

# **Learning Context-Free Grammars from Positive Data and Membership Queries**

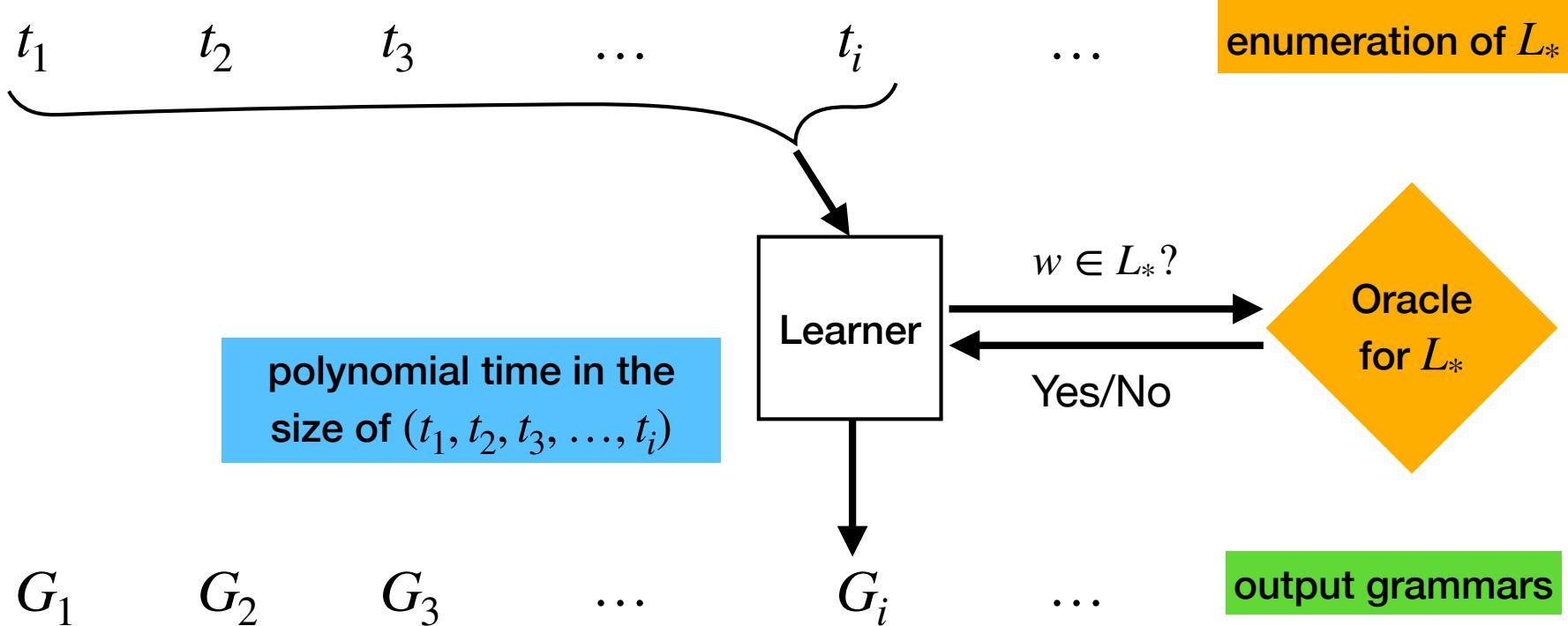
**(based on joint work with Ryo Yoshinaka)**

**Makoto Kanazawa, Hosei University**

# Background

- Algorithmic Learning Theory
  - \* **Grammatical Inference**
    - ▶ Finite automata
      - satisfactory learning algorithms
    - ▶ Context-free grammars / pushdown automata
      - special subclasses
        - deterministic one-counter machines
        - context-deterministic grammars

# Learning from Positive Data and Membership Queries



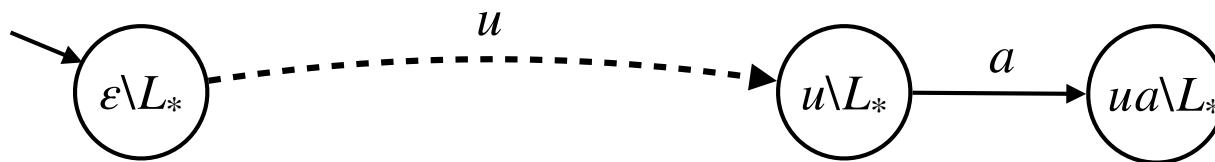
- There exists  $l$  such that  $G_l = G_{l+1} = \dots$  and  $L(G_l) = L_*$ .
- No “delaying trick”.

# Regular Languages

- **Left quotient** of a language  $L \subseteq \Sigma^*$  by a string  $u \in \Sigma^*$ :

$$u \setminus L = \{ x \in \Sigma^* \mid ux \in L \}$$

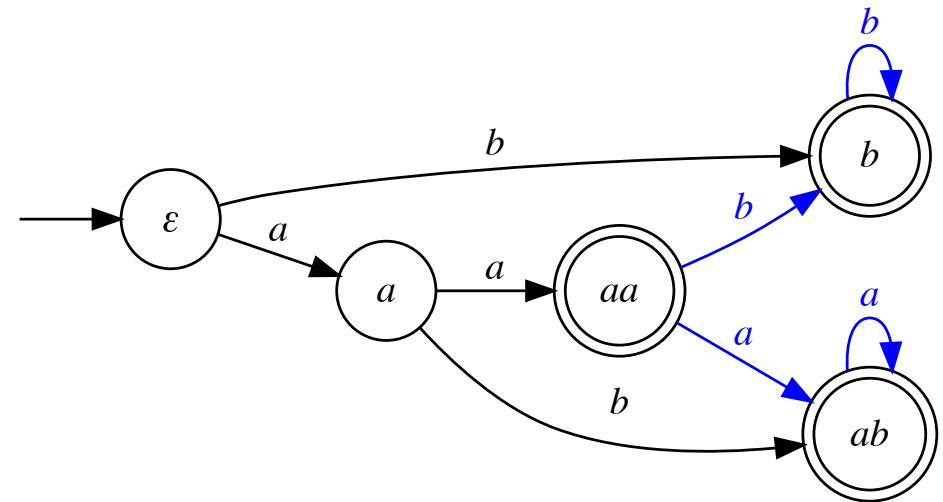
- A language  $L$  is regular if and only if  $\{ u \setminus L \mid u \in \Sigma^* \}$  is finite.
- Every regular language has a canonical **minimal DFA**.



