

THEORY OF COMPUTATION

Grammars - 25

Prof. Dan A. Simovici

UMB

1 Grammars

2 Languages Generated by Grammars

3 Unsolvable Problems Concerning Grammars

Definition

A **grammar** is a semi-Thue process that involves two types of symbols:

- 1 **nonterminal symbols** or **variables** denoted by capital letters, X, Y, Z, S, \dots , and
- 2 **terminal symbols** or **terminals** denoted by small letters, a, b, c, \dots

A special nonterminal symbol S is the **start symbol**.

In addition, for every production $x \rightarrow y$ the left part contains a nonterminal symbol.

A grammar will be denoted as

$$\Gamma = (\mathcal{V}, T, S, P),$$

where

- \mathcal{V} is the set of non-terminals or variables;
- T is the set of terminals;
- $S \in \mathcal{V}$ is the **start symbol**, and
- P is the set of productions.

Definition

The **language generated by Γ** is the set $L(\Gamma) \subseteq T^*$ given by

$$L(\Gamma) = \{u \in T^* \mid S \xrightarrow[\Gamma]{*} u\}.$$

Note that in a grammar **all non-terminal symbols are eliminated in the derivation process** that ends up with a word over the terminal alphabet.

Example

Let $\Gamma = (\{S, X, Y\}, \{a, b\}, S, \{S \rightarrow X, X \rightarrow aX, X \rightarrow 0, X \rightarrow Y, Y \rightarrow bY, Y \rightarrow 0\})$.

Every derivation in Γ that begins with S and ends with a word in T^* has the form

$$\begin{aligned}
 S &\xRightarrow{\Gamma} X \xRightarrow{\Gamma} aX \xRightarrow{\Gamma} aaX \\
 &\xRightarrow{\Gamma} aaaX \xRightarrow{\Gamma} aaaY \xRightarrow{\Gamma} aaabY \\
 &\xRightarrow{\Gamma} aaabbY \xRightarrow{\Gamma} aaabb.
 \end{aligned}$$

Thus, the language $L(\Gamma)$ is $\{a^n b^m \mid n, m \in \mathbb{N}\}$.

Example

Let $\Gamma = (\{S\}, \{a, b\}, S, \{S \rightarrow aSb, S \rightarrow 0\})$.

Every derivation in Γ that begins with S and ends with a word in T^* has the form

$$\begin{aligned} S &\xRightarrow{\Gamma} aSb \xRightarrow{\Gamma} aaSbb \xRightarrow{\Gamma} aaaSbbb \\ &\xRightarrow{\Gamma} aaabbb. \end{aligned}$$

The language generated by this grammar is

$$L(\Gamma) = \{a^n b^n \mid n \in \mathbb{N}\}.$$

Example

Consider the grammar

$$\Gamma = (\{S, X, Y\}, \{a, b, c\}, S, P),$$

where P consists of the following productions:

$$\begin{aligned} \pi_0 & : S \rightarrow abc, & \pi_1 & : S \rightarrow aXbc, \\ \pi_2 & : Xb \rightarrow bX, & \pi_3 & : Xc \rightarrow Ybcc, \\ \pi_4 & : bY \rightarrow Yb, & \pi_5 & : aY \rightarrow aaX, \\ \pi_6 & : aY \rightarrow aa \end{aligned}$$

We will refer later in this lecture to this kind of grammars as **length-increasing grammars** because for each of its productions $x \rightarrow y$ we have $|x| \leq |y|$.

Example cont'd

We claim that $L(\Gamma) = \{a^n b^n c^n \mid n \in \mathbb{P}\}$.

- Any word $\alpha \in \{S, X, Y, a, b, c\}^*$ that occurs in a derivation, $S \xRightarrow{*} \alpha$ contains at most one nonterminal symbol.
- A derivation must end either by applying the production $S \rightarrow abc$ or the production $aY \rightarrow aa$ because only these productions allow us to eliminate a nonterminal symbol.
- If the last production is $S \rightarrow abc$, then the derivation is $S \Rightarrow abc$, and the derived word has the form prescribed.

Otherwise, the symbol Y must be generated starting from S , and the first production applied is $S \rightarrow aXbc$.

