# Theory of Computation 6 Normalforms and Algorithms

**Frank Stephan**

**Department of Computer Science**

**Department of Mathematics**

**National University of Singapore**

**fstephan@comp.nus.edu.sg**

# Repetition 1

Let $(Q_1, \Sigma, \delta_1, s_1, F_1)$ and $(Q_2, \Sigma, \delta_2, s_2, F_2)$ be dfas which recognise $L_1$ and $L_2$, respectively.

Consider $(Q_1 \times Q_2, \Sigma, \delta_1 \times \delta_2, (s_1, s_2), F)$ with $(\delta_1 \times \delta_2)((q_1, q_2), a) = (\delta_1(q_1, a), \delta_2(q_2, a))$. This automaton is called a product automaton and one can choose $F$ such that it recognises the union or intersection or difference of the respective languages.

Union: $F = F_1 \times Q_2 \cup Q_1 \times F_2$;
Intersection: $F = F_1 \times F_2 = F_1 \times Q_2 \cap Q_1 \times F_2$;
Difference: $F = F_1 \times (Q_2 - F_2)$;
Symmetric Difference: $F = F_1 \times (Q_2 - F_2) \cup (Q_1 - F_1) \times F_2$.

# **Repetition 2 and Gaps Filled**

Regular languages are also closed under Kleene star, Kleene plus and concatenation: Use nfas for these and convert to dfas.

Context-free languages are closed under union, Kleene star, Kleene plus, concatenation and intersection with regular languages. They are in general not closed under intersection and complement.

Context-sensitive languages are closed under union, intersection, Kleene star, Kleene plus and concatenation. While these are easy to see, the following result is more difficult: They are also closed under complement (not part of this course).

Recursively enumerable languages are closed under union, intersection, Kleene star, Kleene plus and concatenation; they are not closed under complement.

# Repetition 3: Palindromes

The members of the language $\{x \in \Sigma^* : x = x^{mi}\}$ are called palindromes. A palindrome is a word or phrase which looks the same from both directions.

An example is the German name "OTTO"; furthermore, when ignoring spaces and punctuation marks, a famous palindrome is the phrase "A man, a plan, a canal: Panama." originating from the time when the canal in Panama was built.

The grammar with the rules $S \to aSa|aa|a|\varepsilon$ with $a$ ranging over all members of $\Sigma$ generates all palindromes; so for $\Sigma = \{0, 1, 2\}$ the rules of the grammar would be $S \to 0S0\,|\,1S1\,|\,2S2\,|\,00\,|\,11\,|\,22\,|\,0\,|\,1\,|\,2\,|\,\varepsilon$.

The set of palindromes is not regular.

# Repetition 4: Homomorphisms

## Example

$\mathbf{Ascii(Year\ 2021)} = \mathbf{5965617220323030231}$ represents each letter of "$\mathbf{Year\ 2021}$" by its two-digit hexadecimal ASCII representation.

## Definition 5.28

A homomorphism is a mapping $\mathbf{h}$ with domain $\mathbf{\Sigma}^*$ for some alphabet $\mathbf{\Sigma}$ which preserves concatenation:
$\mathbf{h(v \cdot w) = h(v) \cdot h(w)}$.

## Proposition 5.29

The homomorphism is determined by the images of the single letters and $\mathbf{h(w) = h(a_1) \cdot h(a_2) \cdot \ldots \cdot h(a_n)}$ for a word $\mathbf{w = a_1 a_2 \ldots a_n}$; $\mathbf{h(\varepsilon) = \varepsilon}$.

## Description 5.46

If $\mathbf{h}$ is a homomorphism then $\mathbf{h^{-1}(L) = \{w : h(w) \in L\}}$ is the inverse image of $\mathbf{L}$ under $\mathbf{h}$.

# Repetition 5: Homomorphism-Laws

## Theorem 5.32

The homomorphic images of regular and context-free languages are regular and context-free, respectively.

## Theorem 5.38

Every language generated by some grammar is the homomorphic image of a context-sensitive language; in particular the class of context-sensitive languages is not closed under homomorphism.

## Theorem 5.47

If $L$ is on level $k$ of the Chomsky hierarchy and $h$ is a homomorphism then $h^{-1}(L)$ is also on level $k$.

