# MIX Is Not a Tree-Adjoining Language

Makoto Kanazawa NII,Tokyo



Sylvain Salvati LaBRI/INRIA Bordeaux

## $MIX = \{ w \in \{a,b,c\}^* \mid |w|_a = |w|_b = |w|_c \}$

This talk is about one little language called MIX. It consists of all strings over {a,b,c} that contain the same number of occurrences of each letter.

# $\mathsf{MIX}\not\in\mathsf{TAL}$



We prove that MIX is not a tree-adjoining language.

# Joshi's Conjecture

"TAGs cannot generate this language, although for TAGs the proof is not in hand yet."

Joshi 1985



4

The non-membership of MIX in TAL was conjectured by Joshi almost 30 years ago. He was confident at the time. "It is not known whether TAG ... can generate MIX. This has turned out to be a very difficult problem."

#### Joshi, Vijay-Shanker, and Weir 1991



Six years later, the confidence seemed to have faded.



Of course, the interest in the question comes from the fact that it was not obvious how to prove it.

The pumping lemma for TAL almost certainly cannot be used to prove that MIX is not in TAL.

6

7

MIX and Free Word Order

"[No human language] has ... complete freedom for order." Bach 1981

"[MIX represents] an extremely case of the degree of free word order permitted in a language ... which is linguistically not relevant." Joshi 1985

"... it seems rather unlikely that any natural language will turn out to have a MIX-like characteristic." Gazdar 1985

A couple of respects in which MIX is not like a natural language. 1. Completely free word order.







$$MIX = \{ w \in \{a,b,c\}^* \mid |w|_a = |w|_b = |w|_c \}$$

$$MIX_{k} = \{ w \in \{a_{1}, \dots, a_{k}\}^{*} \mid |w|_{a_{1}} = \dots = |w|_{a_{k}} \}$$

## $\forall k(MIX_k \in TAL) \Rightarrow \forall L(L \in TAL \Rightarrow Pem(L) \in TAL)$ permutation closure

Generalization of MIX.

For TAL, the premise is of course false, but any family of languages that is a "rational cone" or "full trio" and is included in the class of semilinear languages has this property.

8

# MIX Is Not Mildly CS?

"MCSGs capture only certain kinds of dependencies, such as nested dependencies and certain limited kinds of crossing dependencies (for example, in subordinate clause constructions in Dutch or some variations of them, but perhaps not in the so-called MIX ... language ...)"

Joshi, Vijay-Shanker, and Weir 1991



2. Exhibits "unlimited (?) cross-serial dependencies". Joshi et al. 1991 suggested MIX should be excluded from the class of mildly contextsensitive languages because of this property.

# Mildly Context-Sensitive Languages

- Polynomial-time recognition
- Semilinearity/constant growth
- Exhibits limited cross-serial dependencies

Joshi 1985

Three defining conditions of mild context-sensitivity.



### Sylvain Salvati 2011 $MIX \in 2-MCFL$ 2-MCFL • a<sup>n</sup>b<sup>n</sup>c<sup>n</sup>d<sup>n</sup>e<sup>n</sup> ΜΙΧ ΓAL ■ a<sup>n</sup>b<sup>n</sup>c<sup>n</sup>d<sup>n</sup> • *a<sup>n</sup>b<sup>n</sup>c<sup>n</sup>*

Surprisingly, MIX is the language of a 2-multiple context-free grammar or linear context-free rewriting system of fanout 2.

• *a*<sup>*n*</sup>*b*<sup>*n*</sup>

Very difficult proof using algebraic topology.

# MIX ∉ TAL





We use the formalism of head grammar, known to be equivalent to TAG.

