# Theory of Computation 8 Deterministic Membership Testing

**Frank Stephan**

**Department of Computer Science**

**Department of Mathematics**

**National University of Singapore**

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

# Repetition 1: Normal Forms

The normal forms are named after Noam Chomsky (born 1928) and Sheila Greibach (born 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) word 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.
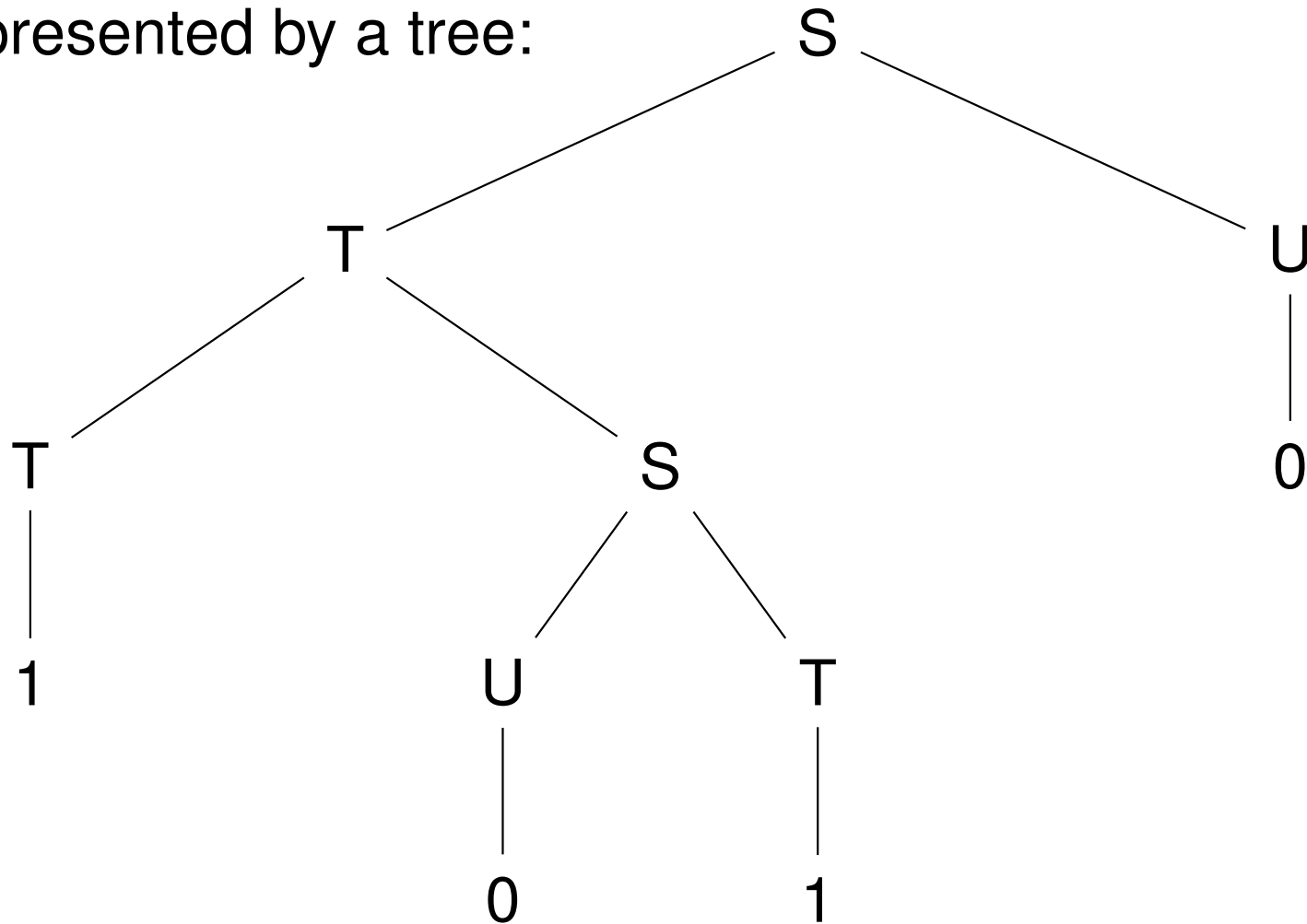
# Repetition 2: Algorithms

Given a context-free grammar as input, algorithms for the following tasks where presented:

- Converting the grammar into Chomsky Normal Form;

- Testing whether the grammar generates a word;

- Testing whether the grammar generates infinitely many words.

