

Pumping and Ogden Properties of Multiple Context-Free Grammars

Makoto Kanazawa
National Institute of Informatics and SOKENDAI
Japan

略歴

- 1994: Ph.D. in Linguistics, Stanford University
- 1994: 千葉大学文学部行動科学科
- 2000: 東京大学情報学環
- 2004: 国立情報学研究所
- 2018: 法政大学理工学部創生科学科

2

Multiple Context-Free Grammars

- Introduced by Seki, Matsumura, Fujii, and Kasami (1987–1991)
- Independently by Vijay-Shanker, Weir, and Joshi (1987)
- **Many** equivalent models
- Often thought to be an adequate formalization of **mildly context-sensitive** grammars (Joshi 1985)

3

Arising from concerns in computational linguistics. CFGs are almost good enough for NL grammars, but not quite; a mild extension of CFGs is needed. Several criteria were put forward as to what constitutes a “mild” extension.

Context-Free Grammars

production

$$A \rightarrow w_0 B_1 w_1 \dots B_n w_n \quad B_i \in N, w_j \in \Sigma^*$$

$$\frac{}{S \Rightarrow_G^* S} \quad \frac{S \Rightarrow_G^* \beta A \gamma \quad A \rightarrow \alpha \in P}{S \Rightarrow_G^* \beta \alpha \gamma}$$

top-down derivation

$$L(G) = \{ w \in \Sigma^* \mid S \Rightarrow_G^* w \}$$

4

Bottom-Up Interpretation

$$\frac{B_i \Rightarrow_G^* v_i \ (i = 1, \dots, n) \quad A \rightarrow w_0 B_1 w_1 \dots B_n w_n \in P}{A \Rightarrow_G^* w_0 v_1 w_1 \dots v_n w_n}$$

$$L(G) = \{ w \in \Sigma^* \mid S \Rightarrow_G^* w \}$$

5

CFGs as Logic Programs on Strings

$$A \rightarrow w_0 B_1 w_1 \dots B_n w_n$$

$$A(w_0 \mathbf{x}_1 w_1 \dots \mathbf{x}_n w_n) \leftarrow B_1(\mathbf{x}_1), \dots, B_n(\mathbf{x}_n)$$

Horn clause

$$L(G) = \{ w \in \Sigma^* \mid G \vdash S(w) \}$$

6

Multiple Context-Free Grammars

$$A(\alpha_1, \dots, \alpha_q) \leftarrow B_1(\mathbf{x}_{1,1}, \dots, \mathbf{x}_{1,q_1}), \dots, B_n(\mathbf{x}_{n,1}, \dots, \mathbf{x}_{n,q_n})$$

$n \geq 0, q, q_i \geq 1,$
 $\alpha_k \in (\Sigma \cup \{ \mathbf{x}_{i,j} \mid i \in [1,n], j \in [1,q_i] \})^*$
 each $\mathbf{x}_{i,j}$ occurs exactly once in $(\alpha_1, \dots, \alpha_q)$

- $q = \dim(A)$ (**dimension** of A)
- $\dim(S) = 1$
- $L(G) = \{ w \in \Sigma^* \mid G \vdash S(w) \}$

It's best to think of an MCFG as a kind of logic program.

Each rule is a definite clause.

Nonterminals are predicates on strings.

$S(\mathbf{x}_1 \# \mathbf{x}_2) \leftarrow D(\mathbf{x}_1, \mathbf{x}_2)$
 $D(\varepsilon, \varepsilon) \leftarrow$
 $D(\mathbf{x}_1 \mathbf{y}_1, \mathbf{y}_2 \mathbf{x}_2) \leftarrow E(\mathbf{x}_1, \mathbf{x}_2), D(\mathbf{y}_1, \mathbf{y}_2)$
 $E(a \mathbf{x}_1 \bar{a}, \bar{a} \mathbf{x}_2 a) \leftarrow D(\mathbf{x}_1, \mathbf{x}_2)$

$\{ w \# w^R \mid w \in D_1^* \}$

2-MCFG
2-ary branching

derivation tree

8

m-MCFG = MCFG with nonterminal dimension not exceeding m

1-MCFG = CFG

Derivation tree for w = proof of S(w)

$S(\mathbf{x}_1 \dots \mathbf{x}_m) \leftarrow A(\mathbf{x}_1, \dots, \mathbf{x}_m)$
 $A(\varepsilon, \dots, \varepsilon) \leftarrow$
 $A(a_1 \mathbf{x}_1 a_2, \dots, a_{2m-1} \mathbf{x}_m a_{2m}) \leftarrow A(\mathbf{x}_1, \dots, \mathbf{x}_m)$

non-branching m-MCFG

Seki et al. 1991

9

The languages of MCFGs form an infinite hierarchy.