# Definition 6.1: Normal Forms

The normal forms are named after Noam Chomsky (born 7.12.1928) and Sheila Greibach (born 6.10.1939).

Assume that the language does not contain $\varepsilon$.

Chomsky Normal Form
All rules are of the form $A \to BC$ or $A \to d$ where $A, B, C$ are non-terminals and $d$ is a terminal.

Greibach Normal Form
All rules are of the form $A \to bw$ where $b$ is a terminal and $A$ a non-terminal and $w$ a (possibly empty) string of non-terminals.

If the language contains $\varepsilon$, one allows in both normal forms $S \to \varepsilon$ for the start symbol $S$ which then is not allowed to appear on the right side of a rule.

# CNF Algorithm Steps 1 and 2

Given context-free grammar $(N_0, \Sigma, P_0, S)$.

1. **Dealing with $\varepsilon$:** Let $N_1 = N_0 \cup \{S'\}$ and
   $P_1 = P_0 \cup \{S' \to S\}$;
   While there are rules $A \to vBw, B \to \varepsilon$ in $P_1$ with
   $A \to vw$ not in $P_1$ Do Begin $P_1 = P_1 \cup \{A \to vw\}$ End;
   Remove all rules $A \to \varepsilon$ with $A \in N_0$ from $P_1$;
   Continue with grammar $(N_1, \Sigma, P_1, S')$.

2. **Dealing with single terminal letters:** Let $N_2 = N_1$, $P_2 = P_1$;
   While there are $a \in \Sigma$ and rule $A \to vaw$ in $P_2$ with
   $vw \neq \varepsilon$ Do Begin Choose a new non-terminal $B \notin N_2$;
   Replace in all rules in $P_2$ all occurrences of $a$ by $B$;
   Update $N_2 = N_2 \cup \{B\}$ and $P_2 = P_2 \cup \{B \to a\}$ End;
   Continue with grammar $(N_2, \Sigma, P_2, S')$.

# CNF Algorithm Steps 3 and 4

**3. Breaking long ride hand sides:** Let $N_3 = N_2$ and $P_3 = P_2$; While there is $A \to Bw$ in $P_3$ with $|Bw| \geq 3$ Do Begin Choose a new non-terminal $C \notin N_3$ and let $N_3 = N_3 \cup \{C\}$; Add the rules $A \to BC, C \to w$ into $P_3$ and remove the rule $A \to Bw$ from $P_3$ End; Continue with grammar $(N_3, \Sigma, P_3, S')$.

**4. Removing rules $A \to B$:** Make a table of all $(A, B), (A, c)$ such that $A, B \in N_3$, $c \in \Sigma$ and, in the grammar $(N_3, \Sigma, P_3, S')$, $A \Rightarrow^* B$ and $A \Rightarrow^* c$, respectively; Let $N_4 = N_3$ and $P_4$ contain the following rules: $S' \to \varepsilon$ in the case that this rule is in $P_3$; $A \to a$ in the case that $(A, a)$ is in the table; $A \to BC$ in the case that there is $D \to EF$ in $P_3$ with $(A, D), (E, B), (F, C)$ in the table; The grammar $(N_4, \Sigma, P_4, S')$ is in Chomsky Normalform.

# Example for Chomsky Normalform

The grammar $(\{S\}, \{0, 1\}, \{S \rightarrow 0S0|1S1|00|11\}, S)$ generates all palindromes of even nonzero length.

Chomsky Normal Form
Non-terminals: $\{S, T, U, V, W\}$; terminals: $\{0, 1\}$; start symbol: $S$; rules: $S \rightarrow TV|UW$, $T \rightarrow VS|0$, $U \rightarrow WS|1$, $V \rightarrow 0$, $W \rightarrow 1$.

The derivation $S \Rightarrow 0S0 \Rightarrow 01S10 \Rightarrow 010010$ in the old grammar is equivalent to
$S \Rightarrow TV \Rightarrow VSV \Rightarrow 0SV \Rightarrow 0S0 \Rightarrow 0UW0 \Rightarrow 0WSW0 \Rightarrow 0WS10 \Rightarrow 01S10 \Rightarrow 01TV10 \Rightarrow 010V10 \Rightarrow 010010$
in the new grammar.