states = left quotients

# Example

$$L_* = aba^* \cup bb^* \cup aa(a^* \cup b^*)$$



$$\varepsilon \setminus L_* = L_*$$

$$a \setminus L_* = ba^* \cup a(a^* \cup b^*)$$

$$b \setminus L_* = b^*$$

$$aa \setminus L_* = a^* \cup b^*$$

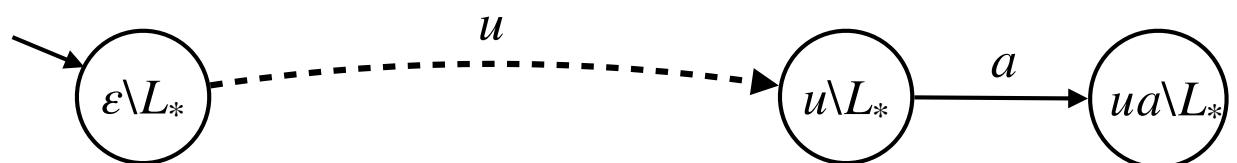
$$ab \setminus L_* = a^*$$

$$b \textcolor{blue}{b} \setminus L_* = b^* = b \setminus L_*$$

$$aaa \textcolor{blue}{a} \setminus L_* = a^* = ab \setminus L_*$$

$$aab \textcolor{blue}{b} \setminus L_* = b^* = b \setminus L_*$$

$$aba \textcolor{blue}{a} \setminus L_* = a^* = ab \setminus L_*$$



states = left quotients

# Right-Linear Grammars

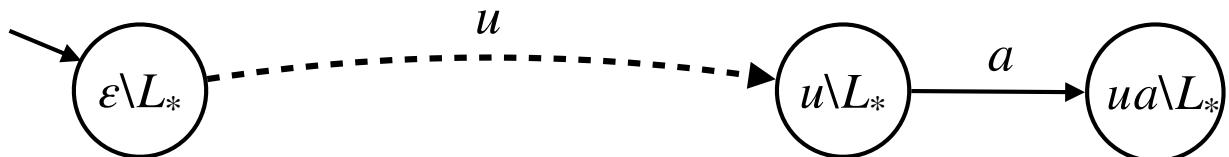
- Right-linear CFG  $G_*$  corresponding to a minimal DFA for  $L_*$ :

$$G_* = (N_*, \Sigma, P_*, S)$$

$$N_* = \{ u \setminus L_* \mid u \in \Sigma^* \}$$

$$P_* = \{ u \setminus L_* \rightarrow a (ua \setminus L_*) \mid u \in \Sigma^*, a \in \Sigma \} \cup \{ u \setminus L_* \rightarrow \varepsilon \mid u \in L_* \}$$

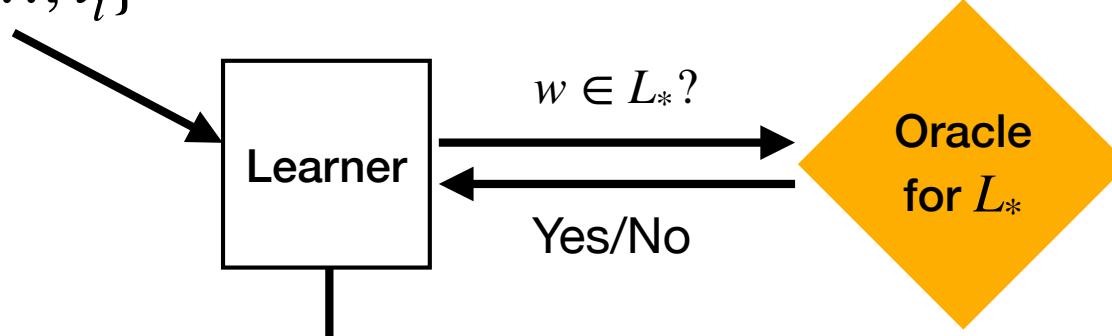
$$S = L_* \quad (= \varepsilon \setminus L_*)$$



states = nonterminals = left quotients

# Inference of Regular Languages

$$T = \{t_1, \dots, t_i\}$$



$$\begin{aligned} \text{Pref}(T) &= \{ u \mid uv \in T \} \\ \text{Suff}(T) &= \{ v \mid uv \in T \} \end{aligned}$$

$$G = (N, \Sigma, P, \langle\langle \varepsilon \rangle\rangle)$$

represents  $u \setminus L_*$

length-lexicographic order on  $\Sigma^*$

$$N = \{ \langle\langle u \rangle\rangle \mid u \in \text{Pref}(T) \wedge \forall v \in \text{Pref}(T) (v < u \rightarrow (\{\varepsilon\} \cup \Sigma) \text{Suff}(T) \cap (v \setminus L_*) \neq (\{\varepsilon\} \cup \Sigma) \text{Suff}(T) \cap (u \setminus L_*) ) \}$$

approximates  $v \setminus L_* \neq u \setminus L_*$

$$P = \{ \langle\langle u \rangle\rangle \rightarrow a \langle\langle v \rangle\rangle \mid \text{Suff}(T) \cap (ua \setminus L_*) = \text{Suff}(T) \cap (v \setminus L_*) \} \cup$$

$$\{ \langle\langle u \rangle\rangle \rightarrow \varepsilon \mid u \in L_* \}$$

approximates  $ua \setminus L_* = v \setminus L_*$

$$N_* = \{ u \setminus L_* \mid u \in \Sigma^* \}$$

represents  $u \setminus L_*$

lexicographic order on  $\Sigma^*$

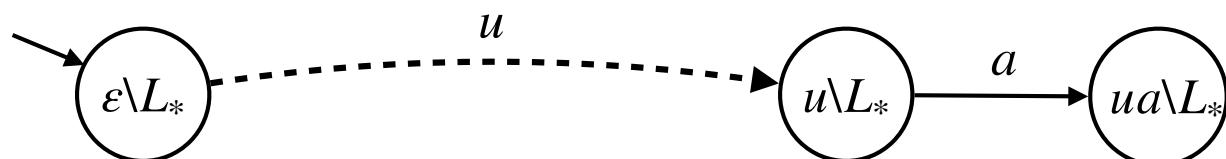
$$N = \{ \langle\langle u \rangle\rangle \mid u \in \text{Pref}(T) \wedge \forall v \in \text{Pref}(T) (v < u \rightarrow (\{\epsilon\} \cup \Sigma) \text{Suff}(T) \cap (v \setminus L_*) \neq (\{\epsilon\} \cup \Sigma) \text{Suff}(T) \cap (u \setminus L_*) \}$$

approximates  $v \setminus L_* \neq u \setminus L_*$

$$\begin{aligned} P_* &= \{ u \setminus L_* \rightarrow a (ua \setminus L_*) \mid u \in \Sigma^*, a \in \Sigma \} \cup \{ u \setminus L_* \rightarrow \epsilon \mid u \in L_* \} \\ &= \{ u \setminus L_* \rightarrow a (v \setminus L_*) \mid ua \setminus L_* = v \setminus L_* \} \cup \{ u \setminus L_* \rightarrow \epsilon \mid u \in L_* \} \end{aligned}$$