Chomsky Hierarchy

Rewriting Grammars	Machines	Logic Programs on Strings	Languages
Unrestricted	Turing	Elementary Formal Systems (Smullyan 1961)	r.e.
Context-Sensitive	LBA	Length-Bounded EFS (Arikawa et al. 1989)	CSL = NSPACE(n)
	Poly-time Turing	Simple LMG (Groenink 1997) / Hereditary EFS (Ikeda and Arimura 1997)	P
		MCFG	MCFL
Context-Free	PDA	Simple EFS (Arikawa 1970)	CFL
Right-Linear	FA		Reg

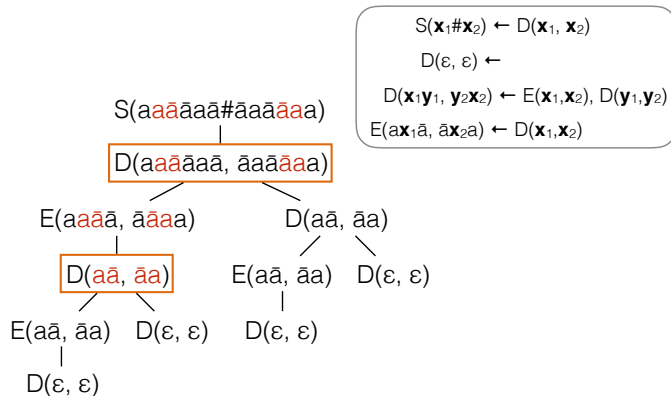
10

Which Properties of CFGs Are Shared by/ Generalize to MCFGs?

- Membership in LOGCFL
- Semilinearity
- ...

11

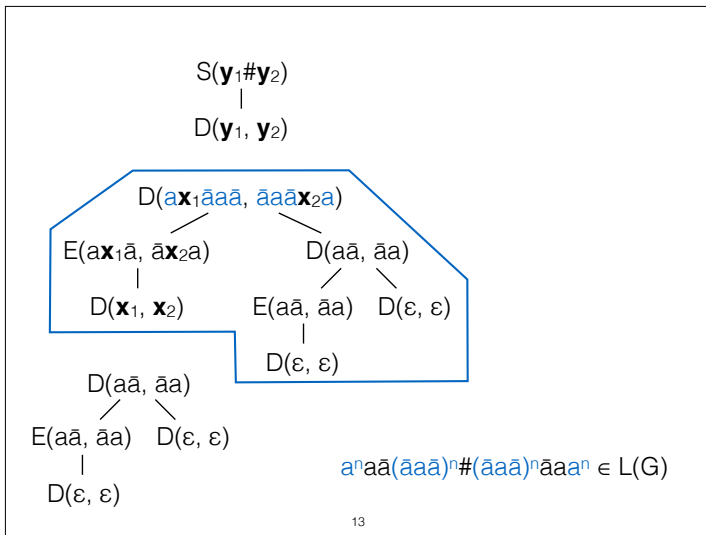
Pumping



12

Derivation trees of MCFGs are similar to those of CFGs.

When the same nonterminal occurs twice on the same path of a derivation tree,...



You can decompose the derivation tree into three parts, and the middle part can be iterated any number of times, including zero times. In the overall derivation tree, the variables x_1, x_2, y_1, y_2 are instantiated by ... The number of iterated substrings (factors) larger than two.

