1}_{\tilde{h}(G)} w$ .

By inductive hypothesis we also know that  $S \xrightarrow{n}_G v'$  with  $\tilde{h}(v') = w'$ .

We miss to show that  $v' \xrightarrow{1}_G v$  with  $\tilde{h}(v) = w$ .

Since  $G$  is context-free, the derivation  $w' \xrightarrow{1}_{\tilde{h}(G)} w$  uses a rule  $(A, \delta) \in \tilde{h}(\mathcal{P})$  of the grammar

$\tilde{h}(G)$ , where  $A \in \mathcal{N}_1$ , such that, for some  $\alpha, \beta \in \Sigma_2^*$ ,  $w = \alpha A \beta \xrightarrow{\tilde{h}(G)} \alpha \delta \beta = w'$ .

Since  $(A, \delta) \in \tilde{h}(\mathcal{P})$ , then there exists  $(A, \delta') \in \mathcal{P}$  such that  $\tilde{h}(A) = A$  and  $\tilde{h}(\delta') = \delta$ .

We have that

$$\tilde{h}(v') = w' = \alpha A \beta = \tilde{h}(\alpha') \tilde{h}(A) \tilde{h}(\beta')$$

for some  $\alpha', \beta' \in \Sigma_1^*$ , and hence the derivation:  $v' = \alpha' A \beta' \xrightarrow{G} \alpha' \delta' \beta' = v$

We further now that  $\tilde{h}(v = \alpha' \delta' \beta') = \tilde{h}(\alpha') \tilde{h}(\delta) \tilde{h}(\beta') = w$ , and hence that  $v' \xrightarrow{G} v$  with  $\tilde{h}(v) = w$ .

Question: why did we need the assumption for  $G$  to be context-free? Can you find a counter-example showing that if  $G$  is an unrestricted grammar (i.e., a formal grammar without restrictions on the left or right-hand side, apart from avoiding  $\varepsilon$  in the left-hand side) then  $h(\llbracket G \rrbracket) \not\supseteq \llbracket \tilde{h}(G) \rrbracket$ ?

6. We want to show that if the language  $\llbracket G \rrbracket$  is regular, then so is  $\llbracket \tilde{h}(G) \rrbracket$ .

If  $\llbracket G \rrbracket$  is regular then there exists a regular grammar  $G'$  such that  $\llbracket G \rrbracket = \llbracket G' \rrbracket$ .

We show that if  $G'$  is a regular grammar, then also  $\tilde{h}(G')$  is a regular grammar. Given a rule in  $(w, w')$  in the production rules of  $\tilde{h}(G')$ , then there is a rule  $(v, v')$  in the production rules of  $G'$  such that  $\tilde{h}(v) = w$  and  $\tilde{h}(v') = w'$ . Since  $G'$  is regular, then  $(v, v')$  can either be in one of the three admissible forms for a rule in a regular grammar. We show that  $(w, w')$  is also of one of such form. The rule  $(v, v')$  can either be such that:

- $v = A$  and  $v' = \varepsilon$ . In this case, we have that  $w = \tilde{h}(A) = A$  and  $w' = \tilde{h}(\varepsilon) = \varepsilon$ ;
- $v = A$  and  $v' = a$ , for some  $a \in \Sigma_1^*$ . In this case  $w = \tilde{h}(A) = A$  and  $w' = \tilde{h}(a)$ . By definition we know that  $\tilde{h}(a) = b$  for some  $b \in \Sigma_2$ ;
- $v = A$  and  $v' = aA$ , for some  $a \in \Sigma_1^*$  and  $A \in \mathcal{N}$ . In this case  $w = \tilde{h}(A) = A$  and  $w' = \tilde{h}(aA) = \tilde{h}(a)\tilde{h}(A) = \tilde{h}(a)A = bA$ , where  $\tilde{h}(a) = b \in \Sigma_2$ .

All the rules in  $\tilde{h}(G')$  are of one of the kind admissible for a regular grammar and thus  $\tilde{h}(G')$  is a regular grammar.

Since  $\llbracket G \rrbracket = \llbracket G' \rrbracket$ , then we have that  $h(\llbracket G \rrbracket) = h(\llbracket G' \rrbracket)$

From 8.5, we also now that  $h(\llbracket G \rrbracket) = \llbracket \tilde{h}(G) \rrbracket$  and  $h(\llbracket G' \rrbracket) = \llbracket \tilde{h}(G') \rrbracket$ .

Thus we have:

$$\llbracket \tilde{h}(G) \rrbracket = h(\llbracket G \rrbracket) = h(\llbracket G' \rrbracket) = \llbracket \tilde{h}(G') \rrbracket$$

Since  $\llbracket \tilde{h}(G) \rrbracket = \llbracket \tilde{h}(G') \rrbracket$ , and  $\llbracket \tilde{h}(G') \rrbracket$  is regular, then also  $\llbracket \tilde{h}(G) \rrbracket$  is regular.