# Exercises 6.4 – 6.6

Bring the following grammars into Chomsky Normal Form.

Exercise 6.4
$(\{S, T\}, \{0, 1\}, \{S \to TTTT, T \to 0T1 | \varepsilon\}, S)$.

Exercise 6.5
$(\{S, T\}, \{0, 1\}, \{S \to ST | T, T \to 0T1 | 01\}, S)$.

Exercise 6.6
$(\{S\}, \{0, 1\}, \{S \to 0SS11SS0, 0110\}, S)$.

Always choose exercises of the right difficulty level: Students who are good in the field should take challenging exercises; students who have difficulties should take exercises they can master. So everyone should learn a bit from the exercises done by oneself; students should also try out exercises they do not present.

# Removing Useless Non-Terminals

Let $(N_0, \Sigma, P_0, S)$ be in Chomsky Normal Form.

1. **Removing non-terminating non-terminals:** Let $N_1 = \{A \in N_0 :$ there is a rule $A \to a$ or $A \to \varepsilon$ in $P_0\}$;
   While there is a rule $A \to BC$ in $P_0$ with $A \in N_0 - N_1$ and $B, C \in N_1$ Do Begin $N_1 = N_1 \cup \{A\}$ End;
   Let $P_1 = \{A \to w$ from $P_0$ with $A \in N_1$ and $w \in N_1 \cdot N_1 \cup \Sigma \cup \{\varepsilon\}\}$;
   If $S \notin N_1$ then terminate with empty grammar else continue with grammar $(N_1, \Sigma, P_1, S)$.

2. **Selecting all reacheable non-terminals:** Let $N_2 = \{S\}$;
   While there is a rule $A \to BC$ in $P_1$ with $A \in N_2$ and $\{B, C\} \not\subseteq N_2$ Do Begin $N_2 = N_2 \cup \{B, C\}$ End;
   Let $P_2 = \{A \to w$ from $P_1 : A \in N_2\}$;
   The grammar $(N_2, \Sigma, P_2, S)$ does not contain any useless non-terminal.

# Quiz 6.8

$\Sigma = \{0\}$; $N = \{Q, R, S, T, U, V, W, X, Y, Z\}$;
rules $S \to TU|UV$, $T \to UT|TV|TW$, $R \to VW|QQ|0$,
$Q \to 0$, $U \to VW|WX$, $V \to WX|XY|0$, $W \to XY|YZ|0$;
start symbol $S$.

Determine the set of reacheable and terminating non-terminals.

# Exercises 6.9 and 6.10

Exercise 6.9: Consider the grammar

$(\{S_0, S_1, \ldots, S_9\}, \{0\}, \{S_0 \to S_0S_0, S_1 \to S_2S_3,$
$S_2 \to S_4S_6|0, S_3 \to S_6S_9, S_4 \to S_8S_2, S_5 \to S_0S_5,$
$S_6 \to S_2S_8, S_7 \to S_4S_1|0, S_8 \to S_6S_4|0, S_9 \to S_8S_7\},$
$S_1)$.

Exercise 6.10: Consider the grammar

$(\{S_0, S_1, \ldots, S_9\}, \{0\}, \{S_0 \to S_1S_1, S_1 \to S_2S_2,$
$S_2 \to S_3S_3, S_3 \to S_0S_0|S_4S_4, S_4 \to S_5S_5,$
$S_5 \to S_6S_6|S_3S_3, S_6 \to S_7S_7|0, S_7 \to S_8S_8|S_7S_7,$
$S_8 \to S_7S_6|S_8S_6, S_9 \to S_7S_8|0\}, S_1)$.

For each of the two grammars, determine the set of reacheable and terminating non-terminals and explain the steps on the way to this set. What is the shortest word generated by the grammar?

# Emptyness Check

The removing of useless non-terminals results in an empty grammar in the case that the grammar does not generate any word. Here an algorithm for a context-free grammar $(N, \Sigma, P, S)$ which is not in any normalform.

**Initialisation:** Let $N' = \emptyset$;

**Loop:** While there are $A \in N - N'$ and a rule $A \to w$ with $w \in (N' \cup \Sigma)^*$
Do Begin $N' = N' \cup \{A\}$ End;

**Decision:** If $S \notin N'$ then the language of the grammar is empty else the language of the grammar constains some word.

