

Ogden's Lemma, Multiple Context-Free Grammars, and the Control Language Hierarchy

Makoto Kanazawa
National Institute of Informatics and SOKENDAI
Japan

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)

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.

Which properties of CFGs are shared by/
generalize to MCFGs?

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

5

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

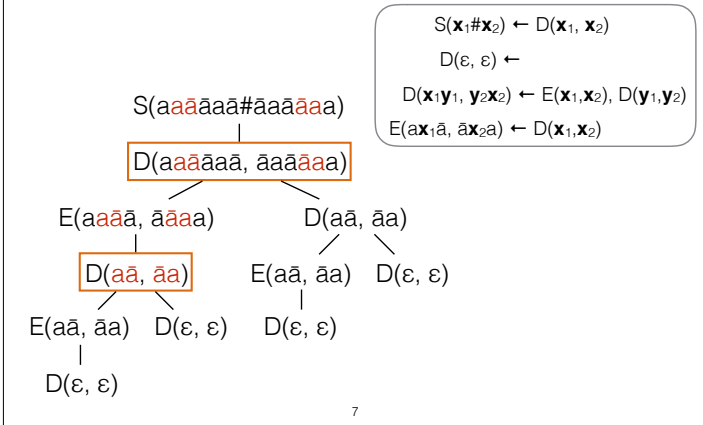
$\{ a_1^n a_2^n \dots a_{2m-1}^n a_{2m}^n \mid n \geq 0 \}$

Seki et al. 1991

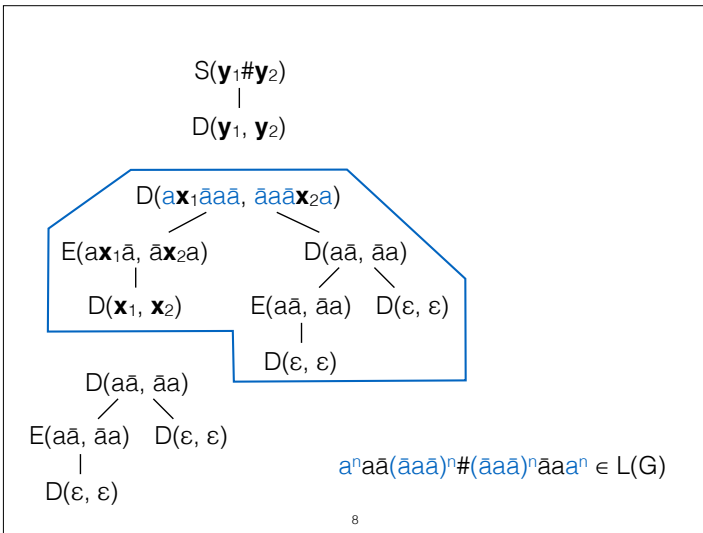
6

The languages of MCFGs form an infinite hierarchy.

Pumping



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 $\mathbf{x}_1, \mathbf{x}_2, \mathbf{y}_1, \mathbf{y}_2$ are instantiated by ...
 The number of iterated substrings (factors) larger than two.

Iterative Properties

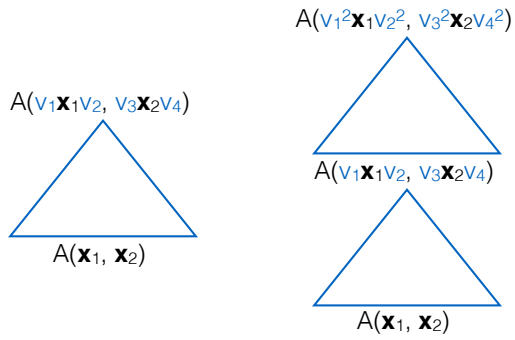
L is **k-iterative** iff $\exists p \forall z \in L (|z| \geq p \Rightarrow$
 $\exists u_1 \dots u_{k+1} \forall v_1 \dots v_k ($
 $z = u_1 v_1 \dots u_k v_k u_{k+1} \wedge$
 $v_1 \dots v_k \neq \varepsilon \wedge$
 $\forall n \geq 0 (u_1 v_1^n \dots u_k v_k^n u_{k+1} \in L))$

