

# Scattered Context Grammars

Alexander Meduna   Jiří Techet

Department of Information Systems  
Faculty of Information Technology  
Brno University of Technology  
Božetěchova 2, Brno 61266, Czech Republic

Based upon the book:

A. Meduna and J. Techet: *Scattered Context Grammars and their Applications*.  
WIT Press, 2009

## Info about the Book

Publisher:	WIT Press
Address of Publisher:	Ashurst Lodge, Ashurst, Southampton, SO40 7AA, U.K.
Website:	<a href="http://www.witpress.com/978-1-84564-426-0.html">www.witpress.com/978-1-84564-426-0.html</a>
ISBN:	978-1-84564-426-0
Published:	2009
Number of Pages:	217

## Organization of the Book

- three parts:
  - I. Introduction
  - II. Theory
  - III. Application and Conclusion
- nine chapters
- exhaustive bibliography

# Scattered Context Grammars and their Applications

## Contents of the Book

<b>Preface</b> .....	vi
<b>I Introduction</b> .....	1
1 Motivation .....	3
2 Definitions .....	7
2.1 Mathematical Background .....	7
2.2 Basics of Formal Language Theory .....	9
2.3 Scattered Context Grammars .....	22
<b>II Theory</b> .....	29
3 Basic Properties .....	31
3.1 Normal Forms .....	31
3.2 Closure Properties .....	34
3.3 Generative Power .....	44

## Contents of the Book

4	Further Properties .....	51
4.1	Terminating Left-Hand Sides .....	51
4.2	Generalized $k$ -Limited Erasing .....	54
5	Restrictions and Extensions .....	71
5.1	$n$ -Limited Derivations .....	71
5.2	Leftmost Derivations .....	85
5.3	Maximal and Minimal Derivations .....	92
5.4	Unordered Scattered Context Grammars .....	101
5.5	Linear Scattered Context Grammars .....	104
5.6	Extended Propagating Scattered Context Gram- mars .....	109

## Contents of the Book

6	Reduction and Economy .....	115
6.1	Reduction .....	115
6.2	Economical Transformations .....	129
7	Parses and their Generators .....	137
7.1	Terminology .....	138
7.2	General Generators .....	141
7.3	Canonical Generators .....	147
7.4	Reduced Generators .....	155

## Contents of the Book

<b>III Applications and Conclusion</b> .....	161
<b>8 Applications in Linguistics</b> .....	163
8.1 Syntax and Related Linguistic Terminology .....	164
8.2 Transformational Scattered Context Grammars ..	168
8.3 Scattered Context in English Syntax .....	171
<b>9 Concluding Remarks</b> .....	183
<b>Bibliography</b> .....	189
<b>Language Family Index</b> .....	197
<b>Symbol Index</b> .....	201
<b>Subject Index</b> .....	205

# Introduction to Scattered Context Grammars



# Scattered Information and Its Grammatical Formalization

- while context-sensitive grammars are suitable for modelling immediate context...

AAAABC AAAA     $BC \rightarrow AA$

... they fail to describe scattered context dependencies efficiently

ABAAAAACA     $BA \rightarrow AA'$ ,  $A'A \rightarrow AA'$ ,  $A'C \rightarrow AA$

- scattered context dependencies are common in real world:

He is interested in football, isn't he?

```
... int i; int j = 10; for (i = 0; i < j; i++) { ...
```

- **scattered context grammars** (introduced by S. Greibach and J. Hopcroft in 1969) are convenient for describing this kind of dependencies

## Scattered Context Grammar

$$G = (V, T, P, S)$$

$V$  is a finite alphabet

$T$  is a set of terminals,  $T \subset V$

$S$  is the start symbol,  $S \in V - T$

$P$  is a finite set of productions of the form

$$(A_1, \dots, A_n) \rightarrow (x_1, \dots, x_n),$$

where  $A_1, \dots, A_n \in V - T$ ,  $x_1, \dots, x_n \in V^*$

## Propagating Scattered Context Grammar

- each  $(A_1, \dots, A_n) \rightarrow (x_1, \dots, x_n)$  satisfies  $x_1, \dots, x_n \in V^+$