Iterative Properties

L is **k-iterative** iff $\exists p \forall z \in L (|z| \geq p \Rightarrow$
 $\exists u_1 \dots u_{k+1} v_1 \dots v_k ($
 $z = u_1 v_1 \dots u_k v_k u_{k+1} \wedge$
 $v_1 \dots v_k \neq \epsilon \wedge$
 $\forall n \geq 0 (u_1 v_1^n \dots u_k v_k^n u_{k+1} \in L))$

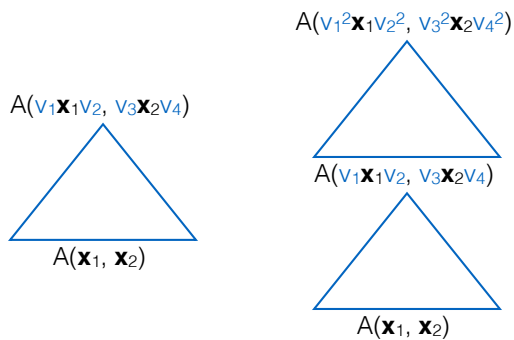
$L \in CFL \Rightarrow L$ is 2-iterative

$L \in m\text{-MCFL} \Rightarrow L$ is $2m$ -iterative ?

wrong claim in 1991

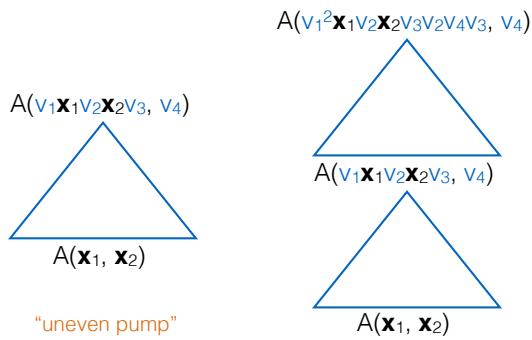
For MCFGs, need to consider a generalized form of the condition of the pumping lemma. Not straightforward; open question for a long time.

Difficulty with Pumping



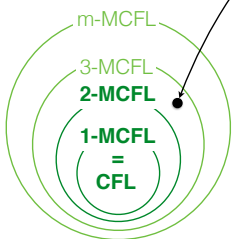
The middle part of the derivation tree may look like this.

Difficulty with Pumping



Or like this.

$S(x_1 \# x_2 \# x_3) \leftarrow A(x_1, x_2, x_3)$
 $A(ax_1, y_1cx_2\bar{c}dy_2\bar{d}x_3, y_3b) \leftarrow A(x_1, x_2, x_3), A(y_1, y_2, y_3)$
 $A(a, \epsilon, b) \leftarrow$

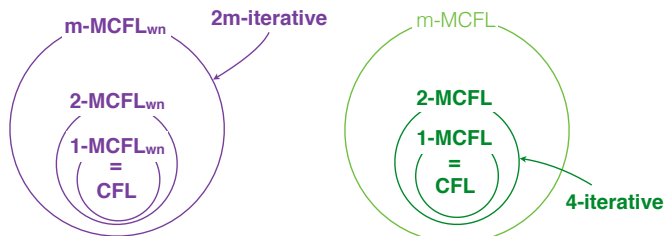


not k-iterative for any k

Kanazawa et al. 2014

The pumping lemma fails for 3-MCFGs.

Pumping Lemma for Subclasses



well-nested MCFGs

Kanazawa 2009

Pumping possible for special cases. Well-nested MCFGs.

Well-Nestedness

$\{ w\#w^R \mid w \in D_1^* \}$

$\{ w\#w \mid w \in D_1^* \}$

$S(\mathbf{x}_1\#\mathbf{x}_2) \leftarrow D(\mathbf{x}_1, \mathbf{x}_2)$

$D(\varepsilon, \varepsilon) \leftarrow$

$D(\mathbf{x}_1\mathbf{y}_1, \mathbf{y}_2\mathbf{x}_2) \leftarrow E(\mathbf{x}_1, \mathbf{x}_2), D(\mathbf{y}_1, \mathbf{y}_2)$



$E(\mathbf{ax}_1\bar{a}, \bar{a}\mathbf{x}_2\mathbf{a}) \leftarrow D(\mathbf{x}_1, \mathbf{x}_2)$

well-nested

$S(\mathbf{x}_1\#\mathbf{x}_2) \leftarrow D(\mathbf{x}_1, \mathbf{x}_2)$

$D(\varepsilon, \varepsilon) \leftarrow$

$D(\mathbf{x}_1\mathbf{y}_1, \mathbf{x}_2\mathbf{y}_2) \leftarrow E(\mathbf{x}_1, \mathbf{x}_2), D(\mathbf{y}_1, \mathbf{y}_2)$



