

Introduction to Tree Language Theory

Hitoshi Ohsaki



National Institute of
Advanced Industrial Science and Technology (AIST)

seminar talk (2/10)

2009

II. Grammar

Grammar

grammar : $\mathcal{G} = (\Sigma, T, N, q_0, \Delta)$

Σ : alphabet

T : set of **terminal symbols** such that $T \subseteq \Sigma$

N : set of **non-terminal symbols** such that $N = \Sigma - T$

q_0 : start symbol such that $q_0 \in N$

Δ : finite set of production rules with the following forms

$$\alpha \rightarrow \beta \quad (\alpha, \beta \in \Sigma^*)$$

Cf. **regular grammar** if production rules are in the following forms:

$$p \rightarrow a \quad p \rightarrow aq \quad p \rightarrow \varepsilon \quad (p, q \in N, a \in T)$$

Generated languages

Given grammar $\mathcal{G} = (\Sigma, T, N, q_0, \Delta)$, define for words u, w over Σ ,

$u \rightarrow_{\mathcal{G}} w$: $\exists \alpha \rightarrow \beta$ in Δ such that

- u is decomposed to words u_1, α, u_2
- w is decomposed to words u_1, β, u_2

$u_0 \rightarrow_{\mathcal{G}}^* u_n$: if $u_0 \equiv u_n$, or

if $u_0 \rightarrow_{\mathcal{G}} u_1$ & $u_1 \rightarrow_{\mathcal{G}}^* u_n$

(the reflexive and transitive closure of $\rightarrow_{\mathcal{G}}$)

$\mathcal{L}(\mathcal{G})$: set of words over T such that for all w in $\mathcal{L}(\mathcal{G})$, $q_0 \rightarrow_{\mathcal{G}}^* w$
(language generated by \mathcal{G})

Note

Question if $w \in \mathcal{L}(\mathcal{G})$ is **undecidable** in general. (Cf. Rice's theorem in seminar 5)

Decidable sub-classes

grammar whose production rules are in the following forms is called

context-sensitive : $\alpha \rightarrow \beta$ ($\alpha, \beta \in \Sigma^*$ such that $|\alpha| \leq |\beta|$,
 $q_0 \rightarrow \varepsilon$ start symbol q_0 never appears
on right-hand side of any rule)

context-free : $p \rightarrow \beta$ ($p \in N, \beta \in \Sigma^*$)

language is called

context-sensitive if it is generated by context-sensitive grammar

context-free if it is generated by context-free grammar

Note

Membership problem, i.e. the question if $w \in \mathcal{L}(\mathcal{G})$, is **decidable** for the classes of context-sensitive grammar and context-free grammar. Why? (**Exercise**)

Example

Define the production rules

$$\Delta_1 : q_0 \rightarrow a q_0 b \quad q_0 \rightarrow \varepsilon$$

$$\begin{aligned} \Delta_2 : q_0 &\rightarrow a q_1 & q_0 &\rightarrow b q_2 & q_0 &\rightarrow \varepsilon \\ q_1 &\rightarrow a q_1 q_1 & q_1 &\rightarrow b q_0 \\ q_2 &\rightarrow a q_0 & q_2 &\rightarrow b q_2 q_2 \end{aligned}$$

for the grammar

$$\mathcal{G}_1 = (\Sigma_1, \{a, b\}, \{q_0\}, q_0, \Delta_1)$$

$$\mathcal{G}_2 = (\Sigma_2, \{a, b\}, \{q_0, q_1, q_2\}, q_0, \Delta_2)$$

then

$$\mathcal{L}(\mathcal{G}_1) = \{a^n b^n \mid n \geq 0\} \quad \mathcal{L}(\mathcal{G}_2) = \{w \in \{a, b\}^* \mid |w|_a = |w|_b\}$$

where $|w|_a$ (resp. $|w|_b$) is the number of occurrences of a (b) in w 5

Eliminating ε -derivations

Given grammar $\mathcal{G} = (\Sigma, T, N, q_0, \Delta)$, for non-terminal symbol q in N , q is called **nullable** if q admits $q \rightarrow_{\mathcal{G}}^* \varepsilon$.