Example cont'd

Note that for every $i \geq 1$ we have

$$a^i X b^i c^i \xrightarrow[\Gamma]{*} a^{i+1} X b^{i+1} c^{i+1}.$$

Indeed, we can write:

$$\begin{array}{ccc}
 a^i X b^i c^i & \xrightarrow[\pi_2]{i} & a^i b^i X c^i & \xrightarrow[\pi_3]{1} & a^i b^i Y b c^{i+1} \\
 & \xrightarrow[\pi_4]{i} & a^i Y b^{i+1} c^{i+1} & \xrightarrow[\pi_5]{1} & a^{i+1} X b^{i+1} c^{i+1}
 \end{array}$$

We claim that a word α contains the infix aY (which allows us to apply the production π_5) and $S \xrightarrow[\Gamma]{*} \alpha$ if and only if α has the form $\alpha = a^i Y b^{i+1} c^{i+1}$ for some $i \geq 1$.

Example cont'd

An easy argument by induction on $i \geq 1$ allows us to show that if $\alpha = a^i Y b^{i+1} c^{i+1}$ then $S \xRightarrow{\Gamma}^* \alpha$. We need to prove only the inverse implication. This can be done by strong induction on the length $n \geq 3$ of the derivation $S \xRightarrow{\Gamma}^* \alpha$.

The shortest derivation that allows us to generate the word containing the infix aY is

$$S \xRightarrow{\Gamma} aXbc \xRightarrow{\Gamma} abXc \xRightarrow{\Gamma} abYbcc \xRightarrow{\Gamma} aYb^2c^2,$$

and this word has the prescribed form.

Example cont'd

Suppose now that for derivations shorter than n the condition is satisfied, and let $S \xrightarrow[G]{*} \alpha$ be a derivation of length n such that α contains the infix aY . By the inductive hypothesis the previous word in this derivation that contains the infix aY has the form $\alpha' = a^j Y b^{j+1} c^{j+1}$. To proceed from α' we must apply the production π_5 and replace Y by X . Thus, we have

$$S \xrightarrow[G]{*} a^j Y b^{j+1} c^{j+1} \Rightarrow_G a^{j+1} X b^{j+1} c^{j+1}.$$

Example cont'd

Next, the symbol X must “travel” to the right using the production π_2 , transform itself into an Y (when in touch with the c s) and Y must “travel” to the left to create the infix aY . This can happen only through the application of the productions π_3 and π_4 , as follows:

$$\begin{array}{ccc}
 a^{j+1}Xb^{j+1}c^{j+1} & \xRightarrow[\pi_2]{j+1} & a^{j+1}b^{j+1}Xc^{j+1} \\
 & \xRightarrow[\pi_3]{1} & a^{j+1}b^{j+1}Ybc^{j+2} \\
 & \xRightarrow[\pi_4]{i} & a^{j+1}Yb^{j+2}c^{j+2},
 \end{array}$$

which proves that α has the desired form. Therefore, all the words in the language $L(\Gamma)$ have the form $a^n b^n c^n$.

Theorem

Let U be a language accepted by a nondeterministic Turing machine \mathcal{M} . Then, there is a grammar Γ such that $U = L(\Gamma)$

Proof

Recall that we defined a semi-Thue process $\Omega(\mathcal{M})$ attached to the TM \mathcal{M} .

We started from \mathcal{M} and defined first the semi-Thue system $\Sigma(\mathcal{M})$ on the alphabet

$$s_0, s_1, \dots, s_K, q_0, q_1, \dots, q_n, q_{n+1}, h$$

containing the following productions:

Quadruple	semi-Thue Production
$q_i s_j s_k q_\ell$	$q_i s_j \rightarrow q_\ell s_k$
$q_i s_j R q_\ell$	$q_i s_j s_k \rightarrow s_j q_\ell s_k, 0 \leq k \leq K$ $q_i s_j h \rightarrow s_j q_\ell s_0 h$
$q_i s_j L q_\ell$	$q_\ell s_k s_j \rightarrow s_0 q_\ell s_k, 0 \leq k \leq K$ $h q_i s_j \rightarrow h q_\ell s_0 s_j$

Proof cont'd

In addition we included in $\Sigma(\mathcal{M})$ the following productions:

- whenever $q_i s_j$ are not the first two symbols of a quadruple of \mathcal{M} we place in $\Sigma(\mathcal{M})$ the production $q_i s_j \rightarrow q_{n+1} s_j$. Thus, q_{n+1} serves as “halt” state.
- Finally, we place in $\Sigma(\mathcal{M})$ the productions:

$$q_{n+1} s_i \rightarrow q_{n+1}, 0 \leq i \leq K,$$

$$q_{n+1} h \rightarrow q_0 h,$$

$$s_i q_0 \rightarrow q_0, 0 \leq i \leq K.$$

Proof cont'd

The system $\Omega(\mathcal{M})$ contains the productions

Quadruple	semi-Thue Production
$q_i s_j s_k q_\ell$	$q_\ell s_k \rightarrow q_i s_j$
$q_i s_j R q_\ell$	$s_j q_\ell s_k \rightarrow q_i s_j s_k, 0 \leq k \leq K$ $s_j q_\ell s_0 h \rightarrow q_i s_j h$
$q_i s_j L q_\ell$	$s_0 q_\ell s_k \rightarrow q_\ell s_k s_j, 0 \leq k \leq K$ $h q_\ell s_0 s_j \rightarrow h q_i s_j$

Proof cont'd

In addition, we have in $\Omega(\mathcal{M})$:

- $q_{n+1}s_j \rightarrow q_i s_j$ when $q_i s_j$ are not the first two symbols of a quadruple of \mathcal{M} , and

-

$$q_{n+1} \rightarrow q_{n+1}s_i, 0 \leq i \leq K,$$

$$q_0 h \rightarrow q_{n+1} h,$$

$$q_0 \rightarrow s_i q_0, 0 \leq i \leq K.$$

Proof cont'd

We construct the grammar Γ by modifying the semi-Thue process $\Omega(\mathcal{M})$ as follows:

- the terminals of Γ are just the letters of the alphabet $T = \{s_1, \dots, s_m\}$ of \mathcal{M} ;
- the non-terminals (variables) of Γ are the symbols of $\Omega(\mathcal{M})$ not in T , $s_0, q_0, \dots, q_n, q_{n+1}, h$;
- two additional non-terminals S and q .

S is the start symbol of Γ .

Proof cont'd

The production of Γ are:

- the productions of $\Omega(\mathcal{M})$;
- $S \rightarrow hq_0h$;
- $hq_1s_0 \rightarrow q$;
- $qs \rightarrow sq$ for each $s \in T$;
- $qh \rightarrow 0$.

Proof cont'd

Suppose \mathcal{M} accepts $u \in T^*$, that is:

$$S \xRightarrow{\Gamma} hq_0h \xRightarrow{\Gamma}^* hq_1s_0uh \xRightarrow{\Gamma} quh \xRightarrow{\Gamma}^* uqh \xRightarrow{\Gamma} u,$$

so that $u \in L(\Gamma)$.

Proof cont'd

Conversely, let $u \in L(\Gamma)$. Then, $u \in T^*$ and $S \xRightarrow{\Gamma}^* u$. Examining the list of productions of Γ this derivation can be written as

$$S \xRightarrow{\Gamma} hq_0h \xRightarrow{\Gamma}^* vqhz \xRightarrow{\Gamma} vz = u.$$

Note that q could be introduced only by using the production $hq_1s_0 \rightarrow q$. Thus, the derivation has the form

$$S \xRightarrow{\Gamma} hq_0h \xRightarrow{\Gamma}^* xhq_1s_0yhz \xRightarrow{\Gamma} xqyhz \xRightarrow{\Gamma}^* xyqhz \xRightarrow{\Gamma} xyz = u,$$

where $xy = v$. Thus, there is a derivation of xhq_1s_0yhz from hq_0h in Γ . This derivation must actually be a derivation in $\Omega(\mathcal{M})$ because the added productions are inapplicable.

Proof cont'd

The productions in $\Omega(\mathcal{M})$ always lead from Post words to Post words, hence xhq_1s_0yhz must be a Post word, which implies $x = z = 0$ and $u = xyz = y$. We conclude that

$$hq_0h \xRightarrow[\Omega(\mathcal{M})]{*} hq_1s_0uh,$$

which implies that \mathcal{M} accepts u .

Let Γ be a grammar having the alphabet

$$\{s_1, \dots, s_n, V_1, \dots, V_k\},$$

where $T = \{s_1, \dots, s_n\}$ is the set of terminals and $\{V_1, \dots, V_k\}$ is the set of variables (nonterminals). We assume that $S = V_1$ is the start symbol.

Assume that the alphabet of Γ is ordered as above and we regard strings on this alphabet as integers in the base $n + k$.