# Finiteness Check

Let $(N, \Sigma, P, S)$ be in Chomsky Normal Form. One computes $N''(A)$ as the set of non-terminals occurring in a derivation of a terminal word from $A$ after at least one step.

**Initialisation 1:** Let $N' = \emptyset$;

**Loop 1:** While there are $A \in N - N'$ and a rule $A \to w$ with $w \in (N' \cup \Sigma)^*$
Do Begin $N' = N' \cup \{A\}$ End;

**Initialisation 2:** For all $A \in N$, let $N''(A) = \emptyset$;

**Loop 2:** While there are $A, B, C, D \in N'$ and a rule $B \to CD$ with $B \in N''(A) \cup \{A\}$ and $\{C, D\} \not\subseteq N''(A)$
Do Begin $N''(A) = N''(A) \cup \{C, D\}$ End;

**Decision:** If there is $A \in N''(S) \cup \{S\}$ with $A \in N''(A)$ then the language of the grammar is infinite else it is finite.

# Exercise 6.13: Time Complexity

Polynomial time means that the algorithm uses time $p(n)$ for inputs of length $n$ to decide some property. If a grammar is in Chomsky Normal Form, $n$ can just be the number of non-terminals.

The checks whether a grammar in CNF generates some word or generates infinitely many words are in polynomial time. How complex is it to not only check whether some word is generated but also to output one such witness? Possible answers are polynomial time, exponential time and double exponential time. Here the time to output a word of length $m$ is $m$ computation steps. Give reasons for the answer.

# Exercises 6.14 – 6.17

For grammar $(\{S, T, U, V, W\}, \{0, 1, 2\}, P, S)$ with below rules determine how many words the grammar generates: (a) None, (b) One, (c) Two, (d) Three, (e) Finitely many and at least four, (f) Infinitely many?

Exercise 6.14: The rules are $S \to TT$, $T \to UU$, $U \to VW | WV$, $V \to 0$, $W \to 1$.

Exercise 6.15: The rules are $S \to ST$, $T \to TU$, $U \to UV$, $V \to VW$, $W \to 0$.

Exercise 6.16: The rules are $S \to UT | TU | 2$, $T \to VV$, $U \to WW$, $V \to 0$, $W \to 1$.

Exercise 6.17: The rules are $S \to SS | TT | UU$, $T \to VV$, $U \to WW$, $V \to 0$, $W \to WW$.

# Derivation Tree

For grammar $(\{S, T, U\}, \{0, 1\}, \{S \to SS|TU|UT,$
$U \to 0|US|SU, T \to 1|TS|ST\}, S)$, a derivation $S \Rightarrow TU \Rightarrow$
$TSU \Rightarrow TUTU \Rightarrow 1UTU \Rightarrow 10TU \Rightarrow 101U \Rightarrow 1010$ can
be represented by a tree:

# Properties of Derivation Trees

Each node has a symbol in it.

If the symbol is a non-terminal then the successor nodes have the symbols of the rule applied at this symbol; if the symbol is a terminal then the node is a leaf.

The derivation does not depend on the order of rules on incomparable nodes are applied; however, the rule of each node has to be applied before it is applied to any successor.

The leaves below a node have the part of the word generated from the corresponding symbol.

Trees are more precise then derivations by lists of rules, as they explicitly say which non-terminal derives into which symbols and these are on the successor nodes.

Sometimes a derivation is ambiguous: For $TU \Rightarrow TSU$, the rule can be $T \to TS$ or $U \to SU$.

# Number of Derivation Trees

As the rules to derive a word are not always unique but can be ambiguous, there can be different derivation trees for the same word. Indeed, one can prove that for some grammars it is impossible to have the derivation tree all the time unique. Consider the grammar
$(\{S, T, U\}, \{0, 1\}, \{S \rightarrow SS|TU|UT, U \rightarrow 0|US|SU,$
$T \rightarrow 1|TS|ST\}, S)$, how many derivation trees has the word below?

Exercise 6.19: How many derivation trees are there for $011001$?

Exercise 6.20: How many derivation trees are there for $000111$?