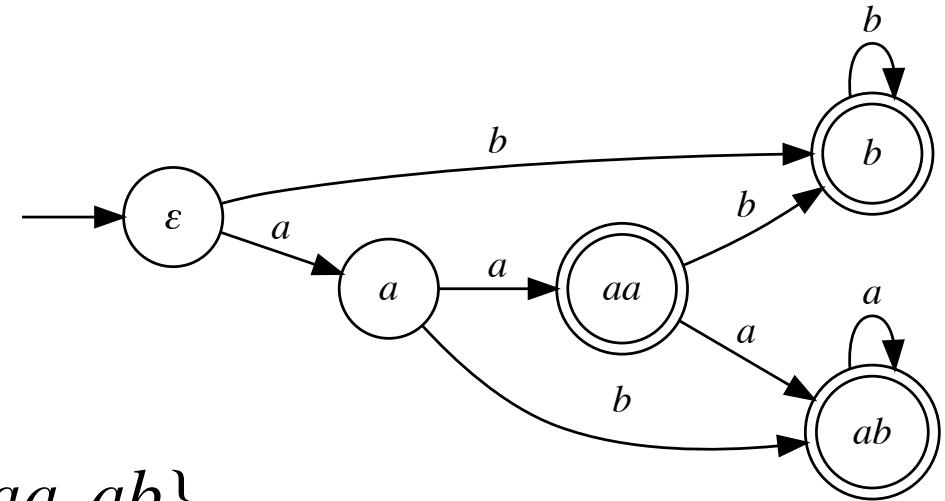
$$\begin{aligned} P &= \{ \langle\langle u \rangle\rangle \rightarrow a \langle\langle v \rangle\rangle \mid \text{Suff}(T) \cap (ua \setminus L_*) = \text{Suff}(T) \cap (v \setminus L_*) \} \cup \\ &\quad \{ \langle\langle u \rangle\rangle \rightarrow \epsilon \mid u \in L_* \} \end{aligned}$$

approximates  $ua \setminus L_* = v \setminus L_*$



# Example

$$L_* = aba^* \cup bb^* \cup aa(a^* \cup b^*)$$



$$T = \{b, aa, ab\}$$

$$\text{Pref}(T) = \{\varepsilon, a, b, aa, ab\}$$

$$\text{Suff}(T) = \{\varepsilon, a, b, aa, ab\}$$

$$\{\varepsilon, a, b\} \text{ Suff}(T) \cap (\varepsilon \setminus L_*) = \{b, aa, ab, bb, aaa, aab\}$$

$$\{\varepsilon, a, b\} \text{ Suff}(T) \cap (a \setminus L_*) = \{a, b, aa, ab, ba, aaa, baa\}$$

$$\{\varepsilon, a, b\} \text{ Suff}(T) \cap (b \setminus L_*) = \{\varepsilon, b, bb\}$$

$$\{\varepsilon, a, b\} \text{ Suff}(T) \cap (aa \setminus L_*) = \{\varepsilon, a, b, aa, aaa, bb\}$$

$$\{\varepsilon, a, b\} \text{ Suff}(T) \cap (ab \setminus L_*) = \{\varepsilon, a, aa, aaa\}$$

represents  $u \setminus L_*$

lexicographic order on  $\Sigma^*$

$$N = \{ \langle\langle u \rangle\rangle \mid u \in \text{Pref}(T) \wedge \forall v \in \text{Pref}(T) (v < u \rightarrow (\{\epsilon\} \cup \Sigma) \text{Suff}(T) \cap (v \setminus L_*) \neq (\{\epsilon\} \cup \Sigma) \text{Suff}(T) \cap (u \setminus L_*)) \}$$

approximates  $v \setminus L_* \neq u \setminus L_*$

$$P = \{ \langle\langle u \rangle\rangle \rightarrow a \langle\langle v \rangle\rangle \mid \text{Suff}(T) \cap (ua \setminus L_*) = \text{Suff}(T) \cap (v \setminus L_*) \} \cup \{ \langle\langle u \rangle\rangle \rightarrow \epsilon \mid u \in L_* \}$$

approximates  $ua \setminus L_* = v \setminus L_*$

$$x \in (\{\epsilon\} \cup \Sigma) \text{Suff}(T) \cap (u \setminus L_*) \iff x \in (\{\epsilon\} \cup \Sigma) \text{Suff}(T) \wedge ux \in L_*$$

- $N$  and  $P$  are determined by  $\text{Pref}(T)$   $(\{\epsilon\} \cup \Sigma) \text{Suff}(T) \cap L_*$ .
- $O(n^4)$  queries to the oracle for  $L_*$  suffice.

# Context-Free Grammars

$$G = (N, \Sigma, P, S)$$

- What do nonterminals correspond to?

$$\frac{\text{left quotients}}{\text{regular languages}} = \frac{\text{??}}{\text{context-free languages}}$$

- The set of terminal strings derived from a nonterminal is included in some **quotient** of the language of the grammar:

$$S \Rightarrow_G^* uAv \quad \text{implies} \quad L_G(A) \subseteq u \setminus L(G) / v = \{ x \in \Sigma^* \mid uxv \in L(G) \}$$

$$\begin{aligned} L_G(A) &= \{ x \in \Sigma^* \mid A \Rightarrow_G^* x \} \\ L(G) &= L_G(S) \end{aligned}$$

- It seems there's nothing further that can be said in general.

# Simple Case: Grammars with Just One Nonterminal

$$\begin{array}{l} S \rightarrow \varepsilon \\ S \rightarrow aSbS \end{array}$$

Dyck language

$$D_1 = \{ x \in \{a, b\}^* \mid |x|_a = |x|_b \wedge \forall uv(uv = x \rightarrow |u|_a \geq |u|_b) \}$$

$$\pi: \quad S \rightarrow w_0 S w_1 \dots S w_k \quad (w_i \in \Sigma^*)$$

When should  $\pi$  be in the hypothesized grammar?

$$\pi \text{ is } \mathbf{valid} \stackrel{\text{def}}{\iff} L_* \supseteq w_0 L_* w_1 \dots L_* w_k$$

Why is this reasonable?

$$\pi: \quad S \rightarrow w_0 S w_1 \dots S w_k \quad (w_i \in \Sigma^*)$$

$$\pi \text{ is } \mathbf{valid} \stackrel{\text{def}}{\iff} L_* \supseteq w_0 L_* w_1 \dots L_* w_k$$

- If  $\pi$  is not valid,  $\pi$  can't be in a correct grammar for  $L_*$ .  
If  $x_1, \dots, x_k \in L_*$ ,  $w_0 x_1 w_1 \dots x_k w_k \notin L_*$ , and  $L_* \subseteq L(G)$ , then

$$\begin{aligned} S &\Rightarrow w_0 S w_1 \dots S w_k \\ &\Rightarrow^* w_0 x_1 w_1 \dots x_k w_k \end{aligned}$$

- If all productions in  $G$  are valid, then  $L(G) \subseteq L_*$ .
- All productions in  $G_*$  are valid.

# Context-Free Grammars

abbreviates three productions:  
 $S \rightarrow aD_1bS, S \rightarrow aA, S \rightarrow bU$

$$\begin{aligned} S &\rightarrow aD_1bS \mid aA \mid bU \\ D_1 &\rightarrow \epsilon \mid aD_1bD_1 \\ A &\rightarrow \epsilon \mid aD_1bA \mid aA \\ U &\rightarrow \epsilon \mid Ua \mid Ub \end{aligned}$$

$$\begin{aligned} \overline{D}_1 = \{ x \in \{a, b\}^* \mid \\ |x|_a \neq |x|_b \vee \exists uv(uv = x \wedge |u|_a < |u|_b) \} \end{aligned}$$

$$L_G(A) = \{ x \in \Sigma^* \mid A \Rightarrow_G^* x \}$$

$(L_G(S), L_G(D_1), L_G(A), L_G(U))$  is the **least fixed point** of the operator  $\Phi_G: (\mathcal{P}(\{a, b\}^*))^4 \rightarrow (\mathcal{P}(\{a, b\}^*))^4$ :

$$\Phi_G(X_S, X_{D_1}, X_A, X_U) =$$

$$(aX_{D_1}bX_S \cup aX_A \cup bX_U, \epsilon \cup aX_{D_1}bX_{D_1}, \epsilon \cup aX_{D_1}bX_A \cup aX_A, \epsilon \cup X_Ua \cup X_Ub)$$

