

# Scattered Context Grammars

A. Meduna   J. Techet

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

2008

- 2008 Masopust T., Techet J.: Leftmost Derivations of Propagating Scattered Context Grammars: A New Proof, *Discrete Mathematics and Theoretical Computer Science*, 10, 39–46
- 2007 Meduna A., Techet J.: Canonical Scattered Context Generators of Sentences with Their Parses, *Theoretical Computer Science*, 389, 73–81
- 2007 Meduna A., Techet J.: Maximal and Minimal Scattered Context Rewriting, *FCT 2007 Proceedings*, Budapest, 412–423
- 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
- 2007 Techet J.:  $k$ -Limited Erasing Performed by Scattered Context Grammars, *WFM 2007*, Hradec nad Moravicí, 227–234
- 2005 Meduna A., Techet J.: Generation of Sentences with Their Parses: the Case of Propagating Scattered Context Grammars, *Acta Cybernetica*, 17, 11–20

- 1 Definitions and Basic Properties
- 2 Conditional Removal of Erasing Productions
- 3 Restrictions and Extensions
  - Non-Context-Free Components of Scattered Context Grammars
  - $n$ -Limited Derivations
  - Leftmost Derivations
  - Maximal and Minimal Rewriting
- 4 Generators of Sentences with Their parses
- 5 Applications in Linguistics
- 6 Further Investigation

- 1 Definitions and Basic Properties
- 2 Conditional Removal of Erasing Productions
- 3 Restrictions and Extensions
  - Non-Context-Free Components of Scattered Context Grammars
  - $n$ -Limited Derivations
  - Leftmost Derivations
  - Maximal and Minimal Rewriting
- 4 Generators of Sentences with Their parses
- 5 Applications in Linguistics
- 6 Further Investigation

# Scattered Information and Its Grammatical Formalization

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

AAAA**BC**AAAA     **BC**  $\rightarrow$  **AA**

... they fail to describe scattered context dependencies efficiently

A**B**AAAAAA**C**A     **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 Grammars

## 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

## 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

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

## Example

Scattered context grammar

$$G = (\{A, B, C, S, a, b, c\}, \{a, b, c\}, P, S),$$

where

$$\begin{aligned} P = \{ & (S) \rightarrow (ABC), \\ & (A, B, C) \rightarrow (aA, bB, cC), \\ & (A, B, C) \rightarrow (\varepsilon, \varepsilon, \varepsilon) \} \end{aligned}$$

Example of derivation

$$S \Rightarrow ABC \Rightarrow aAbBcC \Rightarrow aaAbbBccC \Rightarrow aabbcc$$

Generated language

$$L(G) = \{a^n b^n c^n : n \geq 0\}$$

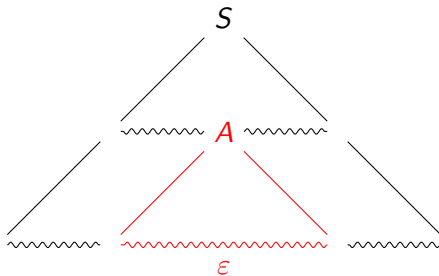
- 1 Definitions and Basic Properties
- 2 Conditional Removal of Erasing Productions
- 3 Restrictions and Extensions
  - Non-Context-Free Components of Scattered Context Grammars
  - $n$ -Limited Derivations
  - Leftmost Derivations
  - Maximal and Minimal Rewriting
- 4 Generators of Sentences with Their parses
- 5 Applications in Linguistics
- 6 Further Investigation

# 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 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  $L(G) = L(\bar{G})$ .*

- 1 Definitions and Basic Properties
- 2 Conditional Removal of Erasing Productions
- 3 Restrictions and Extensions
  - Non-Context-Free Components of Scattered Context Grammars
  - $n$ -Limited Derivations
  - Leftmost Derivations
  - Maximal and Minimal Rewriting
- 4 Generators of Sentences with Their parses
- 5 Applications in Linguistics
- 6 Further Investigation

## 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  $A_i \in (V - T) - \{S\}$ ,  $x_j \in T^*$ ,
  - 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
    - $z_i \in T^*$

## Linear Scattered Context Grammar of Degree $n$

- $(S) \rightarrow (x_1 A_1 \dots x_n A_n x_{n+1}) \in P$  satisfies  $n \geq m$  for all  $(S) \rightarrow (y_1 A_1 \dots y_m A_m y_{m+1}) \in P$

