

# Counting occurrences for a finite set of words: an inclusion-exclusion approach

*Pierre Nicodème*

CNRS - LIX, École polytechnique

joint work with *Frédérique Bassino*, *Julien Clément* and  
*Julien Fayolle*

# Problem setting

Compute **separately** the number of occurrences of a **non-reduced** set of words  $\mathcal{U}$  in a random text under Bernoulli (non-uniform) model

**Reduced set:** no word is factor of another word

Reduced	Non-Reduced
$\mathcal{U} = \{aab, ba, bb\}$	$\mathcal{U} = \{aa, aab, bbaabb\}$

## Methods

- Formal languages manipulations (Régnier-Szpankowski) (**it fails in the non-reduced case**)
- Aho-Corasick (automaton) + Chomsky-Schützenberger
- Inclusion-Exclusion (Goulden-Jackson, Noonan-Zeilberger)

# (Auto)-Correlation Set

## auto-correlation

$$h = ababa \rightsquigarrow \begin{array}{c} ababa \\ ababa| \\ ababa \\ ababa \end{array} \rightsquigarrow \mathcal{C}_{ababa,ababa} = \{\epsilon, ba, baba\}$$

$$\mathcal{C}_{h,h} = \{ w, \quad h.w = r.h \quad \text{and} \quad |w| < |h| \}$$

## correlation

$$\mathcal{C}_{h_1,h_2} = \{ w, \quad h_1.w = r.h_2 \quad \text{and} \quad |w| < |h_2| \}$$

$$h_1 = baba, \quad h_2 = abaaba \longrightarrow \mathcal{C}_{baba,abaaba} = \{aba, baaba\}$$

# Formal Languages Analysis

## (Régnier-Szpankowski)

$$\text{Right } \mathcal{R} = \{ t = u.h \text{ and } \nexists r, s \neq \epsilon, t = r.h.s \}$$

$$\text{Minimal } \mathcal{M} = \{ t \neq \epsilon, h.t = u.h \text{ and } \nexists r, s, h.t = r.h.s \}$$

$$\text{Ultimate } \mathcal{U} = \{ t, \nexists r, s, h.t = r.h.s \}$$

$$\text{Not } \mathcal{N} = \overline{\mathcal{A}^*.h.\mathcal{A}^*} = \{ t, \nexists r, s, t = r.h.s \}$$

$$\mathcal{A}^* = \mathcal{N} + \mathcal{R}.(\mathcal{M})^*.\mathcal{U} \quad \Rightarrow \quad \mathcal{L}_x = \mathcal{N} + \mathcal{R}x.(\mathcal{M}x)^*.\mathcal{U}$$

# Equations over the languages

$$\mathcal{C} = \mathcal{C}_{h,h} \quad \pi_h = \Pr(h) \text{ (Bernoulli model)}$$

$$(I) \mathcal{A}^* = \mathcal{U} + \mathcal{M}\mathcal{A}^* \quad (II) \mathcal{A}^*h = \mathcal{R}\mathcal{C} + \mathcal{R}\mathcal{A}^*.h$$

$$(III) \mathcal{M}^+ = \mathcal{A}^*.h + \mathcal{C} - \epsilon \quad (IV) \mathcal{N}\mathcal{A} = \mathcal{R} + \mathcal{N} - \epsilon$$

solving

$$R(z) = \frac{\pi_h z^{|h|}}{\pi_h z^{|h|} + (1-z)C(z)}$$

$$U(z) = \frac{1}{\pi_h z^{|h|} + (1-z)C(z)}$$

$$N(z) = \frac{C(z)}{\pi_h z^{|h|} + (1-z)C(z)}$$

$$M(z) = 1 + \frac{z-1}{\pi_h z^{|h|} + (1-z)C(z)}$$

$$L(z, x) = \frac{1}{1-z + \pi_h z^{|h|} \frac{1-x}{x + (1-x)C(z)}}$$

## Reduced sets (Régnier)

$$\mathcal{R}_i, \mathcal{M}_{i,j}, \mathcal{U}_i \rightsquigarrow R_i(z), M_{i,j}(z), U_i(z)$$

functions of  $C_{h_1,h_1}(z), C_{h_2,h_2}(z), C_{h_1,h_2}(z), C_{h_2,h_1}(z)$

$$F(z, \mathbf{x}_1, \mathbf{x}_2) = N(z) + (\mathbf{x}_1 R_1(z), \mathbf{x}_2 R_2(z)) \begin{pmatrix} \mathbf{x}_1 M_{1,1}(z) & \mathbf{x}_2 M_{1,2}(z) \\ \mathbf{x}_1 M_{2,1}(z) & \mathbf{x}_2 M_{2,2}(z) \end{pmatrix}^* \begin{pmatrix} U_1(z) \\ U_2(z) \end{pmatrix}$$