Claim 1

One can compute the set N_ε of nullable non-terminal symbols if \mathcal{G} is context-free grammar

Claim 2

If \mathcal{G} is context-free grammar, define $\mathcal{G}' = (\Sigma, T, N, q_0, \Delta')$ where

$$\begin{aligned} \Delta' = \{ & q \rightarrow \alpha_0 \beta_1 \alpha_1 \cdots \beta_n \alpha_n \mid \exists p \rightarrow \alpha_0 p_1 \alpha_1 \cdots p_n \alpha_n \in \Delta : \\ & \exists \alpha_i \in (T \cup N - N_\varepsilon)^*, \exists p_1, \dots, p_n \in N_\varepsilon, \beta_i \in \{p_i, \varepsilon\} \\ & (1 \leq i \leq n) \} \end{aligned}$$

Then

$$\mathcal{L}(\mathcal{G}') = \mathcal{L}(\mathcal{G}) - \{\varepsilon\}$$

... show this claim (**Exercise**)

Proof of Claim 1

We show that for each $q \in N$, q is nullable if and only if $q \in N_\varepsilon$, where N_ε is obtained by the following procedure. This procedure halts on any context-free grammar.

$N_\varepsilon := \emptyset$

while $S = \emptyset$ do

$S := \{q \in N - N_\varepsilon \mid \exists q \rightarrow \alpha \in \Delta \text{ such that } \alpha \in N_\varepsilon^*\};$

$N_\varepsilon := N_\varepsilon \cup S$

od

return N_ε

The “if” part is shown by induction on the number n of loops. For induction step, let S_n be the set S obtained at the n -th loop ($n \geq 1$). For $q \in S_n$, there is $q \rightarrow \alpha$ with $\alpha \in N_\varepsilon^*$. By induction hypothesis, non-terminals of α are all nullable, so $\alpha \xrightarrow{*}_G \varepsilon$, and hence $q \xrightarrow{*}_G \varepsilon$.

The “only if” part is shown by the length of $q \xrightarrow{*}_G \varepsilon$. If the length is 1, Δ contains $q \rightarrow \varepsilon$, and thus, q is added to N_ε at 1st loop. For induction step, suppose $q \xrightarrow{*}_G \varepsilon$ whose length is n ($n \geq 1$) and $\alpha \in N_\varepsilon^*$. By assumption, Δ contains $q \rightarrow \alpha$. If $q \notin N_\varepsilon$, $q \in S$ at some loop, because it satisfies that $q \in N - N_\varepsilon$, $q \rightarrow \alpha \in \Delta$ and $\alpha \in N_\varepsilon^*$. Hence, $q \in N_\varepsilon$. \square

Proposition

For context-free grammar \mathcal{G} with the set T of terminal symbols, one can construct context-free grammar $\mathcal{G}' = (\Sigma, T, N, q_0, \Delta')$ such that $\mathcal{L}(\mathcal{G}') = \mathcal{L}(\mathcal{G})$ and Δ' contains transition rules with the following forms only :

$$q \rightarrow \alpha \quad (q \in N, \alpha \in (\Sigma - \{q_0\})^+),$$

$$q_0 \rightarrow \varepsilon \quad (q_0 \text{ never appears in the right-hand side of any rule})$$

Proof

This is an easy consequence of the previous **Claims 1,2** except the treatment of the start symbol and the production rule of the form $q_0 \rightarrow \varepsilon$. First, replace the start symbol, say p_0 , of \mathcal{G} by q_0 and add $q_0 \rightarrow p_0$ to the set of production rules of \mathcal{G} . After constructing \mathcal{G}' , test if \mathcal{G} generates ε , which is the decidable problem for context-free grammar. If yes, then add $q_0 \rightarrow \varepsilon$ to Δ' of \mathcal{G}' . \square

Corollary

Context-free languages are generated by context-sensitive grammar. (This is not obvious only from grammar definitions on page 4, but from the above Proposition.) δ

Pumping lemma [Bar-Hillel 1961]

Given context-free language L ,

$\exists k \geq 0$: if $z \in L$ and $|z| \geq k$, then z is formed by u, v, w, x, y

as $\overline{u} \quad \overline{v} \quad \overline{w} \quad \overline{x} \quad \overline{y}$

such that $|vx| \geq 1$ and $uv^nwx^ny \in L$ ($n \geq 0$)