# Right-Linear Scattered Context Grammars

## 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  $A_i \in (V - T) - \{S\}$ ,  $x_i \in T^*$ ,
  - 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
    - $z_i \in T^*$

## 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

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

- 1 Definitions and Basic Properties
- 2 Conditional Removal of Erasing Productions
- 3 Restrictions and Extensions
  - Non-Context-Free Components of Scattered Context Grammars
  - $n$ -Limited Derivations
  - Leftmost Derivations
  - Maximal and Minimal Rewriting
- 4 Generators of Sentences with Their parses
- 5 Applications in Linguistics
- 6 Further Investigation

# $n$ -Limited Derivations

- $|x|_A$  denotes the number of occurrences of symbols from  $A$  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  $u \xRightarrow[n]{\text{lim}}_G v$

## $n$ -Limited Derivation

- derivation  $x \xRightarrow[n]{\text{lim}}^*_G y$  in which every derivation step  $u \xRightarrow[j]{\text{lim}}_G v$  satisfies  $j \leq n$

## Language of Order $n$

- $L(G, \text{lim}, n) = \{x \in T^* : S \xRightarrow[n]{\text{lim}}_G^* x\}$

## Language Families

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

## Theorem

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

- 1 Definitions and Basic Properties
- 2 Conditional Removal of Erasing Productions
- 3 **Restrictions and Extensions**
  - Non-Context-Free Components of Scattered Context Grammars
  - $n$ -Limited Derivations
  - **Leftmost Derivations**
  - Maximal and Minimal Rewriting
- 4 Generators of Sentences with Their parses
- 5 Applications in Linguistics
- 6 Further Investigation

# 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

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

## Language Family

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

## Theorem

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

- 1 Definitions and Basic Properties
- 2 Conditional Removal of Erasing Productions
- 3 **Restrictions and Extensions**
  - Non-Context-Free Components of Scattered Context Grammars
  - $n$ -Limited Derivations
  - Leftmost Derivations
  - **Maximal and Minimal Rewriting**
- 4 Generators of Sentences with Their parses
- 5 Applications in Linguistics
- 6 Further Investigation

# Maximal and Minimal Derivation

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

## Maximal Derivation Step

If

1  $u \Rightarrow v [p]$ ,

2 there is no  $r \in P$ ,  $\text{len}(r) > \text{len}(p)$ , such that  $u \Rightarrow w [r]$ ,

where  $p \in P$ , then  $u \text{max} \Rightarrow v [p]$

## Minimal Derivation Step

If

1  $u \Rightarrow v [p]$ ,

2 there is no  $r \in P$ ,  $\text{len}(r) < \text{len}(p)$ , such that  $u \Rightarrow w [r]$ ,

where  $p \in P$ , 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

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

- 1 Definitions and Basic Properties
- 2 Conditional Removal of Erasing Productions
- 3 Restrictions and Extensions
  - Non-Context-Free Components of Scattered Context Grammars
  - $n$ -Limited Derivations
  - Leftmost Derivations
  - Maximal and Minimal Rewriting
- 4 Generators of Sentences with Their parses
- 5 Applications in Linguistics
- 6 Further Investigation

# Production Labels

## 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** uses leftmost derivations in a successful derivation

# 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  $L = \pi(L(G))$ .*

- 1 Definitions and Basic Properties
- 2 Conditional Removal of Erasing Productions
- 3 Restrictions and Extensions
  - Non-Context-Free Components of Scattered Context Grammars
  - $n$ -Limited Derivations
  - Leftmost Derivations
  - Maximal and Minimal Rewriting
- 4 Generators of Sentences with Their parses
- 5 Applications in Linguistics
- 6 Further Investigation

- there are scattered dependencies in natural languages

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

- the subject (he) depends on the predicator (goes):