$L \in \text{CFL} \Rightarrow L$ is 2-iterative

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

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

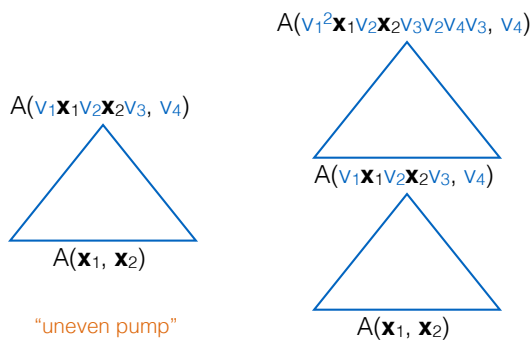
Difficulty with Pumping



10

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

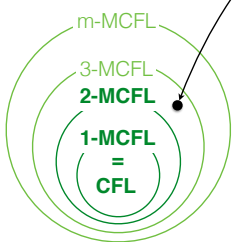
Difficulty with Pumping



11

Or like this.

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



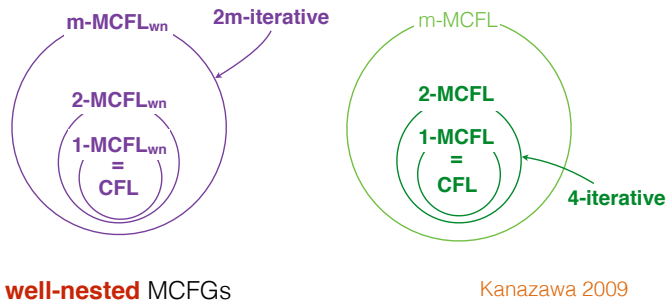
not k -iterative for any k

Kanazawa et al. 2014

12

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

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

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

$S(x_1\#x_2) \leftarrow D(x_1, x_2)$
 $D(\varepsilon, \varepsilon) \leftarrow$
 $D(x_1y_1, y_2x_2) \leftarrow E(x_1, x_2), D(y_1, y_2)$
 $E(ax_1\bar{a}, \bar{a}x_2a) \leftarrow D(x_1, x_2)$

$S(x_1\#x_2) \leftarrow D(x_1, x_2)$
 $D(\varepsilon, \varepsilon) \leftarrow$
 $D(x_1y_1, x_2y_2) \leftarrow E(x_1, x_2), D(y_1, y_2)$
 $E(ax_1\bar{a}, ax_2\bar{a}) \leftarrow D(x_1, x_2)$

well-nested

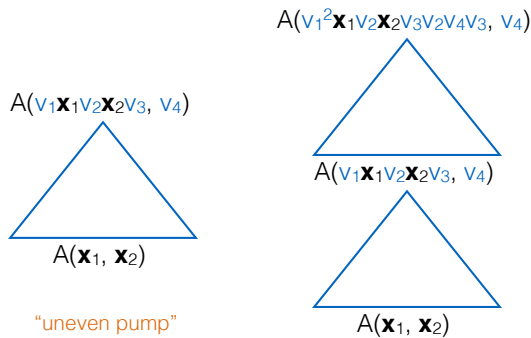
non-well-nested

$$\{ w\#w \mid w \in D_1^* \} \notin MCFL_{wn}$$

Kanazawa and Salvati 2010

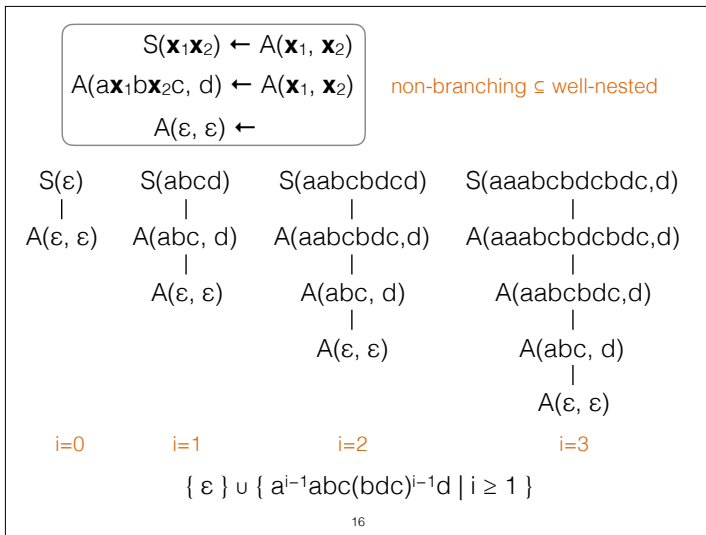
Has a natural equivalent
characterization: $yCFT_{sp}$

Difficulty with Pumping



"uneven pump"