# Derivation Step

## Derivation Step

If  $(A_1, \dots, A_n) \rightarrow (x_1, \dots, x_n) \in P$  and

$$u = u_1 A_1 \dots u_n A_n u_{n+1}$$

$$v = u_1 x_1 \dots u_n x_n u_{n+1},$$

then  $u \Rightarrow v [(A_1, \dots, A_n) \rightarrow (x_1, \dots, x_n)]$

- $\text{alph}(x)$  denotes the set of all symbols appearing in  $x$

## Leftmost Derivation Step

- each  $A_i$  satisfies  $A_i \notin \text{alph}(u_i)$

## Generated Language

- $L(G) = \{x \in T^* : S \Rightarrow^* x\}$

## Language Families

- $\mathcal{L}(SC)$  – scattered context languages
- $\mathcal{L}(PSC)$  – propagating scattered context languages

## Theorem

$$\mathcal{L}(SC) = \mathcal{L}(RE).$$

## Theorem

$$\mathcal{L}(CF) \subset \mathcal{L}(PSC) \subseteq \mathcal{L}(CS).$$

## Theorem

*For every recursively enumerable language  $L$  there exists a propagating scattered context language  $L'$  and a homomorphism  $h$  such that  $h(L') = L$ .*

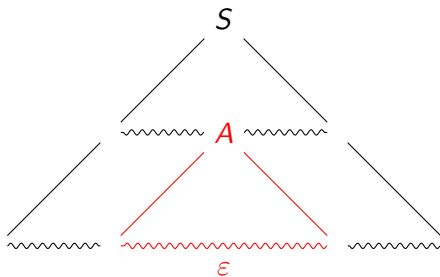
# Results

# Symbols Erased During Derivation

## Symbols Erased During Derivation

A symbol  $A$  is erased during a derivation if the frontier of the subtree rooted at  $A$  is  $\varepsilon$ ;

- if the symbol  $A$  is erased, we write  $\dot{A}$ ,
- if the symbol  $A$  is not erased, we write  $\acute{A}$

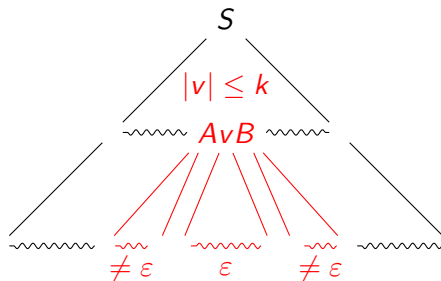


# Nonterminals Erased in a Generalized $k$ -Limited Way

## Nonterminals Erased in a Generalized $k$ -Limited Way

For  $y \in L(G)$ , every sentential form  $x$  in  $S \Rightarrow_G^* y$  satisfies:

- 1 every  $x = uAvBw$ , where  $\hat{A}$ ,  $\hat{B}$ ,  $\hat{v}$ , satisfies  $|\hat{v}| \leq k$ ,
- 2 every  $x = uAw$ , where  $\hat{A}$ , satisfies: if  $\hat{u}$  or  $\hat{w}$ , then  $|u| \leq k$  or  $|w| \leq k$ , respectively





## Nonterminals Erased in Generalized $k$ -Limited Way by SC Grammars

A scattered context grammar  $G$  **erases its nonterminals in a generalized  $k$ -limited way** if  $L(G) = L(G, \varepsilon, k)$ , where

$$L(G, \varepsilon, k) = \{x \in T^* : S \Rightarrow_G^* x, \text{ and } G \text{ erases nonterminals} \\ \text{in a generalized } k\text{-limited way in } S \Rightarrow_G^* x\}$$

### Theorem

*For each  $k \geq 0$  and every scattered context grammar  $G$ , there is a propagating scattered context grammar  $\bar{G}$  such that  $L(G, \varepsilon, k) = L(\bar{G})$ .*

### Corollary

*For every scattered context grammar  $G$  which erases its nonterminals in a generalized  $k$ -limited way, there exists a propagating scattered context grammar  $\bar{G}$  such that  $L(G) = L(\bar{G})$ .*