# Head Grammar



#### Pollard 1984

 $\begin{aligned} A(\mathbf{x}_1, \mathbf{x}_2, \mathbf{y}_1, \mathbf{y}_2) &\leftarrow B(\mathbf{x}_1, \mathbf{x}_2), C(\mathbf{y}_1, \mathbf{y}_2) & \text{left} \\ A(\mathbf{x}_1, \mathbf{x}_2, \mathbf{y}_1, \mathbf{y}_2) &\leftarrow B(\mathbf{x}_1, \mathbf{x}_2), C(\mathbf{y}_1, \mathbf{y}_2) & \text{right} \\ A(\mathbf{x}_1, \mathbf{y}_1, \mathbf{y}_2, \mathbf{x}_2) &\leftarrow B(\mathbf{x}_1, \mathbf{x}_2), C(\mathbf{y}_1, \mathbf{y}_2) \\ A(w_1, w_2) &\leftarrow (w_i \in \Sigma \cup \{ \epsilon \}) \end{aligned}$ 

left concatenation right concatenation wrapping

$$L(G) = \{ w_1 w_2 \mid \vdash_G S(w_1, w_2) \}$$

#### special kind of 2-MCFG or LCFRS(2)

Introduced by Pollard, modified by Roach. Nonterminals stand for binary relations on strings. A grammar is a kind of logic program on strings. Just three operations on pairs of strings.



An example of a head grammar (kind of synchronous CFG). All non-terminating rules are wrapping rules. 14

$$S(\mathbf{x}_{1}\mathbf{y}_{1},\mathbf{y}_{2}\mathbf{x}_{2}) \leftarrow D(\mathbf{x}_{1},\mathbf{x}_{2}), C(\mathbf{y}_{1},\mathbf{y}_{2})$$

$$C(\varepsilon, \#) \leftarrow$$

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

$$I(\mathbf{x}_{1}\mathbf{y}_{1},\mathbf{y}_{2}\mathbf{x}_{2}) \leftarrow F(\mathbf{x}_{1},\mathbf{x}_{2}), D(\mathbf{y}_{1},\mathbf{y}_{2})$$

$$F(\mathbf{x}_{1}\mathbf{y}_{1},\mathbf{y}_{2}\mathbf{x}_{2}) \leftarrow A(\mathbf{x}_{1},\mathbf{x}_{2}), E(\mathbf{y}_{1},\mathbf{y}_{2})$$

$$A(a, a) \leftarrow$$

$$E(\mathbf{x}_{1}\mathbf{y}_{1},\mathbf{y}_{2}\mathbf{x}_{2}) \leftarrow D(\mathbf{x}_{1},\mathbf{x}_{2}), A'(\mathbf{y}_{1},\mathbf{y}_{2})$$

$$A'(\bar{a},\bar{a}) \leftarrow$$

$$HG \equiv well-nested 2-MCFG$$

$$S(\mathbf{x}_{1}\mathbf{y}_{1}, \mathbf{x}_{2}\mathbf{y}_{2}) \leftarrow D(\mathbf{x}_{1}, \mathbf{x}_{2}), C(\mathbf{y}_{1}, \mathbf{y}_{2})$$

$$C(\varepsilon, \#) \leftarrow$$

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

$$D(\mathbf{x}_{1}\mathbf{y}_{1}, \mathbf{x}_{2}\mathbf{y}_{2}) \leftarrow F(\mathbf{x}_{1}, \mathbf{x}_{2}), D(\mathbf{y}_{1}, \mathbf{y}_{2})$$

$$F(\mathbf{x}_{1}\mathbf{y}_{1}, \mathbf{x}_{2}\mathbf{y}_{2}) \leftarrow A(\mathbf{x}_{1}, \mathbf{x}_{2}), E(\mathbf{y}_{1}, \mathbf{y}_{2})$$

$$A(a, a) \leftarrow$$

$$E(\mathbf{x}_{1}\mathbf{y}_{1}, \mathbf{x}_{2}\mathbf{y}_{2}) \leftarrow D(\mathbf{x}_{1}, \mathbf{x}_{2}), A'(\mathbf{y}_{1}, \mathbf{y}_{2})$$

$$A'(\bar{a}, \bar{a}) \leftarrow$$

 $\left\{ w \# w \mid w \in D_{\mathsf{I}} \right\}$ 

non-well-nested 2-MCFG

String operations used in HG are "well-nested".