# Pre-fixed Points of Context-Free Grammars

$$G = (N, \Sigma, P, S)$$

$\Phi_G$ : associated operator

- $(X_B)_{B \in N}$  is a **pre-fixed point** of  $\Phi_G \stackrel{\text{def}}{\iff}$   
 $\Phi_G((X_B)_{B \in N}) \subseteq (X_B)_{B \in N}$   
componentwise inclusion
- $(L_G(B))_{B \in N}$  is the least pre-fixed point of  $\Phi_G$ .
- $(X_B)_{B \in N}$  is a pre-fixed point of  $\Phi_G$  if and only if for every production  $A \rightarrow w_0 B_1 w_1 \dots B_k w_k$  in  $P$ ,

$$X_A \supseteq w_0 X_{B_1} w_1 \dots X_{B_k} w_k.$$

# Fixed Points

$S \rightarrow aD_1bS \mid aA \mid bU$
$D_1 \rightarrow \varepsilon \mid aD_1bD_1$
$A \rightarrow \varepsilon \mid aD_1bA \mid aA$
$U \rightarrow \varepsilon \mid Ua \mid Ub$

$$\begin{aligned}X_S &= aX_{D_1}bX_S \cup aX_A \cup bX_U \\X_{D_1} &= \varepsilon \cup aX_{D_1}bX_{D_1} \\X_A &= \varepsilon \cup aX_{D_1}bX_A \cup aX_A \\X_U &= \varepsilon \cup X_Ua \cup X_Ub\end{aligned}$$

# Pre-fixed Points

$$\begin{array}{l} S \rightarrow aD_1bS \mid aA \mid bU \\ D_1 \rightarrow \varepsilon \mid aD_1bD_1 \\ A \rightarrow \varepsilon \mid aD_1bA \mid aA \\ U \rightarrow \varepsilon \mid Ua \mid Ub \end{array}$$

$$\begin{aligned} X_S &\supseteq aX_{D_1}bX_S \cup aX_A \cup bX_U \\ X_{D_1} &\supseteq \varepsilon \cup aX_{D_1}bX_{D_1} \\ X_A &\supseteq \varepsilon \cup aX_{D_1}bX_A \cup aX_A \\ X_U &\supseteq \varepsilon \cup X_Ua \cup X_Ub \end{aligned}$$

# Pre-fixed Points

$$\begin{array}{l} S \rightarrow aD_1bS \mid aA \mid bU \\ D_1 \rightarrow \varepsilon \mid aD_1bD_1 \\ A \rightarrow \varepsilon \mid aD_1bA \mid aA \\ U \rightarrow \varepsilon \mid Ua \mid Ub \end{array}$$

$$\begin{array}{l} S \rightarrow aD_1bS \\ S \rightarrow aA \\ S \rightarrow bU \\ D_1 \rightarrow \varepsilon \\ D_1 \rightarrow aD_1bD_1 \\ A \rightarrow \varepsilon \\ A \rightarrow aD_1bA \\ A \rightarrow aA \\ U \rightarrow \varepsilon \\ U \rightarrow Ua \\ U \rightarrow Ub \end{array}$$

$$\begin{array}{l} X_S \supseteq aX_{D_1}bX_S \\ X_S \supseteq aX_A \\ X_S \supseteq bX_U \\ X_{D_1} \supseteq \varepsilon \\ X_{D_1} \supseteq aX_{D_1}bX_{D_1} \\ X_A \supseteq \varepsilon \\ X_A \supseteq aX_{D_1}bX_A \\ X_A \supseteq aX_A \\ X_U \supseteq \varepsilon \\ X_U \supseteq X_Ua \\ X_U \supseteq X_Ub \end{array}$$

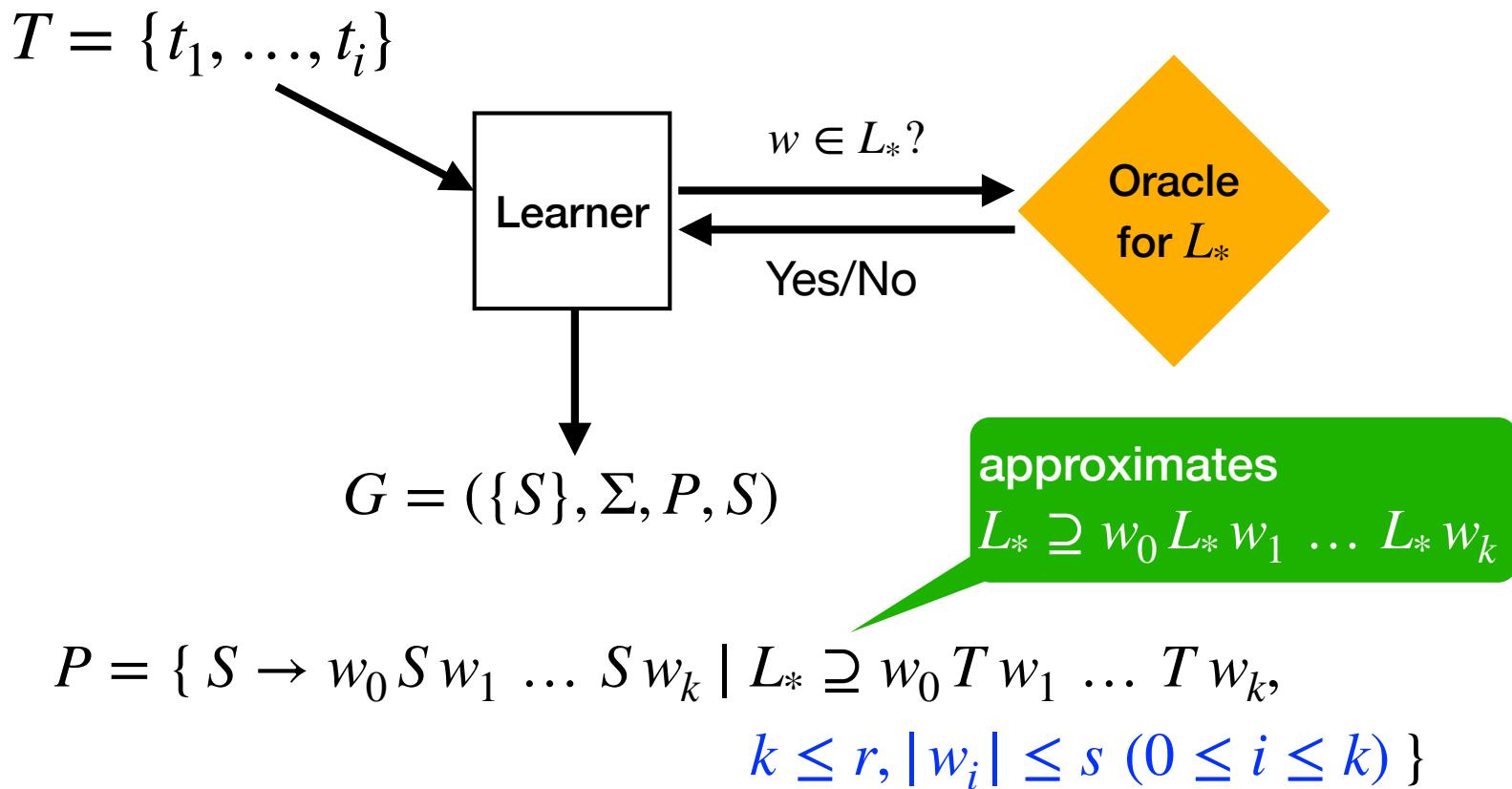
# Simple Case: Grammars with Just One Nonterminal

$$\pi: \quad S \rightarrow w_0 S w_1 \dots S w_k \quad (w_i \in \Sigma^*)$$

$$\pi \text{ is } \mathbf{valid} \stackrel{\text{def}}{\iff} L_* \supseteq w_0 L_* w_1 \dots L_* w_k$$