# Exercise 6.21

One can generalise the notion of derivation tree to all context-free grammars and might then have nodes with more than two successors. Consider the grammar

$$(\{S, T\}, \{0, 1, 2\}, \{S \to TT, T \to 0T1 | 2\}, S).$$

and draw the derivation tree for the word $00211021$. Prove that all words in the language of this grammar have a unique derivation tree.
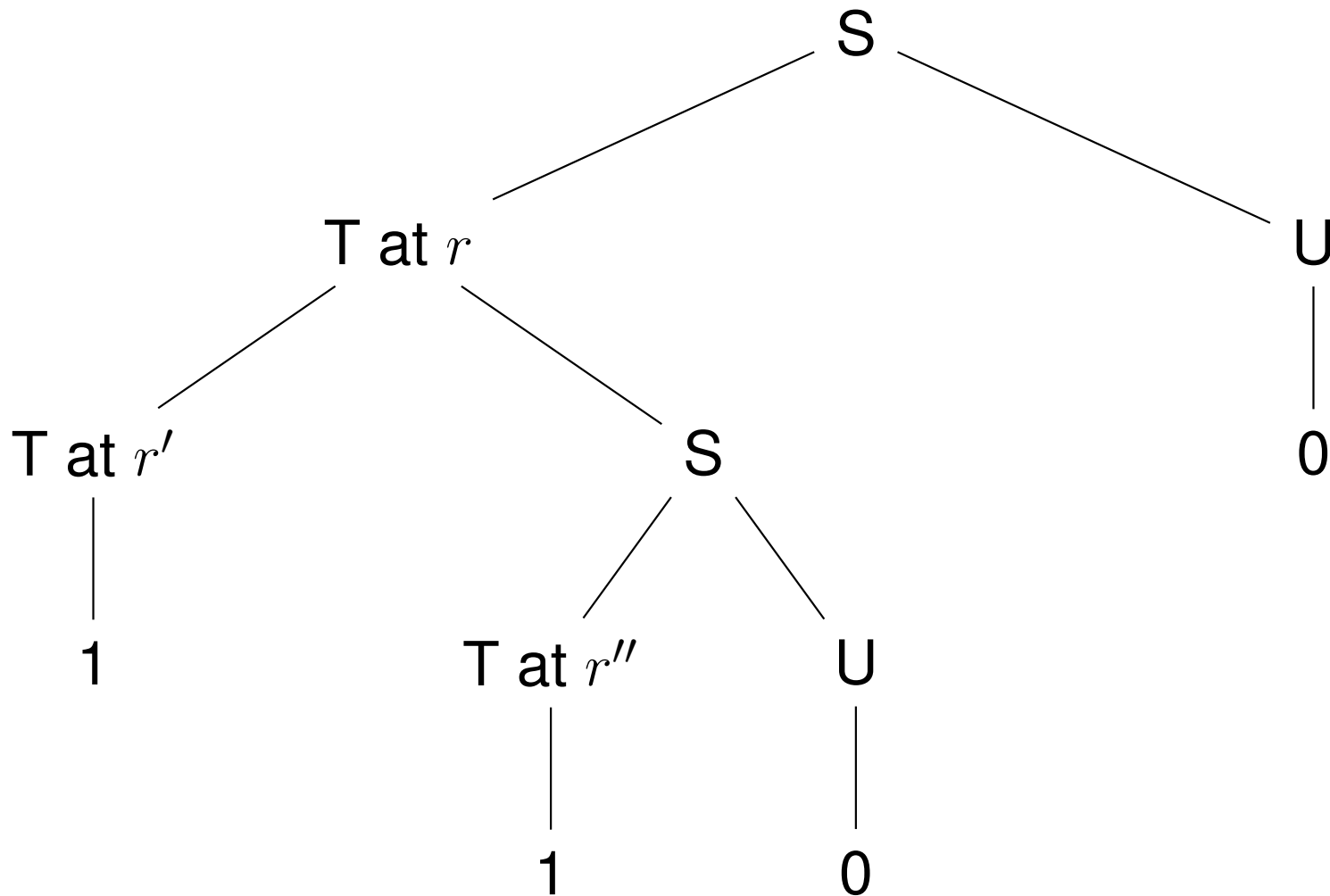
# Pumping Lemma

Theorem 2.15 (b)

Let $L \subseteq \Sigma^*$ be an infinite context-free language generated by a grammar $(N, \Sigma, P, S)$ in Chomsky Normal Form with $h$ non-terminals. Then the constant $k = 2^{h+1}$ satisfies that for every $u \in L$ of length at least $k$ there is a representation $vwxyz = u$ such that $|wxy| \leq k$, $(w \neq \varepsilon$ or $y \neq \varepsilon)$ and $vw^\ell xy^\ell z \in L$ for all $\ell \in \mathbb{N}$.