All these algorithms run in time polynomial in the size of the grammar; for grammars in Chomsky Normal Form, the number of non-terminals can serve as the size.

# Repetition 3: Derivation Tree

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

# Repetition 4

## Marked Symbols

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: $000\,0^\ell\,01\,1^\ell\,111$.

## Ogden's Lemma

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 $\ell \in \mathbb{N}$.

# Membership Testing

For a language $L$ and a word $w$ of length $n$, one wants to decide whether $w \in L$. The following will be shown:

Regular language: Done by a finite automaton, time $O(n)$.

Linear language: Special case of Cocke, Kasami and Younger's algorithm, time $O(n^2)$.

Context-free language: Cocke, Kasami and Younger's algorithm, time $O(n^3)$.

Context-sensitive language: Savitch's algorithm, space $O(n^2)$, time $O(c^{n^2})$ for some $c$.

Recursively enumerable language: No algorithm (undecidable).

# Cocke, Kasami and Younger

Let $(N, \Sigma, P, S)$ be in Chomsky Normal Form and $a_1 a_2 \ldots a_n$ be the input word.

**1. Initialisation:** For all $k$,

$$E_{k,k} = \{A \in N : A \to a_k \text{ is a rule}\}.$$

**2. Loop:** Go through all pairs $(i, j)$ such that they are processed in increasing order of $j - i$ and let

$$E_{i,j} = \{A : \exists \text{ rule } A \to BC \; \exists k$$
$$[i \leq k < j \text{ and } B \in E_{i,k} \text{ and } C \in E_{k+1,j}]\}.$$

**3. Decision:** Word is generated by the grammar iff $S \in E_{1,n}$.

Set $E_{i,j}$ contains all non-terminals generating $a_i \ldots a_j$.
Time $O(n^3)$: $O(n^2)$ values $E_{i,j}$ with $O(n)$ choices of $k$.

# Example 8.2, Word 0011

Grammar $(\{S, T, U\}, \{0, 1\}, \{S \to SS|TU|UT,$
$U \to 0|US|SU, T \to 1|TS|ST\}, S)$. Entries $E_{i,j}$:

$$E_{1,4} = \{S\}$$

$$E_{1,3} = \{U\} \qquad E_{2,4} = \{T\}$$

$$E_{1,2} = \emptyset \qquad E_{2,3} = \{S\} \qquad E_{3,4} = \emptyset$$

$$E_{1,1} = \{U\} \quad E_{2,2} = \{U\} \qquad E_{3,3} = \{T\} \quad E_{4,4} = \{T\}$$

$$0 \qquad\qquad 0 \qquad\qquad 1 \qquad\qquad 1$$

As $S \in E_{1,4}$, the word $0011$ is in the language.

# Example 8.2, Word 0111

Grammar $(\{S, T, U\}, \{0, 1\}, \{S \to SS|TU|UT,$
$U \to 0|US|SU, T \to 1|TS|ST\}, S)$. Entries $E_{i,j}$:

$$E_{1,4} = \emptyset$$
$$E_{1,3} = \{T\} \quad E_{2,4} = \emptyset$$
$$E_{1,2} = \{S\} \quad E_{2,3} = \emptyset \quad E_{3,4} = \emptyset$$
$$E_{1,1} = \{U\} \quad E_{2,2} = \{T\} \quad E_{3,3} = \{T\} \quad E_{4,4} = \{T\}$$
$$0 \qquad\qquad 1 \qquad\qquad 1 \qquad\qquad 1$$

As $S \notin E_{1,4}$, the word $0111$ is not in the language.

Quiz 8.3
Make the table for word $1001$.