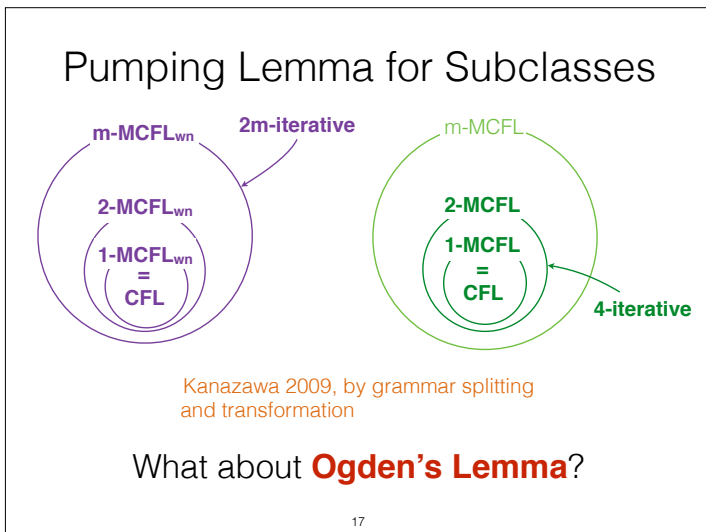
Pumping not easy to prove even for
well-nested MCFGs: this situation can
still arise.



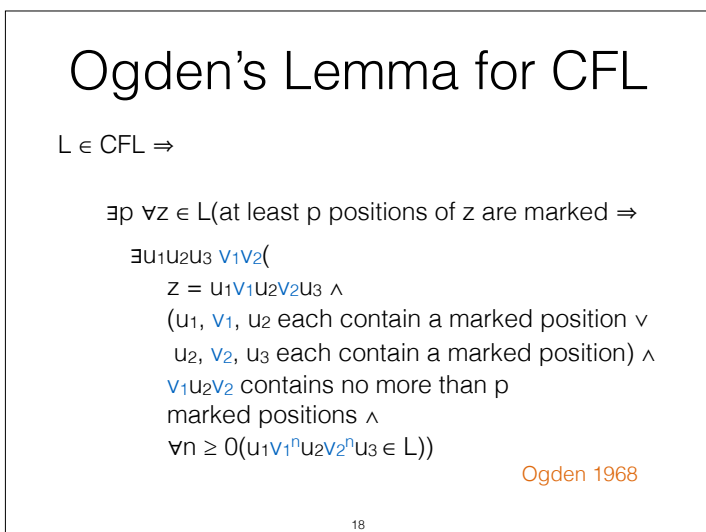
A very simple example.

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.



My proof of the pumping lemma for $m\text{-MCFL}_{wn}$ and 2-MCFL is not straightforward.



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

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

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

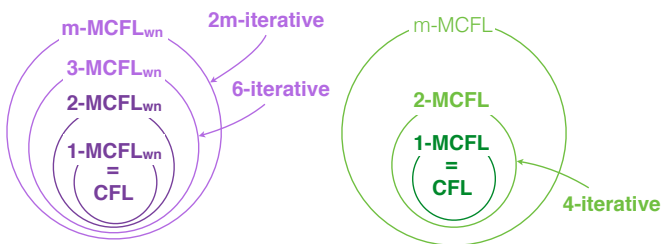
$\exists i (v_i \text{ contains a marked position}) \wedge$

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

19

This is the first new result in this talk.

The Failure of Ogden's Lemma

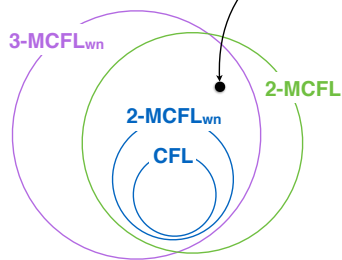


The weak Ogden property **fails** for 3-MCFL_{wn} and 2-MCFL .

20

A language for which the weak Ogden property fails.