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

is grammatically incorrect

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

$(\text{He, goes}) \rightarrow (\text{I, go})$

transforms the original sentence to

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

- these transformations are formalized and studied on simple sentences

# Transformational Scattered Context Grammar

## Transformational Scattered Context Grammar

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

$V$  is the total vocabulary

$T \subset V$  is a set of terminals

$P$  is a set of scattered context productions

$I \subset V$  is the **input vocabulary**

## Transformation $T$ Defined by $G$

- $T(G) = \{(x, y) : x \Rightarrow_G^* y, x \in I^*, y \in T^*\}$
- if  $(x, y) \in T(G)$ , then  $x$  **is transformed to**  $y$  by  $G$
- $x$  and  $y$  are called the **input** and the **output sentence**

# Clauses with *neither* and *nor*

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

$T$  is the set of all words, including all their inflectional forms

$$I = \{\langle x \rangle : x \in T\}$$

$$V = T \cup I$$

$$P = \{(\langle \text{neither} \rangle, \langle \text{nor} \rangle) \rightarrow (\text{both}, \text{and})\} \\ \cup \{(\langle x \rangle) \rightarrow (x) : x \in T - \{\text{neither}, \text{nor}\}\}$$

## Example

$\langle \text{neither} \rangle \langle \text{thomas} \rangle \langle \text{nor} \rangle \langle \text{his} \rangle \langle \text{wife} \rangle \langle \text{went} \rangle \langle \text{to} \rangle \langle \text{the} \rangle \langle \text{party} \rangle$   
 $\Rightarrow_G$  both  $\langle \text{thomas} \rangle$  and  $\langle \text{his} \rangle \langle \text{wife} \rangle \langle \text{went} \rangle \langle \text{to} \rangle \langle \text{the} \rangle \langle \text{party} \rangle$   
 $\Rightarrow_G$  both thomas and  $\langle \text{his} \rangle \langle \text{wife} \rangle \langle \text{went} \rangle \langle \text{to} \rangle \langle \text{the} \rangle \langle \text{party} \rangle$   
 $\Rightarrow_G$  both thomas and his  $\langle \text{wife} \rangle \langle \text{went} \rangle \langle \text{to} \rangle \langle \text{the} \rangle \langle \text{party} \rangle$   
 $\Rightarrow_G^5$  both thomas and his wife went to the party

# Existential Clauses

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

$T$  is the set of all words, including all their inflectional forms

$$I = \{\langle x \rangle : x \in T\}$$

$$V = T \cup I \cup \{X\}$$

$$P = \{(\langle x \rangle, \langle \text{is} \rangle) \rightarrow (\text{there is } xX, \varepsilon), (\langle x \rangle, \langle \text{are} \rangle) \rightarrow (\text{there are } xX, \varepsilon), \\ (\langle x \rangle, \langle \text{was} \rangle) \rightarrow (\text{there was } xX, \varepsilon), (\langle x \rangle, \langle \text{were} \rangle) \rightarrow (\text{there were } xX, \varepsilon) \\ \cup \{(X, \langle x \rangle) \rightarrow (X, x) : x \in T\} \cup \{(X) \rightarrow (\varepsilon)\}$$

## Example

$\langle a \rangle \langle \text{nurse} \rangle \langle \text{was} \rangle \langle \text{present} \rangle \Rightarrow_G \text{there was a } X \langle \text{nurse} \rangle \langle \text{present} \rangle$   
 $\Rightarrow_G \text{there was a } X \text{ nurse } \langle \text{present} \rangle \Rightarrow_G \text{there was a } X \text{ nurse present}$   
 $\Rightarrow_G \text{there was a nurse present}$

- 1 Definitions and Basic Properties
- 2 Conditional Removal of Erasing Productions
- 3 Restrictions and Extensions
  - Non-Context-Free Components of Scattered Context Grammars
  - $n$ -Limited Derivations
  - Leftmost Derivations
  - Maximal and Minimal Rewriting
- 4 Generators of Sentences with Their parses
- 5 Applications in Linguistics
- 6 Further Investigation

# Further Investigation

## Possibilities of Practical Applications

- applications in compilers
- applications in natural language processing

## Theoretical Properties

- our research rises several open problems:
  - Does generalized  $k$ -limited erasing cover all possible kinds of erasing which can be performed by propagating scattered context grammars?
  - For given  $n$ , is it possible to construct a propagating scattered context grammar which rewrites the first  $n$  nonterminals without any additional restrictions?
- other papers dealing with this topic introduced several open problems as well
- the problem of whether  $\mathcal{L}(PSC) = \mathcal{L}(CS)$  remains still open

Thank you