**This collapses in case of non-reduced sets**

# Aho-Corasick

- **Input:** non-reduced set of words  $\mathcal{U}$ .
- **Output:** automaton  $\mathcal{A}_{\mathcal{U}}$  recognizing  $\mathcal{A}^*\mathcal{U}$ .

## Algorithm:

1. build  $\mathcal{T}_{\mathcal{U}}$ , the ordinary **trie** representing the set  $\mathcal{U}$

2. build  $\mathcal{A}_{\mathcal{U}} = (\mathcal{A}, Q, \delta, \epsilon, T)$ :

–  $Q = \text{Pref}(\mathcal{U})$

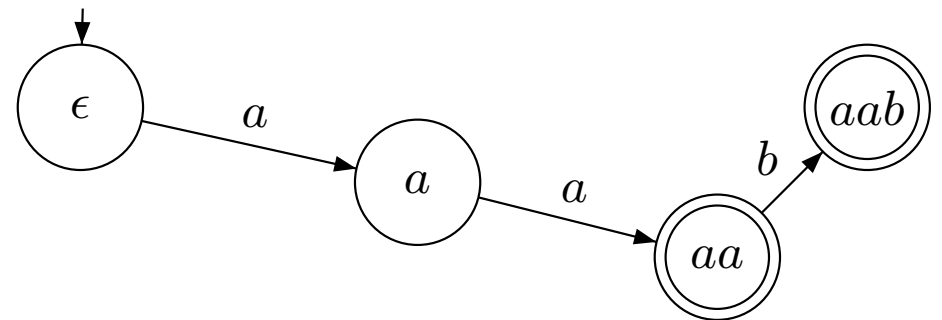
–  $T = \mathcal{A}^*\mathcal{U} \cap \text{Pref}(\mathcal{U})$

–  $\delta(q, x) = \begin{cases} qx & \text{if } qx \in \text{Pref}(\mathcal{U}), \\ \text{Border}(qx) & \text{otherwise,} \end{cases}$

**Border**( $v$ ) = the longest proper suffix of  $v$  which belongs to  $\text{Pref}(\mathcal{U})$  if defined, or  $\epsilon$  otherwise.