$E(\mathbf{ax}_1\bar{a}, \mathbf{ax}_2\bar{a}) \leftarrow D(\mathbf{x}_1, \mathbf{x}_2)$

non-well-nested

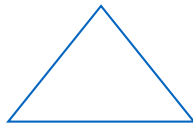
$\{ w\#w \mid w \in D_1^* \} \notin \text{MCFL}_{\text{wn}}$

Kanazawa and Salvati 2010

Has a natural equivalent characterization: $y\text{CFT}_{\text{sp}}$

Difficulty with Pumping

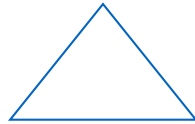
$A(\mathbf{v}_1\mathbf{x}_1\mathbf{v}_2\mathbf{x}_2\mathbf{v}_3, \mathbf{v}_4)$



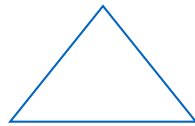
$A(\mathbf{x}_1, \mathbf{x}_2)$

"uneven pump"

$A(\mathbf{v}_1^2\mathbf{x}_1\mathbf{v}_2\mathbf{x}_2\mathbf{v}_3\mathbf{v}_2\mathbf{v}_4\mathbf{v}_3, \mathbf{v}_4)$



$A(\mathbf{v}_1\mathbf{x}_1\mathbf{v}_2\mathbf{x}_2\mathbf{v}_3, \mathbf{v}_4)$



$A(\mathbf{x}_1, \mathbf{x}_2)$

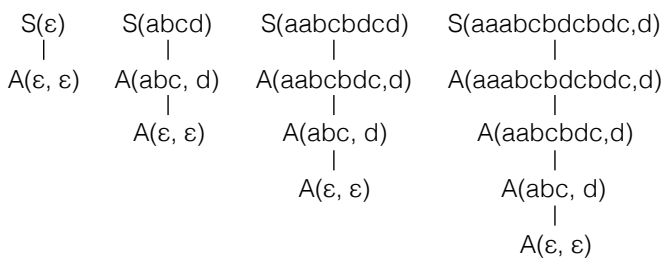
Pumping not easy to prove even for well-nested MCFGs: this situation can still arise.

$S(\mathbf{x}_1\mathbf{x}_2) \leftarrow A(\mathbf{x}_1, \mathbf{x}_2)$

$A(\mathbf{ax}_1\mathbf{b}\mathbf{x}_2\mathbf{c}, \mathbf{d}) \leftarrow A(\mathbf{x}_1, \mathbf{x}_2)$

non-branching \subseteq well-nested

$A(\varepsilon, \varepsilon) \leftarrow$



$i=0$

$i=1$

$i=2$

$i=3$

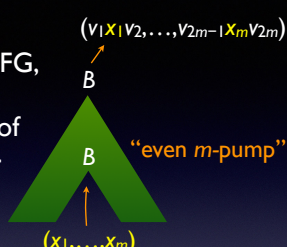
$\{ \varepsilon \} \cup \{ a^{i-1}abc(bdc)^{i-1}d \mid i \geq 1 \}$

A very simple example.

The only choice you can make is the number of times you use the second rule.

Actually 2-iterative, but no straightforward connection between the iterated substrings and parts of derivation trees.

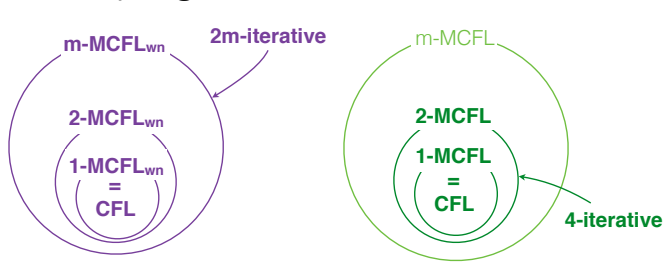
- If G is a well-nested m -MCFG,
 - $\{ T \mid T \text{ is a derivation tree of } G \text{ without even } m\text{-pumps} \}$ may not be finite.
- But there is a well-nested $(m-1)$ -MCFG generating
 - $\{ \text{yield}(T) \mid T \text{ is a derivation tree of } G \text{ without even } m\text{-pumps} \}$.



22

If the derivation tree contains an even m -pump, the string is $2m$ -pumpable. Otherwise, the string is in the language of some w.n. $(m-1)$ -MCFG, and therefore is $2(m-1)$ -pumpable (disregarding finitely many exceptions).
Proof by induction on m .