- If  $\pi$  is not valid,  $\pi$  can't be in a correct grammar for  $L_*$ .
- If all productions in  $G$  are valid, then  $L(G) \subseteq L_*$ .  
 $\therefore L_*$  is a pre-fixed point of  $G$ .
- All productions in  $G_*$  are valid.  
 $\therefore L_* = L(G_*)$  is the least pre-fixed point of  $G_*$ .

# Inference of Context-Free Grammars with Just One Nonterminal



- Need a constant bound  $r$  on  $k$ , since deciding  $L_* \supseteq w_0 T w_1 \dots T w_k$  requires  $|T|^k$  queries to the oracle for  $L_*$ .
- Placing a constant bound  $s$  on  $|w_i|$  makes the set of possible productions finite.

# Learning Context-Free Grammars (with More Than One Nonterminal)

$$\pi: \quad A \rightarrow w_0 B_1 w_1 \dots B_k w_k$$

$$\pi \text{ is valid} \stackrel{\text{def}}{\iff} \llbracket A \rrbracket^{L_*} \supseteq w_0 \llbracket B_1 \rrbracket^{L_*} w_1 \dots \llbracket B_k \rrbracket^{L_*} w_k$$

- A hypothesized nonterminal  $B$  “denotes” a set  $\llbracket B \rrbracket^{L_*}$  relative to the target language  $L_*$ , independently of the rest of the hypothesized grammar.
- In particular, it is not necessarily the case that  $L_G(B) = \llbracket B \rrbracket^{L_*}$  (even in the limit).
- Membership in  $\llbracket B \rrbracket^{L_*}$  **reduces in polynomial time** to membership in  $L_*$ .
- This reduction is uniform across different target languages.

# Simplest Class: Nonterminals Denote Quotients

- The learner uses pairs of strings as nonterminals.

$$\begin{aligned}\llbracket \langle u, v \rangle \rrbracket^{L_*} &= u \setminus L_*/v \\ &= \{ x \in \Sigma^* \mid uxv \in L_* \}\end{aligned}$$

$$\begin{aligned}P = \{ A \rightarrow w_0 B_1 w_1 \dots B_k w_k \mid &\text{ approximates } \llbracket A \rrbracket^{L_*} \supseteq w_0 \llbracket B_1 \rrbracket^{L_*} w_1 \dots \llbracket B_k \rrbracket^{L_*} w_k \\ &\llbracket A \rrbracket^{L_*} \supseteq w_0 (\text{Sub}(T) \cap \llbracket B_1 \rrbracket^{L_*}) w_1 \dots (\text{Sub}(T) \cap \llbracket B_k \rrbracket^{L_*}) w_k, \\ &k \leq r, |w_i| \leq s \ (0 \leq i \leq k) \} &\text{Sub}(T) = \{ x \mid uxv \in T \}\end{aligned}$$

$$S = \langle \varepsilon, \varepsilon \rangle$$

- $L_*$  has infinitely many quotients unless it is regular, so the learner must stop creating new nonterminals.

# Algorithm 1

$T_0 := \emptyset; E_0 := \emptyset; J_0 := \emptyset; G_0 := (\{\langle\langle \varepsilon, \varepsilon \rangle\rangle\}, \Sigma, \emptyset, \langle\langle \varepsilon, \varepsilon \rangle\rangle);$

**for**  $i = 1, 2, 3, \dots$  **do**

$T_i := T_{i-1} \cup \{t_i\};$

**if**  $T_i \subseteq L(G_{i-1})$  **then**

$| J_i := J_{i-1};$

        expand  $J_i$  only when  $T_i \not\subseteq L(G_{i-1})$

**else**

$\text{Con}(T) = \{ (u, v) \mid uwv \in T \}$

$| J_i := \text{Con}(T_i);$

        lexicographic product of  $\prec$   
        with itself

$E = \Sigma^{\leq s} (\text{Sub}(T_i) \Sigma^{\leq s})^{\leq r}$

**end**

$N_i := \{ \langle\langle u, v \rangle\rangle \mid (u, v) \in J_i \wedge$

$\forall (u', v') \in J_i ((u', v') \prec_2 (u, v) \rightarrow E \cap (u' \setminus L_* / v') \neq E \cap (u \setminus L_* / v)) \};$

$P_i := \{ A \rightarrow w_0 B_1 w_1 \dots B_k w_k \mid A, B_1, \dots, B_k \in N_i,$

$\llbracket A \rrbracket^{L_*} \supseteq w_0 (\text{Sub}(T_i) \cap \llbracket B_1 \rrbracket^{L_*}) w_1 \dots (\text{Sub}(T_i) \cap \llbracket B_k \rrbracket^{L_*}) w_k,$

$k \leq r, |w_j| \leq s \ (1 \leq j \leq k) \}$

    output  $G_i := (N_i, \Sigma, P_i, \langle\langle \varepsilon, \varepsilon \rangle\rangle);$

**end**

# CFGs with the Quotient Property

- $G = (N, \Sigma, P, S)$  has the **quotient property**  
 $\overset{\text{def}}{\iff} G$  has a pre-fixed point  $(X_B)_{B \in N}$  with  $X_S = L(G)$  such that  $X_B \in \mathcal{Q}(L(G))$  for all  $B \in N$ .

$$\mathcal{Q}(L) = \{ u \backslash L / v \mid u, v \in \Sigma^* \}$$

## Theorem.

- Algorithm 1 successfully learns  $L_*$  if and only if  $L_*$  has a grammar with the quotient property.
- If Algorithm 1 converges to  $G$ , then  $G$  has the quotient property.