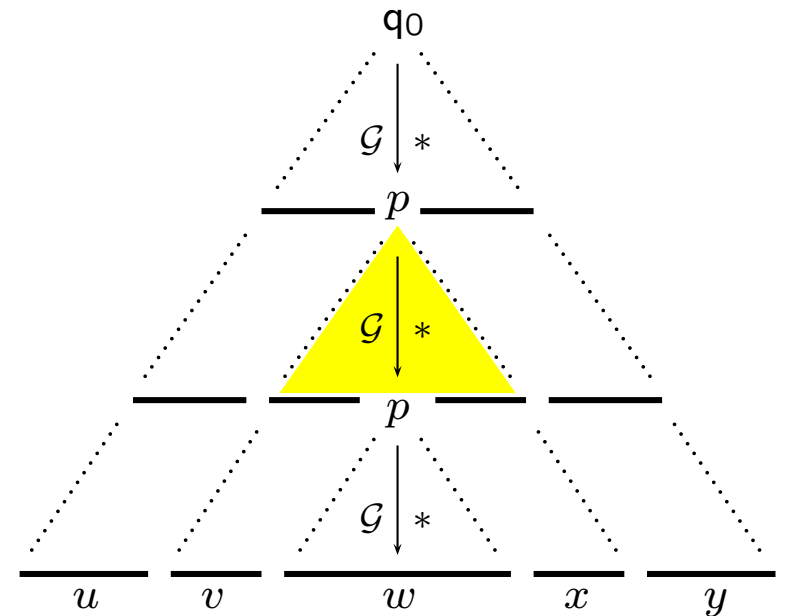
as $\overline{u} \quad \overbrace{\overline{v} \cdots \overline{v}}^n \quad \overline{w} \quad \overbrace{\overline{x} \cdots \overline{x}}^n \quad \overline{y}$

Proof

Suppose $\mathcal{G} = (\Sigma, T, N, q_0, \Delta)$ is context-free grammar whose production rules are already in the forms of $q_0 \rightarrow \varepsilon$ and $q \rightarrow \alpha$ for $\alpha \in (\Sigma - \{q_0\})^+$. For $p, q \in \Sigma$, we say p is a **descendant** of q if $q \rightarrow_{\mathcal{G}}^* upw$. If the length of the derivation $q \rightarrow_{\mathcal{G}}^* upw$ is more than 0, p is a **proper descendant** of q . Observe that the (proper) descendant relation is closed under contexts, i.e. p is a descendant of q such that $q \rightarrow_{\mathcal{G}}^* upw$ if and only if for all s, t in Σ^* , $sq t \rightarrow_{\mathcal{G}}^* supwt$. For every word z generated by \mathcal{G} , one can construct a directed graph of the proper descendant relation associated to z , whose node is marked by an element in Σ and the edge is the proper descendant relation obtained from the derivation $q_0 \rightarrow_{\mathcal{G}}^* z$. (proof cont'd)

Proof (cont'd)

If \mathcal{G} generates a word z whose directed graph of the proper descendant relation admits the path starting from q_0 to a in T such that the length of this path is more than $|N|$, then the path contains a sub-path from p to p . (This situation corresponds to the figure on the right.) If the yellow part is removed, $q_0 \rightarrow_{\mathcal{G}}^* u w y$; if this part is repeated n -times, it turns out $q_0 \rightarrow_{\mathcal{G}}^* u v^n w x^n y$. Since p is the proper descendant of p in this case, by the shape of rules, $|v x| \geq 1$. \square



Proposition 1

Given context-free grammar \mathcal{G} , question if $\mathcal{L}(\mathcal{G}) = \emptyset$ is decidable

Proof

Suppose $\mathcal{G} = (\Sigma, T, N, q_0, \Delta)$, each of whose rules is already in the form of $q_0 \rightarrow \varepsilon$ or $q \rightarrow \alpha$ for $q \in N$, $\alpha \in (\Sigma - \{q_0\})^+$. Let $\ell = \max\{|\alpha| \mid p \rightarrow \alpha \in \Delta\}$ and $k = \ell^{|N|+1}$, then define $L_{\leq k} = \{w \in T^* \mid |w| \leq k\}$. From the proof of the previous lemma, $L_{\leq k} = \emptyset$ if and only if $\mathcal{L}(\mathcal{G}) = \emptyset$. Because $L_{\leq k}$ is finite, question if $L_{\leq k} = \emptyset$ is decidable. \square 10

Example

Using **Pumping Lemma**, we show that $L = \{a^l b^l a^l \mid l \geq 1\}$ is **not** context-free. For leading the contradiction, we suppose below that L is context-free.

Take $z = a^k b^k a^k$ in L . If k is sufficiently large, z can be decomposed to u, v, w, x, y such that $|vx| \geq 1$ and $uv^n wx^n y \in L$ for all $n \geq 0$.