Pumping Lemma for Subclasses



Kanazawa 2009, by grammar splitting and transformation

What about **Ogden's Lemma**?

23

My proof of the pumping lemma for m -MCFL_{wn} and 2-MCFL is not straightforward.

Ogden's Lemma for CFL

$L \in \text{CFL} \Rightarrow$

$\exists p \forall z \in L$ (at least p positions of z are marked \Rightarrow

$\exists u_1 u_2 u_3 v_1 v_2$ (
 $z = u_1 v_1 u_2 v_2 u_3 \wedge$
 $(u_1, v_1, u_2 \text{ each contain a marked position } \vee$
 $u_2, v_2, u_3 \text{ each contain a marked position}) \wedge$
 $v_1 u_2 v_2 \text{ contains no more than } p$
 $\text{marked positions } \wedge$
 $\forall n \geq 0 (u_1 v_1^n u_2 v_2^n u_3 \in L)$)

Ogden 1968

24

Can be used to show inherent ambiguity of some CFLs, e.g., $\{ a^m b^n c^p \mid m = n \vee n = p \}$.

There are various ways of generalizing Ogden's lemma suitable for MCFGs. At least this much should be implied.

L has the **weak Ogden property** iff

$\exists p \forall z \in L$ (at least p positions of z are marked \Rightarrow

$\exists k \geq 1 \exists u_1 \dots u_{k+1} v_1 \dots v_k$

$z = u_1 v_1 \dots u_k v_k u_{k+1} \wedge$

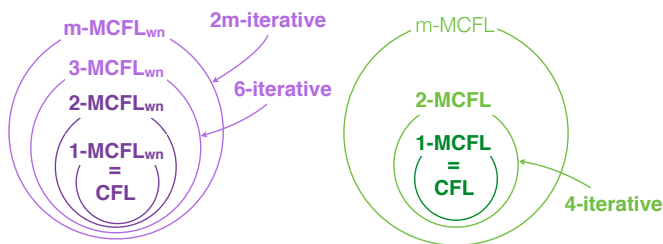
$\exists i (v_i \text{ contains a marked position}) \wedge$

$\forall n \geq 0 (u_1 v_1^n \dots u_k v_k^n u_{k+1} \in L)$

25

This is the first new result in this talk.

The Failure of Ogden's Lemma

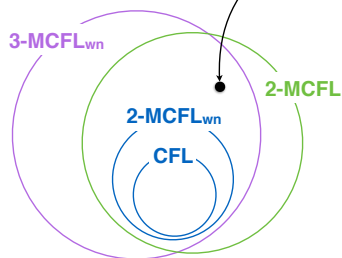


The weak Ogden property **fails** for 3-MCFL_{wn} and 2-MCFL .

26

A language for which the weak Ogden property fails.

$\{ a^{i_1} b^{i_0} a^{i_2} b^{i_1} a^{i_3} b^{i_2} \dots a^{i_n} b^{i_{n-1}} \mid n \geq 3, i_0, \dots, i_n \geq 0 \}$