# Cocke, Kasami and Younger

## Exercise 8.4

Consider the grammar $(\{S, T, U, V, W\}, \{0, 1, 2\}, P, S)$ with $P$ consisting of the rules $S \rightarrow TT$, $T \rightarrow UU|VV|WW$, $U \rightarrow VW|WV|VV|WW$, $V \rightarrow 0$, $W \rightarrow 1$. Make the entries of the Algorithm of Cocke, Kasami and Younger for the words $0011$, $1100$ and $0101$.

## Exercise 8.5

Consider the grammar $(\{S, T, U, V, W\}, \{0, 1, 2\}, P, S)$ with $P$ consisting of the rules $S \rightarrow ST|0|1$, $T \rightarrow TU|1$, $U \rightarrow UV|0$, $V \rightarrow VW|1$, $W \rightarrow 0$. Make the entries of the Algorithm of Cocke, Kasami and Younger for the word $001101$.

# Linear Grammars

Linear languages sit between regular and context-free languages.

## Definition

A grammar is linear iff on every left side of a rule exactly one non-terminal and on every right side of a rule at most one non-terminal and arbitrary many terminals.

## Normal Form

A linear grammar for a language not containing $\varepsilon$ is in normal form iff all rules are of one of the following forms: $A \to Bc$, $A \to cB$, $A \to c$ where $A, B$ are non-terminals and $c$ is a terminal.

If $\varepsilon$ is in the language then one has the rule $S \to \varepsilon$ for the start symbol $S$ and $S$ does not appear on any right side of a rule.

# Parsing Linear Grammars

Let $(N, \Sigma, P, S)$ be a linear grammar in normal form and $a_1 a_2 \ldots a_n$ be an input word.

**1. Initialisation:** For all $k$,

$$E_{k,k} = \{A \in N : A \to a_k \text{ is in } P\}.$$

**2. Loop:** Process all pairs $(i, j)$ with $i < j$ in increasing order of $j - i$ and let

$$E_{i,j} = \{A : \quad \exists \text{ rule } A \to Bc \quad [B \in E_{i,j-1} \text{ and } c = a_j] \text{ or}$$
$$\exists \text{ rule } A \to cB \quad [c = a_i \text{ and } B \in E_{i+1,j}]\}.$$

**3. Decision:** The word is generated by the grammar iff $S \in E_{1,n}$.

# Example 8.8, Word 0110

The grammar

$$(\{S, T, U\}, \{0, 1\}, \{S \rightarrow 0|1|0T|1U, T \rightarrow S0|0, U \rightarrow S1|1\}, S)$$

is a linear grammar for palindromes. For the word $0110$, the entries are

$$E_{1,4} = \{S\}$$

$$E_{1,3} = \emptyset \qquad\qquad E_{2,4} = \{T\}$$

$$E_{1,2} = \{U\} \qquad E_{2,3} = \{S, U\} \qquad E_{3,4} = \{T\}$$

$$E_{1,1} = \{S, T\} \quad E_{2,2} = \{S, U\} \qquad E_{3,3} = \{S, U\} \quad E_{4,4} = \{S, T\}$$

$$0 \qquad\qquad\qquad 1 \qquad\qquad\qquad 1 \qquad\qquad\qquad 0$$

and as $S \in E_{1,4}$, the word is accepted.

# Example 8.8, Word 1110

For processing the word $1110$, one gets the following table:

$$E_{1,4} = \{T\}$$

$$E_{1,3} = \{S, U\} \qquad E_{2,4} = \{T\}$$

$$E_{1,2} = \{S, U\} \qquad E_{2,3} = \{S, U\} \qquad E_{3,4} = \{T\}$$

$$E_{1,1} = \{S, U\} \qquad E_{2,2} = \{S, U\} \qquad E_{3,3} = \{S, U\} \quad E_{4,4} = \{S, T\}$$

$$\quad 1 \qquad\qquad\qquad 1 \qquad\qquad\qquad 1 \qquad\qquad\qquad 0$$

## Exercise 8.9

Make the entries for the word $0110110$ for above grammar.

# Exercises 8.10 and 8.11

Consider the following linear grammar:

$$(\{S, T, U\}, \{0, 1\}, \{S \rightarrow 0T|T0|0U|U0,\ T \rightarrow 0T00|1,$$
$$U \rightarrow 00U0|1\}, S).$$

Convert the grammar into the linear grammar normal form and determine the $E_{i,j}$ for input $00100$.