Theorem

The predicate $u \xRightarrow{\Gamma} v$ is primitive recursive.

Proof.

Let the production of Γ be $x_i \rightarrow y_i$ for $1 \leq i \leq \ell$. For $1 \leq i \leq \ell$ define the predicate $\text{PROD}_i(u, v)$ as

$$(\exists r, s)_{\leq u} [u = \text{CONCAT}(r, x_i, s) \& v = \text{CONCAT}(r, y_i, s)]$$

Since CONCAT is primitive recursive, PROD_i is primitive recursive. Since $u \xRightarrow{\Gamma} v$ if and only if

$$\text{PROD}_1(u, v) \vee \text{PROD}_2(u, v) \vee \dots \vee \text{PROD}_\ell(u, v)$$

the result follows. □

Define the predicate $\text{DERIV}(u, y)$ to mean that for some m
 $y = [u_1, \dots, u_m, 1]$, where u_1, \dots, u_m is a derivation of u from S in Γ , that is,

$$S = u_1 \xRightarrow{\Gamma} u_2 \xRightarrow{\Gamma} \cdots \xRightarrow{\Gamma} u_m = u.$$

u_1 has been added to avoid complications when $u_m = u = 0$.
 Note that the value of S in the base $n + k$ is $n + 1$ (because $S = V_1$ is the $(n + 1)^{\text{st}}$ symbol in the alphabetic list).

Theorem

The predicate $DERIV(u, y)$ is primitive recursive.

Proof.

This follows from the following equivalent statements:

$$\begin{aligned} DERIV(u, y) &\Leftrightarrow (\exists m)_{\leq y} (m + 1 = Lt(y) \\ &\quad \& (y)_1 = n + 1 \& (y)_m = u \& (y)_{m+1} = 1 \\ &\quad \& (\forall j)_{< m} (j = 0 \vee [(y)_j \xrightarrow{r} (y)_{j+1}]) \end{aligned}$$



Note that

- By the definition of $\text{DERIV}(u, y)$ we have

$$S \xRightarrow{\Gamma}^* u \text{ if and only if } (\exists y)\text{DERIV}(u, y).$$

- $S \xRightarrow{\Gamma}^* u$ if and only if $\min_y \text{DERIV}(u, y) \downarrow$.

Therefore, $\{u \mid S \xRightarrow{\Gamma}^* u\}$ is recursively enumerable. Since

$L(\Gamma) = T^* \cap \{u \mid S \xRightarrow{\Gamma}^* u\}$ it follows that $L(\Gamma)$ is r.e.

Corollary

A language U is r.e. if and only if there is a grammar Γ such that $U = L(\Gamma)$.

Putting together previous results we have the following

Theorem

The following are equivalent for a language L :

- 1** *L is r.e.;*
- 2** *L is accepted by a deterministic TM;*
- 3** *L is accepted by a nondeterministic TM;*
- 4** *there is a grammar Γ such that $L = L(\Gamma)$.*

Definition

A grammar Γ is called *length-increasing* if for every production $x \rightarrow y$ we have $|x| \leq |y|$.

An equivalent class of grammars to the class of length-increasing grammars is the class of *context-sensitive grammars*. This equivalence is a topic in the theory of formal languages.

Theorem

If Γ is a length-increasing grammar, then the set $\{u \in (\mathcal{V} \cup \mathcal{T})^ \mid S \xRightarrow{\Gamma}^* u\}$ is recursive.*

Proof

Recall that we have shown that

$$S \xrightarrow[\Gamma]{*} u \text{ if and only if } \min_y \text{DERIV}(u, y) \downarrow$$

It will suffice to obtain a recursive bound for y to establish that $L(\Gamma)$ is recursive.

Note that in every derivation in Γ we have

$$1 = |u_1| \leq |u_2| \leq \dots \leq |u_m| = |u|.$$

Therefore, $u_1, u_2, \dots, u_m = u \leq g(u)$, where $g(u)$ is the smallest number that represents a string of length $|u| + 1$ in the base $n + k$.