27

$A(\varepsilon) \leftarrow$ $A(b\mathbf{x}_1) \leftarrow A(\mathbf{x}_1)$ $B(\mathbf{x}_1, \varepsilon) \leftarrow A(\mathbf{x}_1)$ $B(a\mathbf{x}_1, b\mathbf{x}_2) \leftarrow B(\mathbf{x}_1, \mathbf{x}_2)$ $C(\mathbf{x}_1, \mathbf{x}_2, \varepsilon) \leftarrow B(\mathbf{x}_1, \mathbf{x}_2)$ $C(\mathbf{x}_1, a\mathbf{x}_2, b\mathbf{x}_3) \leftarrow C(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3)$ $C(\mathbf{x}_1\$ \mathbf{x}_2, \mathbf{x}_3, \varepsilon) \leftarrow C(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3)$ $D(\mathbf{x}_1\$ \mathbf{x}_2, \mathbf{x}_3) \leftarrow C(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3)$ $D(\mathbf{x}_1, a\mathbf{x}_2) \leftarrow D(\mathbf{x}_1, \mathbf{x}_2)$ $S(\mathbf{x}_1\$ \mathbf{x}_2) \leftarrow D(\mathbf{x}_1, \mathbf{x}_2)$ <p style="text-align: center; color: orange;">non-branching 3-MCFG</p>	$A(\varepsilon) \leftarrow$ $A(b\mathbf{x}_1) \leftarrow A(\mathbf{x}_1)$ $B(\mathbf{x}_1, \varepsilon) \leftarrow A(\mathbf{x}_1)$ $B(a\mathbf{x}_1, b\mathbf{x}_2) \leftarrow B(\mathbf{x}_1, \mathbf{x}_2)$ $C(\varepsilon, \varepsilon) \leftarrow$ $C(a\mathbf{x}_1, b\mathbf{x}_2) \leftarrow C(\mathbf{x}_1, \mathbf{x}_2)$ $D(\mathbf{x}_1\$ \mathbf{y}_1\mathbf{x}_2, \mathbf{y}_2) \leftarrow B(\mathbf{x}_1, \mathbf{x}_2), C(\mathbf{y}_1, \mathbf{y}_2)$ $D(\mathbf{x}_1\$ \mathbf{y}_1\mathbf{x}_2, \mathbf{y}_2) \leftarrow D(\mathbf{x}_1, \mathbf{x}_2), C(\mathbf{y}_1, \mathbf{y}_2)$ $E(\mathbf{x}_1, \mathbf{x}_2) \leftarrow D(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3)$ $E(\mathbf{x}_1, a\mathbf{x}_2) \leftarrow E(\mathbf{x}_1, \mathbf{x}_2)$ $S(\mathbf{x}_1\$ \mathbf{x}_2) \leftarrow E(\mathbf{x}_1, \mathbf{x}_2)$ <p style="text-align: center; color: orange;">2-MCFG</p>
---	--

$\{ a^{i_1}b^{i_0}\$ a^{i_2}b^{i_1}\$ a^{i_3}b^{i_2}\$ \dots \$ a^{i_n}b^{i_{n-1}} \mid n \geq 3, i_0, \dots, i_n \geq 0 \}$

28

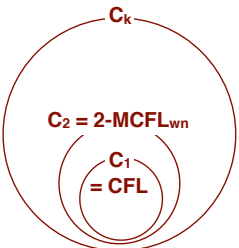
$\{ a^{i_1}b^{i_0}\$ a^{i_2}b^{i_1}\$ a^{i_3}b^{i_2}\$ \dots \$ a^{i_n}b^{i_{n-1}} \mid n \geq 3, i_0, \dots, i_n \geq 0 \}$
 is **2-iterative**

$a\$ a^2b\$ a^3b^2\$ \dots \$ a^{p+1}b^p$

29

Mark the positions of \$.

Weir's (1992) Control Language Hierarchy



$C_k \subseteq 2^{k-1}\text{-MCFL}$
Kanazawa and Salvati 2007

Generalization of Ogden's lemma
Palis and Shende 1995

30

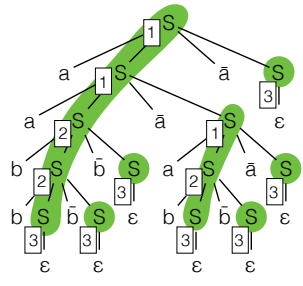
Subclasses of MCFL that are known to have an Ogden property.

Control Grammars

G:
 1: $S \rightarrow aS\bar{a}S$
 2: $S \rightarrow bS\bar{b}S$
 3: $S \rightarrow \epsilon$

CFG with child selection

$K = \{ 1^n 2^n 3 \mid n \geq 0 \}$
 control set



$$L(G, K) = D_2^* n (\{ a^n b^n \mid n \geq 1 \} \bar{b} \{ \bar{a}, \bar{b} \}^*)^*$$

Languages in each level of the control language hierarchy are given by “control grammars”.