## Linear Scattered Context Grammar

- is a scattered context grammar  $G = (V, T, P, S)$
- $P$  is a finite set of productions of the following two forms:
  - 1  $(S) \rightarrow (x_1 A_1 \dots x_k A_k x_{k+1})$ , where  $A_i \in (V - T) - \{S\}$ ,  $x_j \in T^*$  for all  $1 \leq i \leq k$ ,  $1 \leq j \leq k + 1$ , for some  $k \geq 1$ ,
  - 2  $(A_1, \dots, A_k) \rightarrow (z_1, \dots, z_k)$ , where  $A_i \in (V - T) - \{S\}$ , and either
    - $z_i = x_i B_i y_i$ , where  $x_i, y_i \in T^*$ ,  $B_i \in (V - T) - \{S\}$ , or
    - $z_i \in T^*$for all  $1 \leq i \leq k$ , for some  $k \geq 1$

## Linear Scattered Context Grammar of Degree $n$

- every  $(S) \rightarrow (y_1 A_1 \dots y_m A_m y_{m+1}) \in P$  satisfies  $m \leq n$

## Right-Linear Scattered Context Grammar

- is a linear scattered context grammar  $G = (V, T, P, S)$
- $P$  is a finite set of productions of the following two forms:
  - 1  $(S) \rightarrow (x_1 A_1 \dots x_k A_k)$ , where  $A_i \in (V - T) - \{S\}$ ,  $x_i \in T^*$  for all  $1 \leq i \leq k$ , for some  $k \geq 1$ ,
  - 2  $(A_1, \dots, A_k) \rightarrow (z_1, \dots, z_k)$ , where  $A_i \in (V - T) - \{S\}$ , and either
    - $z_i = x_i B_i$ , where  $x_i \in T^*$ ,  $B_i \in (V - T) - \{S\}$ , or
    - $z_i \in T^*$for all  $1 \leq i \leq k$ , for some  $k \geq 1$

## Language Families

- $\mathcal{L}(SC, LIN, n)$  – linear scattered context grammars of degree  $n$
- $\mathcal{L}(SC, RLIN, n)$  – right-linear scattered context grammars of degree  $n$

## Theorem

For each  $n \geq 1$ ,

$$\begin{aligned}\mathcal{L}(SC, LIN, n) &\subset \mathcal{L}(SC, LIN, n + 1), \\ \mathcal{L}(SC, RLIN, n) &\subset \mathcal{L}(SC, RLIN, n + 1), \\ \mathcal{L}(SC, RLIN, n) &\subset \mathcal{L}(SC, LIN, n).\end{aligned}$$

- $\mathcal{L}(SC, LIN) = \bigcup_{n=1}^{\infty} \mathcal{L}(SC, LIN, n)$
- $\mathcal{L}(SC, RLIN) = \bigcup_{n=1}^{\infty} \mathcal{L}(SC, RLIN, n)$

## Theorem

$$\begin{aligned}\mathcal{L}(CF) - \mathcal{L}(SC, LIN) &\neq \emptyset, \quad \mathcal{L}(CF) - \mathcal{L}(SC, RLIN) \neq \emptyset, \\ \mathcal{L}(SC, RLIN) &\subset \mathcal{L}(SC, LIN) \subset \mathcal{L}(PSC).\end{aligned}$$

# $n$ -Limited Derivations

- $|x|_W$  denotes the number of occurrences of symbols from set  $W$  in  $x$

## $n$ -Limited Derivation Step

If  $(A_1, \dots, A_k) \rightarrow (x_1, \dots, x_k) \in P$ ,

$$u = u_1 A_1 u_2 \dots u_k A_k u_{k+1},$$

$$v = u_1 x_1 u_2 \dots u_k x_k u_{k+1},$$

and  $u$  satisfies

$$|u_1 A_1 \dots u_k A_k|_{V-T} \leq n,$$

then  $u \stackrel{n}{\text{lim}} \Rightarrow_G v$

## $n$ -Limited Derivation

- derivation  $x \stackrel{n}{\text{lim}} \Rightarrow_G^* y$  in which every derivation step  $u \stackrel{j}{\text{lim}} \Rightarrow_G v$  satisfies  $j \leq n$