# Examples of CFGs with the Quotient Property

- CFGs with just one nonterminal.
- Right-linear grammars corresponding to minimal DFAs of regular languages.
- $L = \{ a^n b^n \mid n \geq 0 \} \{ a^n b^n \mid n \geq 0 \}$

$$S \rightarrow AA$$

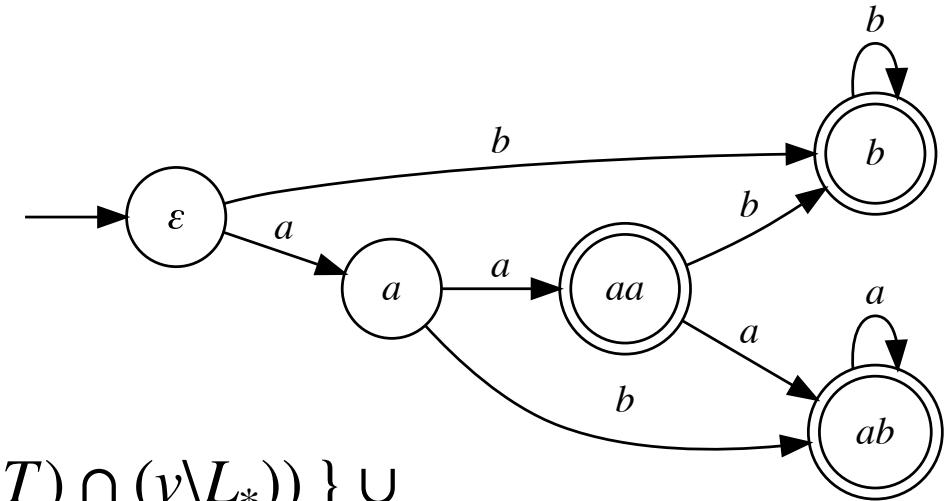
$$A \rightarrow \epsilon \mid aAb$$

$$S = \epsilon \backslash L / \epsilon \quad (= L),$$

$$A = a \backslash L / bab$$

# An Analogous Algorithm for Regular Languages

$$L_* = aba^* \cup bb^* \cup aa(a^* \cup b^*)$$

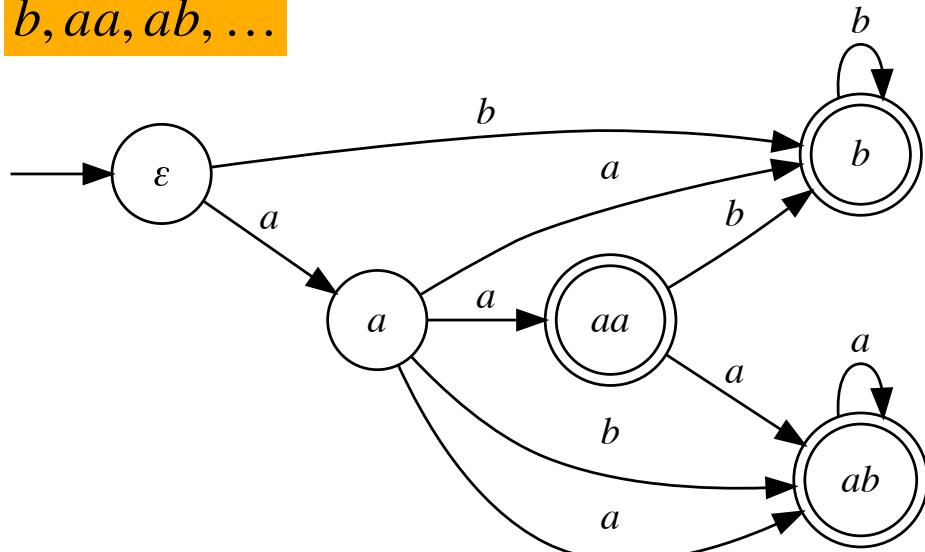


$$P = \{ \langle\langle u \rangle\rangle \rightarrow a \langle\langle v \rangle\rangle \mid u \setminus L_* \supseteq a(\text{Suff}(T) \cap (v \setminus L_*)) \} \cup$$

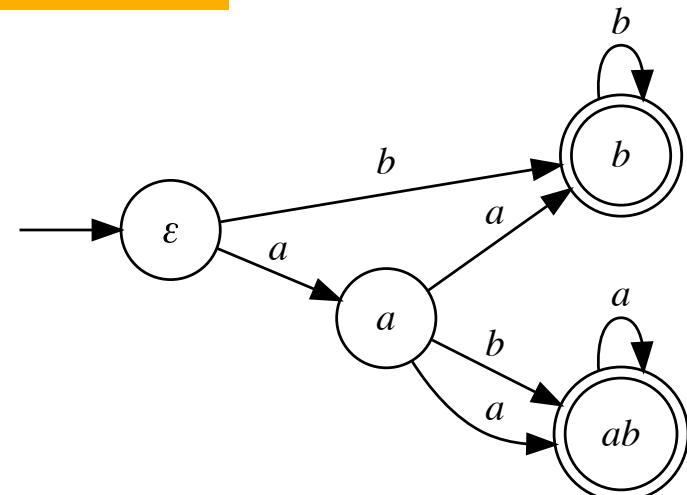
$$\{ \langle\langle u \rangle\rangle \rightarrow \epsilon \mid u \in L_* \}$$

approximates  $u \setminus L_* \supseteq a(v \setminus L_*) \iff ua \setminus L_* \supseteq v \setminus L_*$

$b, aa, ab, \dots$



$b, aba, \dots$



# $\Gamma$ -closure

- $\Gamma$ : finite set of operations on  $\mathcal{P}(\Sigma^*)$  (of variable arity)
- For  $\mathcal{L} \subseteq \mathcal{P}(\Sigma^*)$ ,  
$$\Gamma(\mathcal{L}) = \{ f(L_1, \dots, L_m) \mid f: (\mathcal{P}(\Sigma^*))^m \rightarrow \mathcal{P}(\Sigma^*), f \in \Gamma, L_1, \dots, L_m \in \mathcal{L} \}$$
$$\Gamma^0(\mathcal{L}) = \mathcal{L}$$
$$\Gamma^{n+1}(\mathcal{L}) = \mathcal{L} \cup \Gamma(\Gamma^n(\mathcal{L}))$$

$\Gamma$ -closure of  $\mathcal{Q}(L)$

- Sets in  $\bigcup_{t \geq 0} \Gamma^t(\mathcal{Q}(L))$  can be represented by expressions built from  $\langle\langle u, v \rangle\rangle$  and symbols for operations in  $\Gamma$ .