Proof Idea

In the derivation tree of a word $u$ longer than $2^{h+1}$, find the lowest node having a non-terminal $A$ which is also on a branch somewhere below. Below $A$, no branch repeats a non-terminal and therefore each branch has at most the length $h$. Thus there are at most $2^{h+1}$ leaves below $A$. Since $A$ itself repeats, one has that $S \Rightarrow^* vAz \Rightarrow^* vwAyz \Rightarrow^* vwxyz$ and $|wxy| \leq 2^{h+1}$ and $A \Rightarrow^* wAy$ and $wy \neq \varepsilon$.

# Example of Derivation Tree



Lowest node $\mathbf{r}$ with non-terminal repeated below at $\mathbf{r'}$ or $\mathbf{r''}$.
Possible Pumpings: $\mathbf{1(10)^{\ell}0}$, $\mathbf{1^{\ell}10^{\ell}0}$.

# Ogden's Lemma

Mark the first four $1$ in word $0000011111$ as $000001{\color{red}1111}$.
The word can be pumped such that at least one but at most four marked symbols are pumped or between the pumped parts: $000\,0^\ell\,0{\color{red}1}\,{\color{red}1}^\ell\,{\color{red}111}$.

## Theorem 6.22

Let $L \subseteq \Sigma^*$ be an infinite context-free language generated by a grammar $(N, \Sigma, P, S)$ in Chomsky Normal Form with $h$ non-terminals. Then the constant $k = 2^{h+1}$ satisfies that for every $u \in L$ with at least $k$ marked symbols, there is a representation $vwxyz = u$ such that $wxy$ contains at most $k$ marked symbols, $wy$ contains at least $1$ marked symbol and $vw^\ell xy^\ell z \in L$ for all $\ell \in \mathbb{N}$.

# Example 6.23

## The Language

Let $L$ be the language of all words $w \in 1^+(0^+1^+)^+$ with no two runs of zeroes of equal lengths. That is, words in $L$ start and end with $1$ and do not contain a subword $10^h1$ twice for any $h > 0$.

$L$ contains words $1000101$ and $1100110001$ and $1111111$. $L$ neither contains $100111001$ nor $10100110100001$.

## Traditional Context-Free Pumping Lemma Satisfied

More precisely: Each word longer than $2$ symbols in $L$ can somewhere be pumped by a single one-symbol pump.

If $w \in 1^+0^+1^+$ then pump some $0$ in the middle.

If $w$ contains a border of at least two $1$ then pump a $1$ in this border.

If $w$ has a longest run of $0$ separated by a single $1$ from another run of $0$ then pump this separating $1$.

# L does not satisfy Ogden's Lemma

If $k$ is a constant supposed to work for Ogden's Pumping Lemma then consider $u = 1010^2 10^3 1 \ldots 10^{4k} 1$ and mark the zeroes of the subword $10^k 1$.

Let $vwxyz$ split $u$ into five parts.

If $w$ or $y$ in $\{0, 1\}^+ - \{0\}^* - \{1\}^*$ then $vw^4 xy^4 z$ contains a subword $10^h 1$ at least twice.

One of $w$ and $y$ must be in $0^+$ and contain some of the marked zeroes.

In the word $vwwxyyz$ pumping has replaced $10^k 1$ by $10^{k+h} 1$ for $h \in \{1, 2, \ldots, k\}$ and thus the subword $10^{k+h} 1$ in $u$ must also be pumped and replaced by $10^{k+h+\ell} 1$ for $\ell \in \{1, 2, \ldots, k+h\}$. However, $10^{k+h+\ell} 1$ is also a subword of $u$ which occurs in $vwwxyyz$ and thus $vwwxyyz \notin L$ (as $h \leq k$ and $\ell \leq 2k$ and $h + k + \ell \leq 4k$).