Control Language Hierarchy

$$C_1 = \text{CFL}$$

$$C_{k+1} = \{ L(G, K) \mid K \in C_k \}$$

Ogden’s Lemma for C_k

$L \in C_k \Rightarrow$

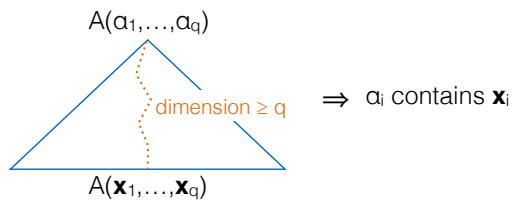
$\exists p \forall z \in L$ (at least p positions of z are marked \Rightarrow
 $\exists u_1 \dots u_{2k+1} v_1 \dots v_{2k}$
 $z = u_1 v_1 \dots u_{2k} v_{2k} u_{2k+1} \wedge$
 $\exists i (u_i, v_i, u_{i+1} \text{ each contain a marked position}) \wedge$
 $v_{2k-1} u_{2k-1} v_{2k-1+1}$ contains no more than p
 marked positions \wedge
 $\forall n \geq 0 (u_1 v_1^n u_2 v_2^n \dots u_{2k} v_{2k}^n u_{2k+1} \in L)$

Palis and Shende 1995

- $k = 1$ gives Ogden’s (1968) original lemma.

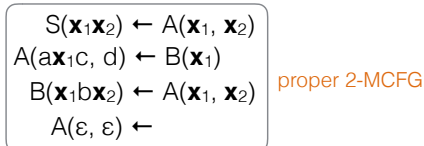
It is quite straightforward to prove an Ogden’s lemma for C_k .

Proper MCFGs

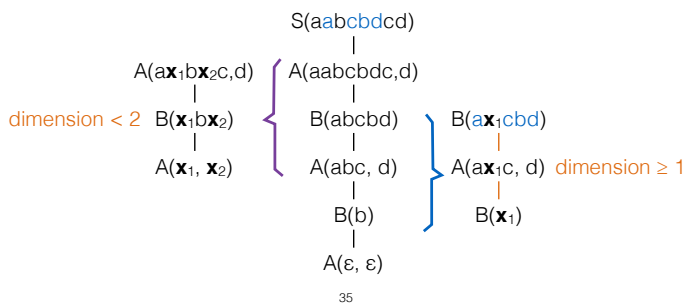


34

Approach this existing result from the MCFG formalism.
Sufficient condition for an Ogden property.



Slight variation of an earlier example.



35

Ogden's Lemma for m -MCFL_{prop}

$L \in m$ -MCFL_{prop} \Rightarrow

$\exists p \forall z \in L$ (at least p positions of z are marked) \Rightarrow

$\exists u_1 \dots u_{2m+1} v_1 \dots v_{2m}$ (
 $Z = u_1 v_1 \dots u_{2m} v_{2m} u_{2m+1} \wedge$
 $\exists i (u_i, v_i, u_{i+1}$ each contain a marked position) \wedge
 $v_1 u_2 v_2 v_3 u_4 v_4, \dots, v_{2m-1} u_{2m} v_{2m}$ together contain
 no more than p marked positions \wedge
 $\forall n \geq 0 (u_1 v_1^n u_2 v_2^n \dots u_{2m} v_{2m}^n u_{2m+1} \in L)$)

- $m = 1$ gives Ogden's (1968) original lemma.

36

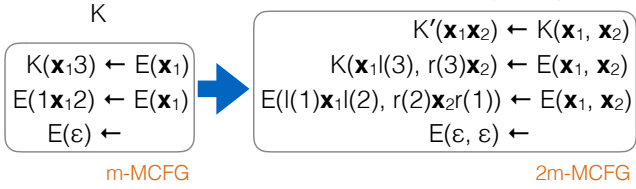
Constrain all of v_1, \dots, v_{2m}

$C_k \subseteq 2^{k-1}\text{-MCFL}_{\text{prop}}$

G: $1: S \rightarrow aS\bar{a}S$ $l(1) = a, r(1) = \bar{a}S$
 $2: S \rightarrow bS\bar{b}S$ $l(2) = b, r(2) = \bar{b}S$ *homomorphisms*
 $3: S \rightarrow \epsilon$ $l(3) = \epsilon, r(3) = \epsilon$