If $|v| = 0$, then $x = a^i$ or $x = b^i$ with $i \leq k$. However, $uv^2 wx^2 y$ can not be an element in L . Thus, $|v| \neq 0$.

If $|v| \geq 1$, then $v = a^i$ or $v = b^i$ with $i \leq k$. If $|x| = 0$, then it immediately leads to the contradiction, so $|x| \neq 0$.

If $|x| \geq 1$ also, then $x = a^j$ or $x = b^j$ with $j \leq k$. However, in any combination of the possibility of v and x , $uv^2 wx^2 y$ can not be an element in L , leading to the contradiction.

Proposition 2

Given two context-free grammars $\mathcal{G}_1, \mathcal{G}_2$, the intersection-emptiness problem, i.e. the question if $\mathcal{L}(\mathcal{G}_1) \cap \mathcal{L}(\mathcal{G}_2) = \emptyset$, is **undecidable**.

Proof

Use the reduction from PCP. For an instance of PCP $I = \{ \langle u_1, w_1 \rangle, \dots, \langle u_n, w_n \rangle \}$ over Σ , define the grammar $\mathcal{G}_I = (\Sigma \cup \{q_0\}, \Sigma, \{q_0\}, q_0, \Delta_I)$, provided that $q_0 \notin \Sigma$:

$$\Delta_I : \quad q_0 \rightarrow \bar{u}_1 q_0 w_1 \quad \dots \quad q_0 \rightarrow \bar{u}_n q_0 w_n \quad q_0 \rightarrow \varepsilon$$

where \bar{u}_i denotes the reverse of word u_i , e.g. $\overline{abcd} = dcba$. Moreover, if $\Sigma = \{a_1, \dots, a_k\}$, define the grammar $\mathcal{G}_{\text{eq}} = (\Sigma \cup \{q_0, q_1\}, \Sigma, \{q_0, q_1\}, q_0, \Delta_{\text{eq}})$, where

$$\begin{aligned} \Delta_{\text{eq}} : \quad & q_0 \rightarrow a_1 q_1 a_1 \quad \dots \quad q_0 \rightarrow a_k q_1 a_k \\ & q_1 \rightarrow a_1 q_1 a_1 \quad \dots \quad q_1 \rightarrow a_k q_1 a_k \quad q_1 \rightarrow \varepsilon \end{aligned}$$

Note that $\mathcal{L}(\mathcal{G}_{\text{eq}})$ does not contain the empty word ε . Then, by construction, $\mathcal{L}(\mathcal{G}_I) \cap \mathcal{L}(\mathcal{G}_{\text{eq}}) \neq \emptyset$ if and only if I has a solution. However, question if an instance of PCP has a solution is undecidable. □ 12

Closure properties (1)

Let $C(CF_T)$ be the set of context-free languages over T of terminal symbols

Proposition

The class $C(CF_T)$ is closed under union. However, the class $C(CF_T)$ is **not** closed under intersection or complement if $|T| \geq 2$

Proof

For union, let $\mathcal{G}_1 = (\Sigma_1, T, N_1, p_0, \Delta_1)$ $\mathcal{G}_2 = (\Sigma_2, T, N_2, q_0, \Delta_2)$ be context-free grammar. We suppose $N_1 \cap N_2 = \emptyset$. Let r_0 be fresh symbol, then define $\mathcal{G} = (\Sigma_1 \cup \Sigma_2 \cup \{r_0\}, T, N_1 \cup N_2 \cup \{r_0\}, r_0, \Delta_1 \cup \Delta_2 \cup \{r_0 \rightarrow p_0, r_0 \rightarrow q_0\})$. By construction, for word $w \in T^*$, \mathcal{G} generates w if and only if \mathcal{G}_1 or \mathcal{G}_2 generates w .

For not being closed under intersection, it suffices to show that $L_1 = \{a^\ell b^\ell a^m \mid \ell, m \geq 1\}$ and $L_2 = \{a^m b^\ell a^\ell \mid \ell, m \geq 1\}$ are generated by context-free grammar.

For not being closed under complement, use **de Morgan's law**, because $L_1 \cap L_2 = ((L_1)^c \cup (L_2)^c)^c$. \square

Closure properties (2)

Let $C(CS_T)$ be the set of context-sensitive languages over T of terminal symbols

Proposition

The class $C(CS_T)$ is closed under union, intersection, complement

Proof

For union, one can apply the same proof of the previous proposition.