## Exercise 8.11

Which two of the following languages are linear? Provide linear grammars for these two languages:

- $L = \{0^n 1^m 2^k : n + k = m\}$;

- $H = \{0^n 1^m 2^k : n + m = k\}$;

- $K = \{0^n 1^m 2^k : n \neq m \text{ or } m \neq k\}$.

# Kleene Star

Algorithm 8.12

Let $L$ be generated by a linear grammar and $a_1 a_2 \ldots a_n$ be a word. To check whether $a_1 a_2 \ldots a_n \in L^*$, do the following:

**First Part:** Compute for each $i, j$ with $1 \leq i \leq j \leq n$ the set $E_{i,j}$ of all non-terminals which generate $a_i a_{i+1} \ldots a_j$.

**Initialise Loop for Kleene Star:** Let $F_0 = 1$.

**Loop for Kleene Star:** For $m = 1, 2, \ldots, n$ Do Begin If there is a $k < m$ with $S \in E_{k+1,m}$ and $F_k = 1$ Then let $F_m = 1$ Else let $F_m = 0$ End.

**Decision:** $w \in L^*$ iff $F_n = 1$.

First Part is $O(n^2)$, Loop for Kleene Star is $O(n^2)$ .

# Other Combinations

Assume that $H, K, L$ are linear languages and that one has computed for $a_1 \ldots a_n$ the entries $E_{i,j}^L, E_{i,j}^H, E_{i,j}^K$ say whether the word $a_i \ldots a_j$ is in $L, H, K$, respectively. This information can be computed in $O(n^2)$. Complete the algorithm to do the following check in $O(n^2)$ time:

Exercise 8.13
Is $a_1 a_2 \ldots a_n \in L \cdot H \cdot K$?

Exercise 8.14
Is $a_1 a_2 \ldots a_n \in (L \cap H)^* \cdot K$?

# Regular Closure of CTF

A language $H$ is said to be in the regular closure of the context-free languages iff it is obtained by combining finitely many context-free languages with intersection, union, concatenation, set-difference, Kleene Star and Kleene Plus. An example is

$$H = (L_1 \cap L_2)^* \cdot L_3 \cdot (L_1 \cap L_2)^* - L_4$$

and here $L_1, L_2, L_3, L_4$ are context-free.

Exercise 8.15
Prove by structural induction that every language $H$ which is in the regular closure of the context-free languages has an $O(n^3)$ decision algorithm.

# Counting Derivation Trees

One can modify the algorithm of Cocke, Kasami and Younger to count derivation trees.

Given grammar in CNF $(N, \Sigma, P, S)$, a word $w = a_1 a_2 \ldots a_n$ and all $A \in N$, let $E_{i,j}$ denote the set of all nonterminals which generate the characters $a_i \ldots a_j$ and $D_{i,j,A}$ denote the number of derivation trees which can, with root $A$, derive a word $a_i \ldots a_j$. Let $P$ be the set of rules. For each $A \in N$, if $A \in E_{i,i}$, there is exactly one tree with $A \to a_i$ and so $D_{i,i,A} = 1$ else $D_{i,i,A} = 0$. If $i < j$ then

$$D_{i,j,A} = \sum_{(B,C):\ A \to BC\ \in\ P}\ \sum_{k:\ i \le k < j} D_{i,k,B} \cdot D_{k+1,j,C}$$

and the overall number of derivation trees is $D_{1,n,S}$.

# Exercises 8.17 - 8.20

Let $P$ contain the rules $V \to VV|WW|0$ and $W \to VW|WV|1$. Consider the grammars $G = (\{V, W\}, \{0, 1\}, P, W)$ and $H = (\{U, V, W\}, \{0, 1, 2\}, P \cup \{U \to VU|UV|2\}, U)$.

## Exercise 8.17
How many derivation trees has $0011100$ in $G$?

## Exercise 8.18
How many derivation trees has $0000111$ in $G$?

## Exercise 8.19
How many derivation trees has $021111$ in $H$?

## Exercise 8.20
How many derivation trees has $010012$ in $H$?

# Polynomial Space

An algorithm can be measured by

- The time needed to do the computation;

- The space (size of variables and arrays and ...) needed to do the computation.