# Square-Containing Words

Theorem 6.24 [Ehrenfeucht and Rozenberg 1983; Ross and Winklmann 1982]
The language $L$ of square-containing words over $\{0, 1, 2\}$ is not context-free.

Proposition 6.25
The language $L$ of square-containing words over $\{0, 1, 2\}$ satisfies Ogden's Lemma with constant $6$.

If $u = 012012102012012102$ then one can see that $012(01)^*2102012012102 \subseteq L$. Details on next page.

# Proof of Proposition

For given word $u \in L$ with at least six letters marked, consider those splittings $vwxyz = u$ where the following conditions are met:

- $v$ ends with the same letter $a$ with which $z$ starts;

- $w$ contains at least one marked letter;

- $x$, $y$ are $\varepsilon$.

Such splittings exist obviously when there is $a \in \{0, 1, 2\}$ such that there are three marked letters $a$ or there are two marked letters $a$ with one other marked letter in between.

The remaining case is wlog of form $r_0 0 r_1 0 r_2 1 r_3 1 r_4 2 r_5 2 r_6$. If $r_1$ is in $0^*$ then one let $w = 1$ and $v$ everything left of the first $1$ and $z$ everything right of it. If $r_1$ contains a $b \neq 0$ then one let $w$ be all the part right of this $b$ and left of the first marked $b$ and $w$ contains the marked letter $0$.

# Exercises 6.27 and 6.28

## Exercise 6.27

Prove that the language

$$L = \{a^h \cdot w : a \in \{0, 1, 2\}, w \in \{0, 1, 2\}^*, w \text{ is}$$
$$\text{square-free and } h \in \mathbb{N}\}$$

satisfies Theorem 2.15 (b) but does not satisfy Ogden's Pumping Lemma. The fact that there are infinitely many square-free words can be used without proof.

## Exercise 6.28

Use the Block Pumping Lemma to prove the following variant of Ogden's Lemma for regular languages: If a language $L$ satisfies the Block Pumping Lemma with constant $k + 1$ then one can, for each word $u$ having at least $k$ marked symbols, find a splitting of the word into parts $x, y, z$ such that $u = xyz$ and $xy^*z \subseteq L$ and $y$ contains at least $1$ and at most $k$ marked symbols.

# Greibach Normal Form

Example 6.29

Consider $(\{S\}, \{0, 1\}, \{S \rightarrow 0S0|1S1|00|11|0|1\}, S)$. For the Greibach Normal Form, one needs two additional non-terminals $T, U$ and updates the rules as follows:

$$S \rightarrow 0ST|1SU|0T|1U|0|1, \ T \rightarrow 0, \ U \rightarrow 1.$$

Consider $(\{S\}, \{0, 1\}, \{S \rightarrow SS|0S1|1S0|10|01\}, S)$. For the Greibach Normal Form, one needs two additional non-terminals $T, U$ and updates the rules as follows:

$$S \rightarrow 0SU|0U|1ST|1T, \ T \rightarrow 0|0S, \ U \rightarrow 1|1S.$$

Exercises 6.30 and 6.31
Construct Greibach Normal Form for the intersection of $0*1*0*1*$ with each of these two languages, respectively.

# Languages and Derivatives

Exercise 6.32. Prove the following rules of the derivative with $a \in \Sigma$ and $x \in \Sigma^*$:

- $(L \cup H)_x = L_x \cup H_x$ and $(L \cap H)_x = L_x \cap H_x$;
- If $\varepsilon \in L$ then $(L \cdot H)_a = L_a \cdot H \cup H_a$ else $(L \cdot H)_a = L_a \cdot H$;
- $(L^*)_a = L_a \cdot L^*$.

The following theorem characterises the context-free languages by stating that all derivatives have to be formed from a finite set of languages using concatenation and union; furthermore, all the derivatives of these finitely many languages satisfy the same condition.

Theorem 6.33. A language $L$ is context-free iff there is a finite list of languages $H_1, H_2, \ldots, H_n$ with $L = H_1$ such that for every word $x$ and every $H_m$, $(H_m)_x$ is a finite union of finite products of some $H_k$.

Exercise 6.34. Prove this theorem using Exercise 6.32 and Existence of Greibach Normal Form.