For intersection, let $\mathcal{G}_1 = (\Sigma_1, T, N_1, p_0, \Delta_1)$ and $\mathcal{G}_2 = (\Sigma_2, T, N_2, q_0, \Delta_2)$ are context-sensitive grammar. Assume that production rules in $\mathcal{G}_1, \mathcal{G}_2$ are **Kuroda normal form** except the rules $p_0 \rightarrow \varepsilon, q_0 \rightarrow \varepsilon$ if each of them exists in Δ_1 or Δ_2 , respectively. This assumption does not lose the generality of the discussion. (See **Exercise**). Define the grammar $\mathcal{G} = (\Sigma_3, T, N, \langle p_0, q_0 \rangle, \Delta)$. The set N of non-terminal symbols is defined below, and production rules in Δ are defined in the next slide.

$$N : N_1 \cup \{ \langle p, q \rangle \mid p \in N_1, q \in N_2 \}$$

(proof cont'd) 14

Proof (cont'd)

Δ :	$\langle p_1, q_0 \rangle \rightarrow \langle q_1, q_0 \rangle r_1$	if	$p_1 \rightarrow q_1 r_1$	in	Δ_1		
	$\langle p_1, q_0 \rangle q_1 \rightarrow \langle r_1, q_0 \rangle s_1$	if	$p_1 q_1 \rightarrow r_1 s_1$	in	Δ_1		
	$\langle p_1, q_0 \rangle \rightarrow \langle q_1, q_0 \rangle$	if	$p_1 \rightarrow q_1$	in	Δ_1		
	$\langle p_1, q_0 \rangle q_1 \rightarrow p_1 \langle q_1, q_0 \rangle$	if	p_1, q_1	in	N_1		
	$p_1 \langle q_1, q_0 \rangle \rightarrow \langle p_1, q_0 \rangle q_1$	if	p_1, q_1	in	N_1		
	$\langle p_1, p_2 \rangle q_1 \rightarrow \langle p_1, q_2 \rangle \langle q_1, r_2 \rangle$	if	p_1, q_1	in	$N_1, p_2 \rightarrow q_2 r_2$	in	Δ_2
	$\langle p_1, p_2 \rangle \langle q_1, q_2 \rangle \rightarrow \langle p_1, r_2 \rangle \langle q_1, s_2 \rangle$	if	p_1, q_1	in	$N_1, p_2 q_2 \rightarrow r_2 s_2$	in	Δ_2
	$\langle p_1, p_2 \rangle \rightarrow \langle p_1, q_2 \rangle$	if	p_1	in	$N_1, p_2 \rightarrow q_2$	in	Δ_2
	$\langle p_1, p_2 \rangle \rightarrow a$	if	$p_1 \rightarrow a$	in	$\Delta_1, p_2 \rightarrow a$	in	Δ_2
	$\langle p_0, q_0 \rangle \rightarrow \varepsilon$	if	$p_0 \rightarrow \varepsilon$	in	$\Delta_1, q_0 \rightarrow \varepsilon$	in	Δ_2

One can show that for each p_1, p_2, \dots, p_n in N_1 and q_1, q_2, \dots, q_n in N_2 ,

$$p_0 \rightarrow_{\mathcal{G}_1}^* p_1 p_2 \cdots p_n \text{ iff } \langle p_0, q_0 \rangle \rightarrow_{\mathcal{G}}^* \langle p_1, q_0 \rangle p_2 \cdots p_n$$

$$q_0 \rightarrow_{\mathcal{G}_2}^* q_1 q_2 \cdots q_n \text{ iff } \langle p_1, q_0 \rangle p_2 \cdots p_n \rightarrow_{\mathcal{G}}^* \langle p_1, q_1 \rangle \langle p_2, q_2 \rangle \cdots \langle p_n, q_n \rangle$$

Corollary

For the class of $C(CS_T)$, question if $\mathcal{L}(\mathcal{G}) = \emptyset$ is **undecidable**
 (Cf. The intersection-emptiness problem for $C(CF_T)$ on page 12)

Exercise

1. Show that for the class of context-sensitive grammar, membership problem is decidable.

2. Show Claim 2.

3. **Chomsky normal form** if production rules are in the following forms:

$$p \rightarrow qr \quad p \rightarrow a \quad q_0 \rightarrow \varepsilon \quad (p \in N, q, r \in N - \{q_0\}, a \in T)$$

Show that a language is generated by a grammar in Chomsky normal form if and only if it is generated by a context-free grammar.