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



21

$A(\varepsilon) \leftarrow$ $A(bx_1) \leftarrow A(x_1)$ $B(x_1, \varepsilon) \leftarrow A(x_1)$ $B(ax_1, bx_2) \leftarrow B(x_1, x_2)$ $C(x_1, x_2, \varepsilon) \leftarrow B(x_1, x_2)$ $C(x_1, ax_2, bx_3) \leftarrow C(x_1, x_2, x_3)$ $C(x_1 \$, x_2, x_3, \varepsilon) \leftarrow C(x_1, x_2, x_3)$ $D(x_1 \$, x_2, x_3) \leftarrow C(x_1, x_2, x_3)$ $D(x_1, ax_2) \leftarrow D(x_1, x_2)$ $S(x_1 \$, x_2) \leftarrow D(x_1, x_2)$ <p style="text-align: center; color: orange;">non-branching 3-MCFG</p>	$A(\varepsilon) \leftarrow$ $A(bx_1) \leftarrow A(x_1)$ $B(x_1, \varepsilon) \leftarrow A(x_1)$ $B(ax_1, bx_2) \leftarrow B(x_1, x_2)$ $C(\varepsilon, \varepsilon) \leftarrow$ $C(ax_1, bx_2) \leftarrow C(x_1, x_2)$ $D(x_1 \$, y_1 x_2, y_2) \leftarrow B(x_1, x_2), C(y_1, y_2)$ $D(x_1 \$, y_1 x_2, y_2) \leftarrow D(x_1, x_2), C(y_1, y_2)$ $E(x_1, x_2) \leftarrow D(x_1, x_2, x_3)$ $E(x_1, ax_2) \leftarrow E(x_1, x_2)$ $S(x_1 \$, x_2) \leftarrow E(x_1, x_2)$ <p style="text-align: center; color: orange;">2-MCFG</p>
---	---

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

22

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

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

23

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

24

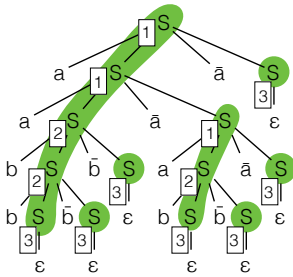
Subclasses of MCFL that are known to have an Ogden property.

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



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

25

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

Control Language Hierarchy

$$C_1 = \text{CFL}$$

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

26

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

$$\exists u_1 \dots u_{2k+1} v_1 \dots v_{2k} ($$

$$z = u_1 v_1 \dots u_{2k} v_{2k} u_{2k+1} \wedge$$

$$\exists i (u_i, v_i, u_{i+1} \text{ each contain a marked position}) \wedge$$

$$v_{2k-1} u_{2k-1} v_{2k-1+1} \text{ contains no more than } p$$

$$\text{marked positions} \wedge$$

$$\forall n \geq 0 (u_1 v_1^n u_2 v_2^n \dots u_{2k} v_{2k}^n u_{2k+1} \in L)$$

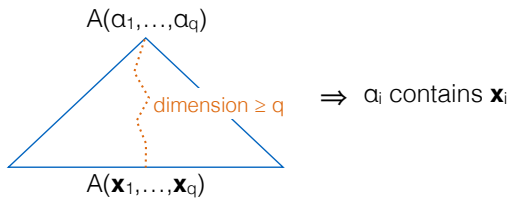
Palis and Shende 1995

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

27

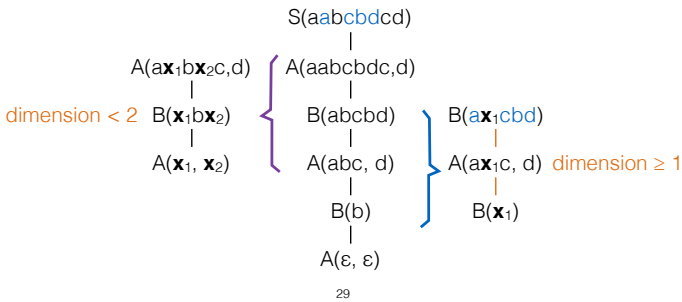
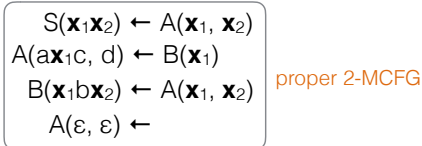
It is quite straightforward to prove an Ogden's lemma for C_k .

Proper MCFGs



28

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



29

Slight variation of an earlier example.

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

$\exists u_1 \dots u_{2m+1} v_1 \dots v_{2m}$ (
 $Z = u_1 v_1 \dots u_{2m} v_{2m} u_{2m+1} \wedge$
 $\exists i (u_i, v_i, u_{i+1} \text{ each contain a marked position}) \wedge$
 $v_1 u_2 v_2 v_3 u_4 v_4 \dots v_{2m-1} u_{2m} v_{2m}$ 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.