Proof cont'd

Note that:

- $g(u)$ is the value in the base $n + k$ of a string consisting of $|u| + 1$ repetitions of 1, so $g(u) = \sum_{i=0}^{|u|} (n + k)^i$, which is primitive recursive because $|u|$ is primitive recursive.
- We may assume that the derivation

$$S = u_1 \xRightarrow{\Gamma} u_2 \xRightarrow{\Gamma} \cdots \xRightarrow{\Gamma} u_m = u$$

contains no repetitions because given a sequence of steps

$$z = u_i \xRightarrow{\Gamma} u_{i+1} \xRightarrow{\Gamma} \cdots \xRightarrow{\Gamma} u_{i+l} = z$$

we could eliminate the steps u_{i+1}, \dots, u_l .

Thus, the length of the derivation is bounded by the total number of strings of length less or equal to $|u|$ on the alphabet with $n + k$ symbols, which is just the number $g(u)$.

Hence,

$$[u_1, \dots, u_m, 1] = \prod_{i=1}^m p_i^{u_i} \cdot p_{m+1} \leq h(u),$$

where

$$h(u) = \prod_{i=1}^{g(u)} p_i^{g(u)} \cdot p_{g(u)+1}.$$

Finally, we have $S \xrightarrow{\Gamma}^* u$ if and only if $(\exists y)_{\leq h(u)} \text{DERIV}(u, y)$, which gives the result.

Theorem

If Γ is a length-increasing grammar, then $L(\Gamma)$ is recursive.

Proof.

By the previous theorem, the set $\{u \in (\mathcal{V} \cup \mathcal{T})^* \mid S \xrightarrow{\Gamma}^* u\}$ is recursive. Since

$$L(\Gamma) = \{u \in (\mathcal{V} \cup \mathcal{T})^* \mid S \xrightarrow{\Gamma}^* u\} \cap \mathcal{T}^*,$$

and \mathcal{T}^* is recursive, it follows that $L(\Gamma)$ is recursive. □

Let \mathcal{M} be a TM and let u be a word in the alphabet of \mathcal{M} . The grammar Γ_u is constructed as follows:

- The variables of Γ_u are the entire alphabet of $\Sigma(\mathcal{M})$ together with S (the start symbol) and a new nonterminal symbol V . There is just one terminal symbol a .
- The production of Γ_u are all productions of $\Sigma(\mathcal{M})$ together with

$$S \rightarrow hq_1s_0uh, hq_0h \rightarrow V, V \rightarrow aV, V \rightarrow a$$

We have $S \xrightarrow[\Gamma_u]{*} V$ if and only if \mathcal{M} accepts u .

Lemma

If \mathcal{M} accepts u , then $L(\Gamma_u) = \{a^i \mid i \neq 0\}$; if \mathcal{M} does not accept u , then $L(\Gamma_u) = \emptyset$.

Proof.

The fact that \mathcal{M} accepts u means that:

$$S \xRightarrow[\Gamma_u]^* hq_1s_0uh \xRightarrow{\Gamma_u} hq_0h \xRightarrow{\Gamma_u} V \xRightarrow[\Gamma_u]^* a^{n-1}V \xRightarrow{\Gamma_u} a^n,$$

If \mathcal{M} does not accept u , then the word hq_0u cannot be generated, so $L(\Gamma_u) = \emptyset$. □

Select \mathcal{M} such that the language accepted by it is not recursive. Then, there is no algorithm for determining for given u whether \mathcal{M} accepts u . The lemma implies that

$$\begin{aligned}\mathcal{M} \text{ accepts } u &\Leftrightarrow L(\Gamma_u) \neq \emptyset \\ &\Leftrightarrow L(\Gamma_u) \text{ is infinite} \\ &\Leftrightarrow a \in L(\Gamma_u).\end{aligned}$$

The above prove the following:

Theorem

There is no algorithm to determine of a given grammar Γ whether

- 1** $L(\Gamma)$ is empty;
- 2** $L(\Gamma)$ is infinite;
- 3** $v_0 \in L(\Gamma)$ for a fixed word v_0 .

Theorem

There is no algorithm for determining of a given pair of grammars Γ_1 and Γ_2 whether

- 1 $L(\Gamma_1) \subseteq L(\Gamma_2)$;
- 2 $L(\Gamma_1) = L(\Gamma_2)$.

Proof

Let Γ_1 be the grammar whose productions are

$$S \rightarrow aS, S \rightarrow a$$

We have $L(\Gamma_1) = \{a^i \mid i \neq 0\}$. Thus, by the previous theorem, \mathcal{M} accepts u if and only if $L(\Gamma_1) = L(\Gamma_u)$ if and only if $L(\Gamma_1) \subseteq L(\Gamma_u)$.