# Extended Regular Closure

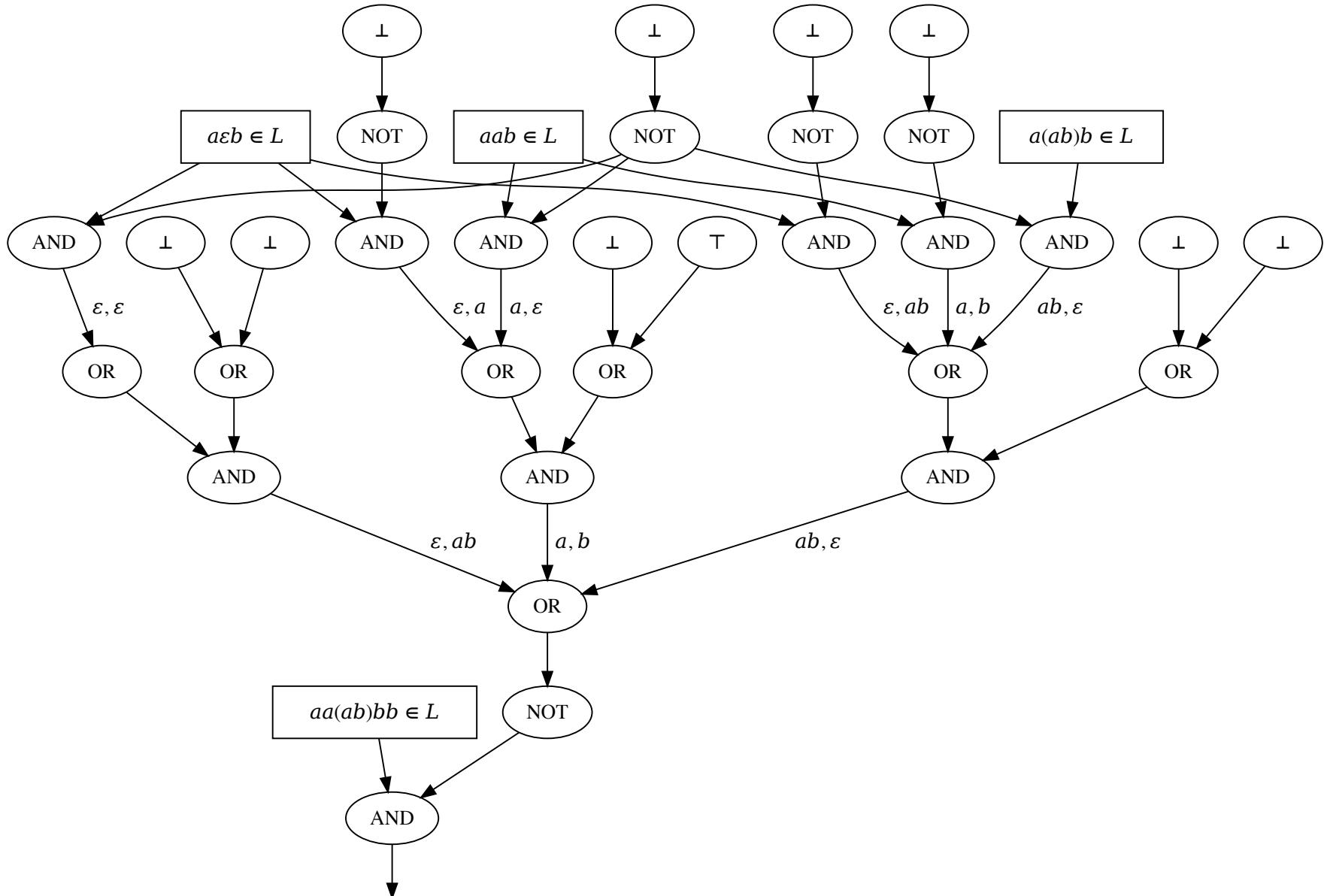
$$\Gamma = \{\cap, \overline{\cdot}, \cup\} \cup \{\emptyset, \varepsilon\} \cup \Sigma \cup \{\text{concatenation}, *\}$$

extended regular expression over query atoms

$$\begin{aligned} & [[\langle\langle aa, bb \rangle\rangle \cap \overline{(\langle\langle a, b \rangle\rangle \overline{\emptyset}) (a \cup b)}]]^L \\ &= (aa \backslash L / bb) \cap (\{a, b\}^* - ((a \backslash L / b)(\{a, b\}^* - \emptyset))(\{a\} \cup \{b\})) \\ &= (aa \backslash L / bb) \cap (\{a, b\}^* - (a \backslash L / b)\{a, b\}^*\{a, b\}) \\ &= \{x \mid x \in aa \backslash L / bb \wedge \text{no proper prefix of } x \text{ is in } a \backslash L / b\} \end{aligned}$$

- If  $e$  is an extended regular expression over query atoms, then  $[[e]]^L$  **reduces in polynomial time** to  $L$ .

$$[[e]]^L \leq_{tt}^P L$$



$$ab \in [[\langle\langle aa, bb \rangle\rangle \cap \overline{(\langle\langle a, b \rangle\rangle \emptyset)}(a \cup b)]]^L$$

# $\Gamma$ -closure Property

- A CFG  $G = (N, \Sigma, P, S)$  has the  **$\Gamma^t$ -property**  $\overset{\text{def}}{\iff} G$  has a pre-fixed point  $(X_B)_{B \in N}$  with  $X_S = L(G)$  such that  $X_B \in \Gamma^t(\mathcal{Q}(L(G)))$  for all  $B \in N$ .
- $G$  has the  **$\Gamma$ -closure property**  $\overset{\text{def}}{\iff} G$  has the  $\Gamma^t$ -property for some  $t$ .

# Learning CFGs with the Extended Regular Closure Property

- The learner uses extended regular expressions over query atoms as nonterminals.

$$P = \{ A \rightarrow w_0 B_1 w_1 \dots B_k w_k \mid \text{approximates } \llbracket A \rrbracket^{L_*} \supseteq w_0 \llbracket B_1 \rrbracket^{L_*} w_1 \dots \llbracket B_k \rrbracket^{L_*} w_k \\ \llbracket A \rrbracket^{L_*} \supseteq w_0 (\text{Sub}(T) \cap \llbracket B_1 \rrbracket^{L_*}) w_1 \dots (\text{Sub}(T) \cap \llbracket B_k \rrbracket^{L_*}) w_k, \\ k \leq r, |w_i| \leq s \ (0 \leq i \leq k) \}$$

$$S = \langle\!\langle \varepsilon, \varepsilon \rangle\!\rangle$$

# Algorithm 2

$T_0 := \emptyset; E_0 := \emptyset; J_0 := \emptyset; G_0 := (\{\langle\!\langle \varepsilon, \varepsilon \rangle\!\rangle\}, \Sigma, \emptyset, \langle\!\langle \varepsilon, \varepsilon \rangle\!\rangle);$

**for**  $i = 1, 2, 3, \dots$  **do**

$T_i := T_{i-1} \cup \{t_i\};$

**if**  $T_i \subseteq L(G_{i-1})$  **then**

|  $J_i := J_{i-1};$

expand  $J_i$  only when  $T_i \not\subseteq L(G_{i-1})$

**else**

|  $J_i := \text{Con}(T_i);$

$\text{Con}(T) = \{ (u, v) \mid uwv \in T \}$

**end**

$Q_i := \{ \langle\!\langle u, v \rangle\!\rangle \mid (u, v) \in J_i \wedge$

lexicographic product of  $\prec$  with itself

$E = \Sigma^{\leq s} (\text{Sub}(T_i) \Sigma^{\leq s})^{\leq r}$

$\forall (u', v') \in J_i ((u', v') \prec_2 (u, v) \rightarrow E \cap (u' \backslash L_* / v') \neq E \cap (u \backslash L_* / v)) \};$

$N_i :=$  the set of  $\Gamma^t$ -expressions over  $Q_i$ ;

$P_i := \{ A \rightarrow w_0 B_1 w_1 \dots B_k w_k \mid A, B_1, \dots, B_k \in N_i,$

$\llbracket A \rrbracket^{L_*} \supseteq w_0 (\text{Sub}(T_i) \cap \llbracket B_1 \rrbracket^{L_*}) w_1 \dots (\text{Sub}(T_i) \cap \llbracket B_k \rrbracket^{L_*}) w_k,$

$k \leq r, |w_j| \leq s \ (1 \leq j \leq k) \}$

output  $G_i := (N_i, \Sigma, P_i, \langle\!\langle \varepsilon, \varepsilon \rangle\!\rangle);$

**end**

$$\{ \cap \} \subseteq \Gamma \subseteq \{\cap, \overline{\cdot}, \cup\} \cup \{\emptyset, \varepsilon\} \cup \Sigma \cup \{\text{concatenation}, *\}$$

## Theorem.

- Algorithm 2 successfully learns  $L_*$  if and only if  $L_*$  has a grammar with the  $\Gamma^t$ -property.
- If Algorithm 2 converges to  $G$ , then  $G$  has the  $\Gamma^t$ -property.