## Language of Order $n$

- $L(G, \text{lim}, n) = \{x \in T^* : S \stackrel{n}{\text{lim}} \Rightarrow_G^* x\}$

## Language Families

- $\mathcal{L}(\text{PSC}, \text{lim}, n)$
- $\mathcal{L}(\text{PSC}, \text{lim}, \infty) = \bigcup_{i=1}^{\infty} \mathcal{L}(\text{PSC}, \text{lim}, i)$

## Theorem

$$\mathcal{L}(\text{CF}) = \mathcal{L}(\text{PSC}, \text{lim}, 1) \subset \dots \subset \mathcal{L}(\text{PSC}, \text{lim}, \infty) \subset \mathcal{L}(\text{CS}).$$

# Leftmost Derivations

- much simplified proof of the result proved by V. Virkkunen in 1973

## Propagating Scattered Context Grammar which Uses Leftmost Derivations

- propagating scattered context grammar  $G = (V, T, P, S)$  whose language is defined as

$$L(G, \text{lm}) = \{x \in T^* : S \xrightarrow{\text{lm}}_G^* x\}$$

## Language Family

- $\mathcal{L}(\text{PSC}, \text{lm})$

## Theorem

$$\mathcal{L}(\text{PSC}, \text{lm}) = \mathcal{L}(\text{CS}).$$

# Maximal and Minimal Derivation

■  $\text{len}((A_1, \dots, A_n) \rightarrow (x_1, \dots, x_n)) = |A_1 \dots A_n| = n$

## Maximal Derivation Step

Let  $p \in P$ . If

1  $u \Rightarrow v [p]$

2 for every  $r \in P$  such that  $u \Rightarrow w [r] : \text{len}(p) \geq \text{len}(r)$

then  $u \text{max} \Rightarrow v [p]$

## Minimal Derivation Step

Let  $p \in P$ . If

1  $u \Rightarrow v [p]$

2 for every  $r \in P$  such that  $u \Rightarrow w [r] : \text{len}(p) \leq \text{len}(r)$

then  $u \text{min} \Rightarrow v [p]$



## Maximal and Minimal Languages

- $L(G, \text{max}) = \{x \in T^* : S \xrightarrow{\text{max}}^* x\}$
- $L(G, \text{min}) = \{x \in T^* : S \xrightarrow{\text{min}}^* x\}$

## Language Families

- $\mathcal{L}(PSC, \text{max})$
- $\mathcal{L}(PSC, \text{min})$

## Theorem

$$\mathcal{L}(CS) = \mathcal{L}(PSC, \text{max}).$$

## Theorem

$$\mathcal{L}(CS) = \mathcal{L}(PSC, \text{min}).$$

## Production Label

- for every grammar  $G$ ,  $lab(G)$  denotes the set of its production labels
- each  $p \in lab(G)$  uniquely identifies one production:

$$p : (A_1, \dots, A_n) \rightarrow (x_1, \dots, x_n)$$

## Derivation Made by Productions

- if  $x \Rightarrow y$  by  $p : (A_1, \dots, A_n) \rightarrow (x_1, \dots, x_n)$ , we write

$$x \Rightarrow y [p]$$

- if  $x \Rightarrow^* y$  by productions labeled with  $p_1, \dots, p_n$ , we write

$$x \Rightarrow^* y [p_1 \dots p_n]$$

# Proper Generator of Its Sentences With Their Parses

## Parse (Szilard Word, Control Word)

If

$$S \Rightarrow^* x [\rho],$$

where  $x \in T^*$ ,  $\rho \in \text{lab}(G)^*$ , then  $x$  is a sentence generated according to parse  $\rho$

## Proper Generator of Its Sentences With Their Parses

- $G = (V, T, P, S)$ , where  $\text{lab}(G) \subset T$ , which satisfies

$$L(G) = \{x : x = y\rho, y \in (T - \text{lab}(G))^*, \rho \in \text{lab}(G)^*, S \Rightarrow^* x [\rho]\}$$