#### $a^{5}b^{14}a^{19}c^{29}b^{15}a^{5}$ has no 2-decomposition.

Outline of the proof



A decomposition is similar to an HG derivation but independent of any grammar. Any HG derivation can be turned into a decomposition by stripping off nonterminals.

#### decomposition of w



Formally, a decomposition is a binary tree whose nodes are labeled by pairs of strings. Every HG derivation of w gives a decomposition of w.



19

The parameter n measures how unbalanced the occurrence counts of the three letters can be at any node in a decomposition.



$$w \in MIX \Leftrightarrow \psi(w) = (0, 0)$$

 $u_0v_1u_1v_2u_2, u_0v_1'u_1v_2'u_2 \in MIX \Rightarrow \psi(v_1v_2) = \psi(v_1'v_2')$ 

Properties of the function  $\boldsymbol{\psi}$ 



#### $a^{5}b^{14}a^{19}c^{29}b^{15}a^{5}$ has no 2-decomposition.

Prove the first implication.

Suppose L(G) = MIX.



### $\vdash_{G} A(v_1, v_2), \vdash_{G} A(v_1', v_2') \Rightarrow \Psi(v_1v_2) = \Psi(v_1'v_2')$

Whenever the same nonterminal holds of two pairs of strings, their  $\psi$  value is the same.

Let  $n = \max\{ \|\Psi(v_1v_2)\|_{\infty} \mid \vdash_G A(v_1, v_2) \text{ for some } A \}.$  $\{ \Psi(v_1v_2) \mid \vdash_G A(v_1, v_2) \text{ for some } A \} \subseteq [-n, n] \times [-n, n]$ 



A head grammar G determines the value of n. Every decomposition that comes from a derivation according to G is an n-decomposition.



#### $\forall w \in MIX(w \text{ has } 2\text{-decomposition})$

#### $a^{5}b^{14}a^{19}c^{29}b^{15}a^{5}$ has no 2-decomposition.

The first implication done. Now prove the second implication.







The proof uses a homomorphism that repeats each letter n times.



2*n*-decomposition of  $\gamma_n(w)$ 



The goal is to show that  $w \in MIX$  has a 2-decomposition.

Start with an n-decomposition of  $\gamma_n(w)$ .

Turn it into a "neat" decomposition.

Primed strings start and end at block boundaries.

26



 $\psi_i(u_1') - \psi_i(u_1) \in [-n+1, n-1] \quad \psi_i(u_2') - \psi_i(u_2) \in [-n+1, n-1]$ 

The "unbalancedness" may increase after the transformation.





 $\Psi(u_1'u_2') = (pn, qn)$ 

 $\Psi(u_1'u_2') \in [-2n, 2n] \times [-2n, 2n]$ 



2-decomposition of w



Invert the homomorphism  $\gamma_n$  and obatin a 2-decomposition of w.



#### $\forall w \in MIX(w \text{ has } 2\text{-decomposition})$

#### $a^{5}b^{14}a^{19}c^{29}b^{15}a^{5}$ has no 2-decomposition.

#### $a^{5}b^{14}a^{19}c^{29}b^{15}a^{5}$ has no 2-decomposition.



computer verification

mathematical proof

Can show this in two ways. Computer program is available from my web site. There's no easy way to explain the mathematical proof.

### $a^{5}b^{14}a^{18}c^{28}b^{14}a^{5}$ has a 2-decomposition.



We believe our counterexample is minimal. A slight modification leads to a 2-decomposable string.

#### $a^{5}b^{14}a^{18}c^{28}b^{14}a^{5}$ has a 2-decomposition.



The "slope" too close to the corner to the right of the origin.

33

### $a^{5}b^{13}a^{18}c^{28}b^{15}a^{5}$ has a 2-decomposition.



Another simplification.

## $a^{5}b^{14}a^{18}c^{28}b^{14}a^{4}$ has a 2-decomposition.



Yet another simplification.



Now we know much more about MIX than before.



MIX belongs to 2-MCFL - TAL. There are two other such languages. RESP belongs to MCFL<sub>wn</sub>. What about MIX? 37