Let $n$ be a parameter to measure the size of the input.

If one can be computed in time $F(n)$ then, under certain assumptions to the machine model, it can also be done in space $F(n)$.

However, for the converse, only a rough estimate is known: If something can be computed in space $F(n)$ then one can compute it in time $2^{F(n)}$.

# Theorem of Savitch

Algorithm 8.21

Context-senstive grammar $(N, \Sigma, P, S)$, input word $w$.

**Recursive Call:** Function $\mathbf{Check}(u, v, t)$
  Begin If $u = v$ or $u \Rightarrow v$ Then Return$(1)$;
  If $t \leq 1$ and $u \neq v$ and $u \not\Rightarrow v$ Then Return$(0)$;
  Let $t' = t/2$; Let $r' = 0$;
  For all $u' \in (N \cup \Sigma)^*$ with $|u| \leq |u'| \leq |v|$ Do
  Begin If $\mathbf{Check}(u, u', t') = 1$ and $\mathbf{Check}(u', v, t') = 1$
  Then $r' = 1$ End; Return$(r')$ End.

**Decision:** If $\mathbf{Check}(S, w, k^n) = 1$ Then $w \in L$ Else $w \notin L$.

Space Complexity, per call $O(n)$, in total $O(n^2)$;
Value of $t$: $k^n/2^h$ in depth $h$ of recursion $(k = |\Sigma| + |N| + 1)$;
Number of nested calls: $O(\log(k^n)) = O(\log(k) \cdot n)$.
Runtime: $O(c^{n^2})$ for any $c > (2k)^{\log(k)}$.

# Example

Check whether $S \Rightarrow^* w$ within $8$ steps. One call $\mathbf{Check}(\mathbf{S}, \mathbf{w}, \mathbf{8})$.

Is there a word $\mathbf{v}$ such that $S \Rightarrow^* v$ and $v \Rightarrow^* w$ both within $4$ steps? For each $\mathbf{v}$, two calls $\mathbf{Check}(\alpha, \beta, \mathbf{4})$ one after each other with $(\alpha, \beta) = (\mathbf{S}, \mathbf{v}), (\mathbf{v}, \mathbf{w})$.

Is there a word $\mathbf{v}'$ such that $\alpha \Rightarrow^* v'$ and $v' \Rightarrow^* \beta$ both within $2$ steps? For each word $\mathbf{v}'$, two calls $\mathbf{Check}(\alpha', \beta', \mathbf{2})$ one after each other with $(\alpha', \beta') = (\alpha, \mathbf{v}'), (\mathbf{v}', \beta')$.

Is there a word $\mathbf{v}''$ such that $\alpha' \Rightarrow^* v''$ and $v'' \Rightarrow^* \beta'$ both within $1$ steps? For each word $\mathbf{v}''$, two calls $\mathbf{Check}(\alpha'', \beta'', \mathbf{1})$ one after each other with $(\alpha'', \beta'') = (\alpha', \mathbf{v}''), (\mathbf{v}'', \beta')$.

Is $\alpha'' = \beta''$ or $\alpha'' \Rightarrow \beta''$ true? No further call needed, bottom of recursion reached.

# Example 8.22

Consider grammar $(\{S, T, U, V, W\}, \{0, 1\}, P, S)$ with rules $P$ being the following:

$$S \to 0S \,|\, U, \quad U \to V \,|\, 0, \quad 0V \to 1U, \quad V \to 1, \quad 1V \to WU,$$
$$1W \to W0, \quad 0W \to 10.$$

This gives a binary counter with $S \to 0S$ generating enough digits before doing $S \to U$. $U$ stands for last digit $0$, $V$ stands for last digit $1$, $W$ stands for a digit $0$ still having a carry bit to pass on.

Deriving a binary number $k$ needs at least $k$ steps; as the length $n$ of $k$ is logarithmic in $k$, the derivation length can be exponential in $n$. In particular deriving $1^n$ needs more than $2^n$ steps.