# Extended Regular Closure Property

$$\begin{array}{l} S \rightarrow aD_1bS \mid aA \mid bU \\ D_1 \rightarrow \varepsilon \mid aD_1bD_1 \\ A \rightarrow \varepsilon \mid aD_1bA \mid aA \\ U \rightarrow \varepsilon \mid Ua \mid Ub \end{array}$$

$$\begin{aligned} \overline{D_1} &= \{ x \in \{a,b\}^* \mid \\ &\quad |x|_a \neq |x|_b \vee \exists uv(uv = x \wedge |u|_a < |u|_b) \} \\ &= \{ x \in \{a,b\}^* \mid \\ &\quad |x|_a > |x|_b \vee \exists uv(uv = x \wedge |u|_a < |u|_b) \} \end{aligned}$$

$$\begin{aligned} A &= \{ x \in \{a,b\}^* \mid \forall uv(x = uv \rightarrow |u|_a \geq |u|_b) \} \\ &= \{ x \in \{a,b\}^* \mid \text{no prefix of } x \text{ is in } D_1b \} \\ &= \overline{D_1b\{a,b\}^*} \\ &= \left[ \left[ \overline{\langle \varepsilon, \varepsilon \rangle} b (a \cup b)^* \right] \right]^{\overline{D_1}} \end{aligned}$$

# Boolean Closure Property

$$\begin{array}{l} S \rightarrow D_1 b D_1 \mid D_1 a D_1 \mid D_1 b S \mid S a D_1 \\ D_1 \rightarrow \varepsilon \mid a D_1 b D_1 \end{array}$$

$$\overline{D_1} = \{ x \in \{a, b\}^* \mid \exists mn(m + n > 0 \wedge \text{nf}(x) = b^m a^n) \}$$

normal form of  $x$  under the rewriting  $ab \rightarrow \varepsilon$

$$\begin{aligned} S &\Rightarrow^* D_1 b \dots D_1 b D_1 a D_1 \dots a D_1 \\ &\Rightarrow^* x_1 b \dots x_m b y a z_1 \dots a z_n \end{aligned}$$

$$D_1 = [\overline{\langle \varepsilon, \varepsilon \rangle}]^{\overline{D_1}}$$

# Extended Regular Closure Property

$$S \rightarrow aA \mid aD_1bS \mid bB \mid bD_1^R aS$$

$$A \rightarrow \varepsilon \mid aA \mid aD_1bA$$

$$D_1 \rightarrow \varepsilon \mid aD_1bD_1$$

$$B \rightarrow \varepsilon \mid bB \mid bD_1^R aB$$

$$D_1^R \rightarrow \varepsilon \mid bD_1^R aD_1^R$$

$$\overline{O_1} = \{ x \in \{a,b\}^* \mid |x|_a \neq |x|_b \}$$

$$\begin{aligned} A &= \{ x \in \{a,b\}^* \mid \forall uv(x = uv \rightarrow |u|_a \geq |u|_b) \} \\ &= \{ x \in \{a,b\}^* \mid \neg \exists uv(x = uv \wedge |u|_a + 1 = |u|_b) \} \end{aligned}$$

$$= \overline{O_1} b \{a,b\}^*$$

$$= [\overline{\langle\langle \varepsilon, \varepsilon \rangle\rangle} b (a \cup b)^*]^{O_1}$$

$$D_1 = A \cap O_1$$

$$= [\overline{\langle\langle \varepsilon, \varepsilon \rangle\rangle} b (a \cup b)^* \cap \overline{\langle\langle \varepsilon, \varepsilon \rangle\rangle}]^{O_1}$$

**Theorem.**  $\overline{O_1}$  does not have a grammar with the Boolean closure property.

The pumping lemma applied to a long string  $a^p$  gives

$$\begin{aligned} S &\Rightarrow^* a^{l_1} A a^{r_1} \\ A &\Rightarrow^* a^{l_2} A a^{r_2} \quad (l_2 + r_2 > 0) \\ A &\Rightarrow^* a^m \end{aligned}$$

If  $(X_B)_{B \in N}$  is a pre-fixed point with  $X_S = \overline{O_1}$ , then

$$\{ a^{nl_2+m+nr_2} \mid n \geq 0 \} \subseteq X_A \subseteq \bigcap_{n \geq 0} a^{l_1+nl_2} \setminus \overline{O_1} / a^{nr_2+r_1}$$

It follows that  $\{ |x|_a - |x|_b \mid x \in X_A \}$  is both infinite and co-infinite.

But  $\{ |x|_a - |x|_b \mid x \in u \setminus \overline{O_1} / v \} = \mathbb{Z} - \{ -(|uv|_a - |uv|_b) \}$  is a co-finite set.

## CFLs That Have No Grammar with the Extended Regular Closure Property

$$L = \{ a^l b^m a^n b^q \mid l, m, n, q > 0 \wedge (l = n \vee m > q) \}$$

- $L$  is **inherently ambiguous**.
- $L$  does not have a grammar with the extended regular closure property.

**Question.** Are there any CFLs that are not inherently ambiguous that have no grammar with the extended regular closure property?

# Star-Free Closure Property

$$S \rightarrow aA \mid aD_1bS \mid bB \mid bD_1^R aS$$

$$A \rightarrow \varepsilon \mid aA \mid aD_1bA$$

$$D_1 \rightarrow \varepsilon \mid aD_1bD_1$$

$$B \rightarrow \varepsilon \mid bB \mid bD_1^R aB$$

$$D_1^R \rightarrow \varepsilon \mid bD_1^R aD_1^R$$

$$\overline{O_1} = \{ x \in \{a,b\}^* \mid |x|_a \neq |x|_b \}$$

$$\begin{aligned} A &= \{ x \in \{a,b\}^* \mid \forall uv(x = uv \rightarrow |u|_a \geq |u|_b) \} \\ &= \{ x \in \{a,b\}^* \mid \neg \exists uv(x = uv \wedge |u|_a + 1 = |u|_b) \} \\ &= \overline{O_1} b \{a,b\}^* \\ &= \left[ \overline{\langle\langle \varepsilon, \varepsilon \rangle\rangle b (a \cup b)^*} \right]^{\overline{O_1}} \\ &= \left[ \overline{\langle\langle \varepsilon, \varepsilon \rangle\rangle b \overline{\emptyset}} \right]^{\overline{O_1}} \end{aligned}$$

star-free expression over query atoms

# Extended Regular Closure vs. Star-Free Closure

**Question.** Are there any CFLs that have a grammar with the extended regular closure property but have no grammar with the star-free closure property?

$$L = \{ a^n b^m c^l \mid (n \text{ is odd} \wedge n > m) \vee (n \text{ is even} \wedge n > l) \}$$