- **leftmost generator** makes every successful derivation in a **leftmost** way

# Results

- $G = (V, P, S, T)$  is a proper generator of its sentences with their parses
- weak identity  $\pi$  from  $V^*$  to  $(V - lab(G))^*$ :
  - $\pi(a) = a$  for each  $a \in (V - lab(G))$
  - $\pi(p) = \epsilon$  for each  $p \in lab(G)$

## Theorem

*For every recursively enumerable language  $L$ , there exists a **propagating** scattered context grammar  $G$  such that  $G$  is a proper generator of its sentences with their parses and  $L = \pi(L(G))$ .*

## Theorem

*For every recursively enumerable language  $L$ , there exists a **propagating** scattered context grammar  $G = (V, T, P, S)$  such that  $G$  is a proper **leftmost** generator of its sentences with their parses,  $|V - T| \leq 6$ , and  $L = \pi(L(G))$ .*

# Applications in Linguistics

- there are scattered dependencies in natural languages

He usually, but not always, goes to work early.

- there is a scattered dependency between the subject (he) and the predicator (goes):

I usually, but not always, goes to work early.

- the dependency can be easily captured by productions of scattered context grammars:

(He, goes)  $\rightarrow$  (I, go)

transforms the original sentence to

I usually, but not always, go to work early.

# Example

Consider the language  $L$  consisting of these grammatical English sentences:

*Your grandparents are all your grandfathers and all your grandmothers.*

*Your **great**-grandparents are all your **great**-grandfathers and all your **great**-grandmothers.*

*Your **great-great**-grandparents are all your **great-great**-grandfathers and all your **great-great**-grandmothers.*

⋮

In brief,

$$L = \{ \text{your } \{\mathbf{great-}\}^i \text{grandparents are all your } \{\mathbf{great-}\}^i \text{grandfathers and all your } \{\mathbf{great-}\}^i \text{grandmothers} : i \geq 0 \}.$$

# Example

Introduce the scattered context grammar  $G = (V, T, P, S)$ , where

$T = \{\text{all, and, are, grandfathers, grandmothers, grandparents, great-, your}\}$ ,

$V = T \cup \{S, \#\}$ , and  $P$  consists of these three productions:

$(S) \rightarrow (\text{your } \# \text{grandparents are all your } \# \text{grandfathers}$   
 $\text{and all your } \# \text{grandmothers}),$   
 $(\#, \#, \#) \rightarrow (\# \text{great-}, \# \text{great-}, \# \text{great-}),$   
 $(\#, \#, \#) \rightarrow (\varepsilon, \varepsilon, \varepsilon).$

Obviously, this scattered context grammar generates  $L$ ; formally,  $L = L(G)$ .



# Conclusion

## Possibilities of Practical Applications

- applications in compilers
- applications in natural language processing

## Main Open Problem

- $\mathcal{L}(PSC) = \mathcal{L}(CS)?$

## Journal Articles by the Authors

- 2009 A. Meduna and J. Techet: An infinite hierarchy of language families generated by scattered context grammars with  $n$ -limited derivations, *Theoretical Computer Science*, 2009, in press
- 2008 Masopust T., Techet J.: Leftmost Derivations of Propagating Scattered Context Grammars: A New Proof, *Discrete Mathematics and Theoretical Computer Science*, 10, 39–46
- 2008 A. Meduna and J. Techet: Scattered context grammars that erase nonterminals in a generalized  $k$ -limited way, *Acta Informatica*, 45(7), 593–608
- 2007 Meduna A., Techet J.: Canonical Scattered Context Generators of Sentences with Their Parses, *Theoretical Computer Science*, 389, 73–81
- 2005 Meduna A., Techet J.: Generation of Sentences with Their Parses: the Case of Propagating Scattered Context Grammars, *Acta Cybernetica*, 17, 11–20

- 2007 Meduna A., Techet J.: Maximal and Minimal Scattered Context Rewriting, *FCT 2007 Proceedings*, Budapest, 412–423, (LNCS)
- 2007 Meduna A., Techet J.: Reduction of Scattered Context Generators of Sentences Preceded by Their Leftmost Parses, *Proceedings of DCFS 2007*, High Tatras, 178–185