30

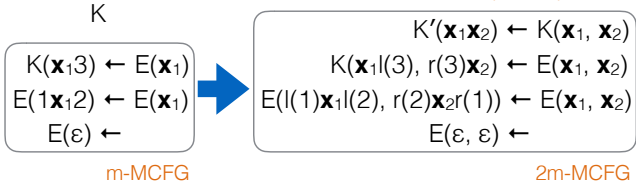
Constrain all of v_1, \dots, v_{2m}

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

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*
 $3: S \rightarrow \epsilon$ $l(3) = \epsilon, r(3) = \epsilon$

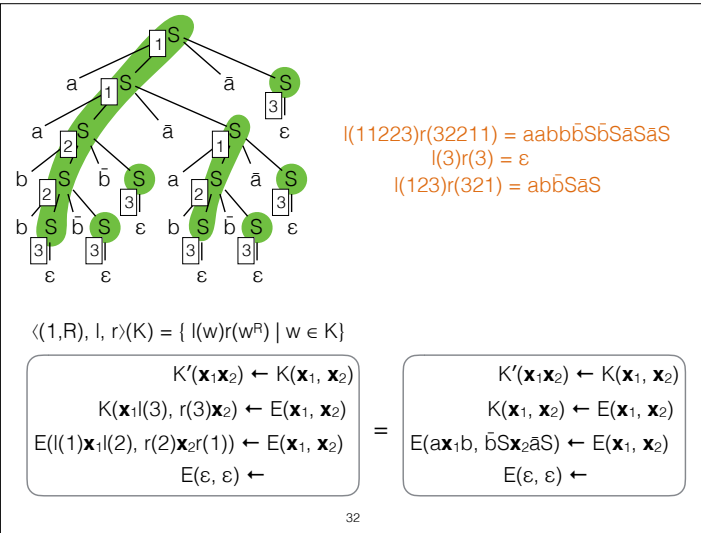
$K = \{ 1^n 2^n 3 \mid n \geq 0 \}$

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



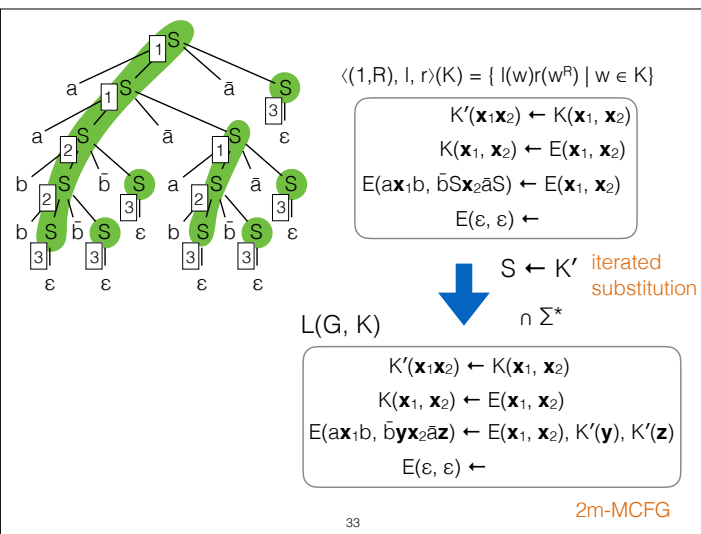
31

The earlier result is subsumed by the present result.



Generates strings with nonterminals.

32

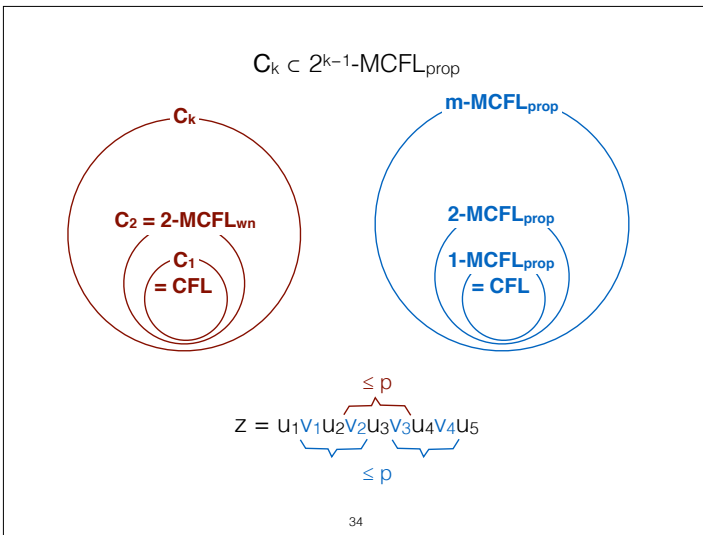


The construction doubles the dimension, preserves properness.

33

2m-MCFG

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