# Exercises 8.23 and 8.24

Exercise 8.23

Give a proof that there are $k^n$ or less words of length up to $n$ over the alphabet $\Sigma \cup N$ with $k-1$ symbols.

Exercise 8.24

Modify Savitch's Algorithm such that it computes the length of the shortest derivation of a word $w$ in the context-sensitive grammar, provided that such derivation exists. If it does not exist, the algorithm should return the special value $\infty$.

# Naive Algorithm

Exercise 8.25

What is the time and space complexity of this naive algorithm?

**Recursive Call:** Function $\mathbf{Check}(\mathbf{u}, \mathbf{w}, \mathbf{t})$

Begin If $\mathbf{u} = \mathbf{w}$ or $\mathbf{u} \Rightarrow \mathbf{w}$ Then Return$(\mathbf{1})$;

If $\mathbf{t} \leq \mathbf{1}$ and $\mathbf{u} \neq \mathbf{v}$ and $\mathbf{u} \not\Rightarrow \mathbf{w}$ Then Return$(\mathbf{0})$;

Let $\mathbf{r}' = \mathbf{0}$;

For all $\mathbf{v} \in (\mathbf{N} \cup \mathbf{\Sigma})^*$ with $\mathbf{u} \Rightarrow \mathbf{v}$ and $|\mathbf{v}| \leq |\mathbf{w}|$ Do

Begin If $\mathbf{Check}(\mathbf{v}, \mathbf{w}, \mathbf{t} - \mathbf{1}) = \mathbf{1}$ Then $\mathbf{r}' = \mathbf{1}$ End;

Return$(\mathbf{r}')$ End;

**Decision:** If $\mathbf{Check}(\mathbf{S}, \mathbf{w}, \mathbf{k^n}) = \mathbf{1}$ Then $\mathbf{w} \in \mathbf{L}$ Else $\mathbf{w} \notin \mathbf{L}$.

# Growing Grammars

Definition [Dahlhaus and Warmuth 1986]
A grammar $(N, \Sigma, P, S)$ is growing context-sensitive iff $|l| < |r|$ for all rules $l \to r$ in the grammar.

Theorem [Dahlhaus and Warmuth 1986]
Given a growing context-senstive grammar there is a polynomial time algorithm which decides membership of the language generated by this growing grammar.

In this result, polynomial time means here only with respect to the words in the language, the dependence on the size of the grammar is not polynomial time. So if one asks the uniform decision problem for an input consisting of a pair of a grammar and a word, no polynomial time algorithm is known for this problem. As the problem is NP-complete, the algorithm is unlikely to exist.

# Example 8.28

Consider the grammar

$$(\{S, T, U\}, \{0, 1\}, \{S \rightarrow 011 | T11, T \rightarrow T0U | 00U, U0 \rightarrow 0UU, U1 \rightarrow 111\}, S)$$

which is growing. This grammar has derivations like $S \Rightarrow T11 \Rightarrow 00U11 \Rightarrow 001111$ and $S \Rightarrow T11 \Rightarrow T0U11 \Rightarrow 00U0U11 \Rightarrow 00U01111 \Rightarrow 000UU1111 \Rightarrow 000U111111 \Rightarrow 00011111111$. The language of the grammar is

$$\{0^n 1^{2^n} : n > 0\} = \{011, 001111, 0^3 1^8, 0^4 1^{16}, 0^5 1^{32}, \ldots\}$$

and not context-free, as infinite languages satisfying the context-free pumping lemma can only have constant gaps (sequence of lengths without a word). This grammar has growing gaps.

# Exercises 8.29 - 8.31

## Exercise 8.29

Show that every context-free language is the union of a language generated by a growing grammar and a language containing only words up to length $1$.

## Exercise 8.30

Modify the proof of Theorem 6.11 to prove that every recursively enumerable language, that is, every language generated by some grammar is the homomorphic image of a language generated by a growing context-sensitive grammar.

## Exercise 8.31

Construct a growing grammar for the language $\{1^{2^n}0^{2n}1^{2^n} : n > 0\}$ which is the "palindromisation" of the language from Example 8.28.