4. A grammar $\mathcal{G} = (\Sigma, T, N, q_0, \Delta)$ is called **Kuroda normal form** if production rules are in the following forms:

$$pq \rightarrow rs \quad p \rightarrow qr \quad p \rightarrow q \quad p \rightarrow a \quad (p, q, r, s \in N, a \in T)$$

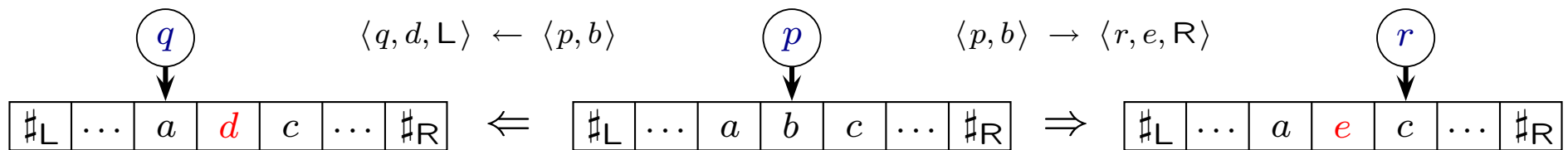
Show that languages generated by grammar in Kuroda normal form are context-sensitive, and conversely, context-sensitive languages which does not include the empty word are generated by grammar in Kuroda normal form.

Exercise (cont'd)

5. [Pumping lemma for regular grammar] Given regular language L , there exists $k \geq 0$ such that if $z \in L$ and $|z| \geq k$, then z is formed by u, v, w with $|v| \geq 1$ and $uv^nw \in L$ ($n \geq 0$). Verify this statement.
6. Show that the language $L = \{a^\ell b^\ell \mid \ell \geq 0\}$ is not generated by regular grammar.
7. Show that context-free languages over a singleton set of a symbol are generated by regular grammar.
8. For each of languages $L_1 = \{a^\ell b^\ell a^m \mid \ell, m \geq 1\}$ and $L_2 = \{a^m b^\ell a^\ell \mid \ell, m \geq 1\}$, construct context-free grammar, respectively.
9. [Liu & Weiner 1973] For each $n \geq 0$, let $C(CF_T^n)$ be the class of languages over T defined as follows: $C(CF_T^0)$ is the set of context-free languages whose set of terminal symbols is T ; $C(CF_T^{n+1})$ is $\{L_1 \cap L_2 \mid L_1 \in C(CF_T^0), L_2 \in C(CF_T^n)\}$. Show that for all $i \geq 0$, $C(CF_T^i) \subsetneq C(CF_T^{i+1})$.
10. Show that $C(CF_T^n) \subsetneq C(CS_T)$ for all $n \geq 0$.

Appendix : Linear-bounded automata (LBA)

LBA is a non-deterministic Turing machine with rewritable length-bounded single-tape as input : $(\Sigma \cup \{\#_L, \#_R\}, Q, q_0, Q_{fin}, \Delta)$, where $\#_L, \#_R$ are special tape symbols meaning left- right-ends. Transition rules in Δ are in the form of $\langle p, a \rangle \rightarrow \langle q, b, x \rangle$ ($p, q \in Q, a, b \in \Sigma, x \in \{L, R\}$).



The class of languages accepted by linear-bounded automata, denoted by $\text{NSPACE}(n)$, and $\text{C}(\text{CS}_T)$ coincide [1]. From space complexity observation, it can be shown that $\text{NSPACE}(n) = \text{co-NSPACE}(n)$ [2]. This implies that for all $L \in \text{C}(\text{CS}_T)$, $(L)^c \in \text{C}(\text{CS}_T)$, stating that $\text{C}(\text{CS}_T)$ is closed under complement.

-
- [1] S.-Y. Kuroda: *Class of Languages and Linear-Bounded Automata*, Information and Control 7, pp.207–223, 1964
- [2] N. Immerman: *Nondeterministic Space is Closed Under Complementation*, SIAM Journal of Computing 17, pp.935–938, 1988

Copyright (version Jul-01-2009) © 2009 Hitoshi Ohsaki

National Institute of Advanced Industrial Science and
Technology (AIST) – Senri-site, AIST Kansai.

Office: Shin-Senri Nishi 1–2–14 (MSK bldg. 5th floor),
Toyonaka, Osaka 560–0083, Japan

URL: <http://staff.aist.go.jp/hitoshi.ohsaki/>

All rights reserved.

No part of this lecture material may be reproduced in
any form or by any means, electronic, mechanical, pho-
tocopying, or otherwise, without the prior consent of the
author.