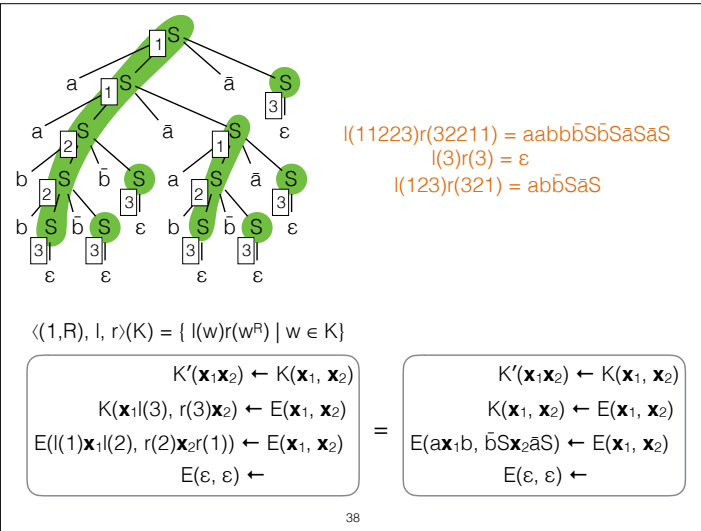
$K = \{ 1^n 2^n 3 \mid n \geq 0 \}$

$\langle (1,R), l, r \rangle(K) = \{ l(w)r(w^R) \mid w \in K \}$
homomorphic replication



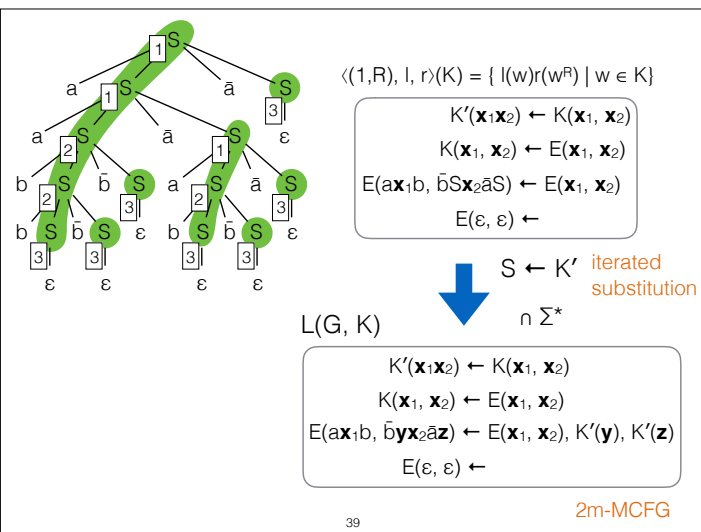
37

The earlier result is subsumed by the present result.



Generates strings with nonterminals.

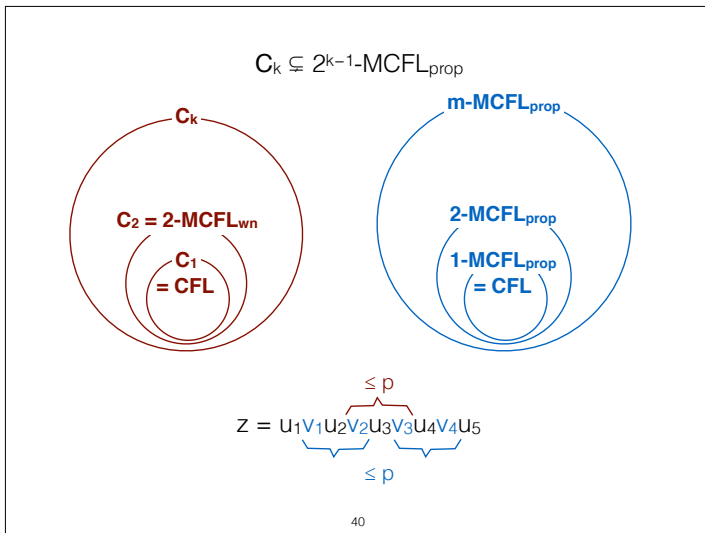
38



The construction doubles the dimension, preserves properness.

39

The different requirement shows the properness of the inclusion.



Summary

- Pumping doesn't imply Ogden: There is no Ogden-like theorem for $3\text{-MCFL}_{\text{wn}} \cap 2\text{-MCFL}$
- There is a natural Ogden's lemma for $m\text{-MCFL}_{\text{prop}}$
- Covers Weir's control language hierarchy

41

Reference: Makoto Kanazawa. 2016. Ogden's lemma, multiple context-free grammars, and the control language hierarchy. In Adrian-Horia Dediu, Jan Janoušek, Carlos Martín-Vide, and Bianca Truthe, editors, Language and Automata Theory and Applications: 10th International Conference, LATA 2016, pages 371-383. Lecture Notes in Computer Science 9618. Cham: Springer.