# Example

$$\mathcal{U} = \{aab, aa\}$$

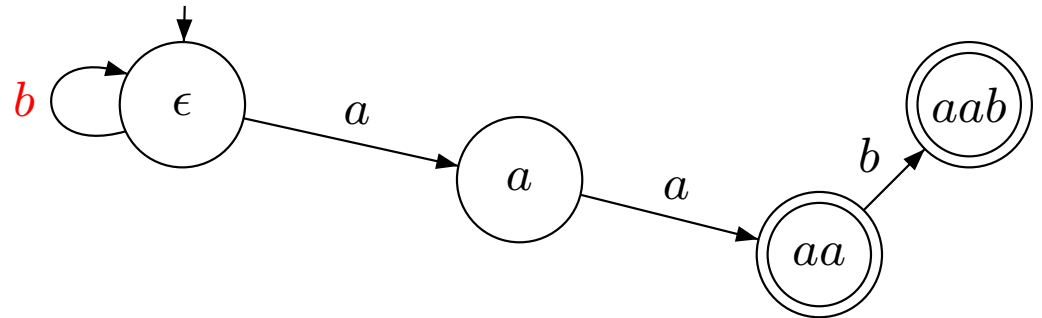


Trie  $\mathcal{T}_{\mathcal{U}}$  of  $\mathcal{U}$



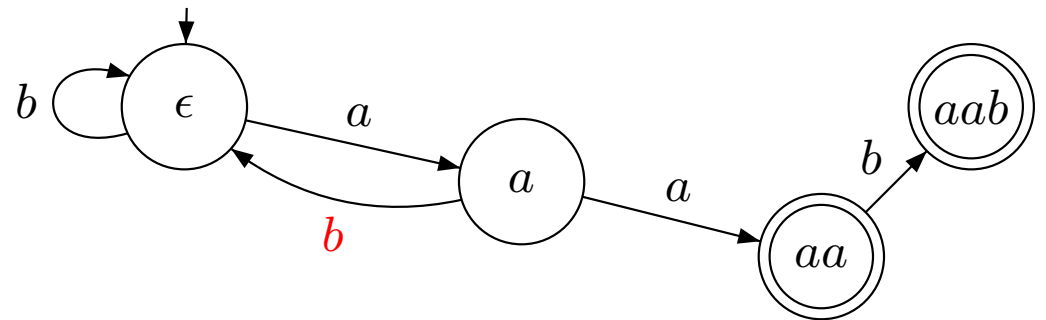
# Example

$$\mathcal{U} = \{aab, aa\} \quad \delta(\epsilon, b) = \text{Border}(b) = \epsilon$$



# Example

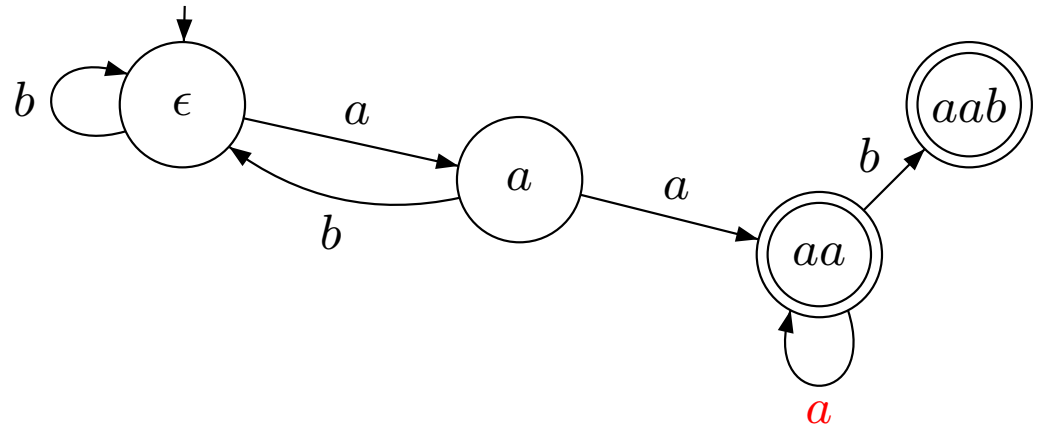
$$\mathcal{U} = \{aab, aa\} \quad \delta(a, b) = \text{Border}(a.b) = \epsilon$$



# Example

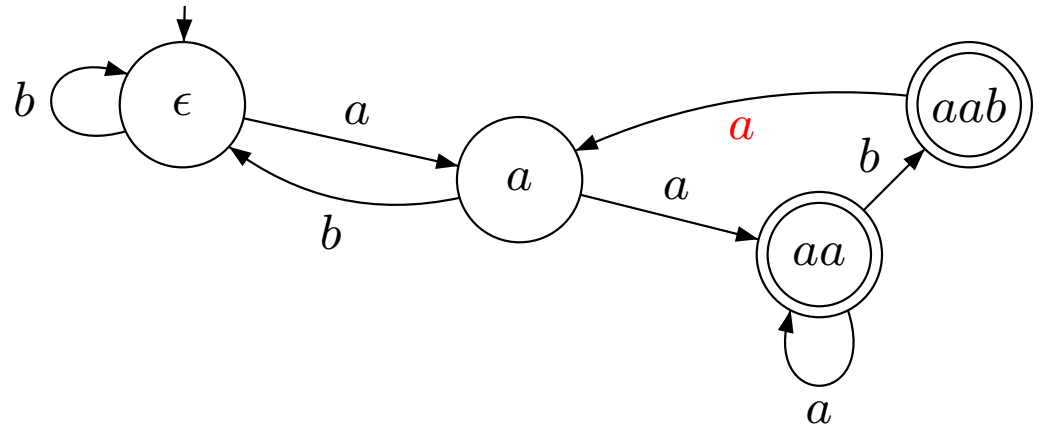
$$\mathcal{U} = \{aab, aa\}$$

$$\delta(aa, a) = \text{Border}(aa.a) = aa$$



# Example

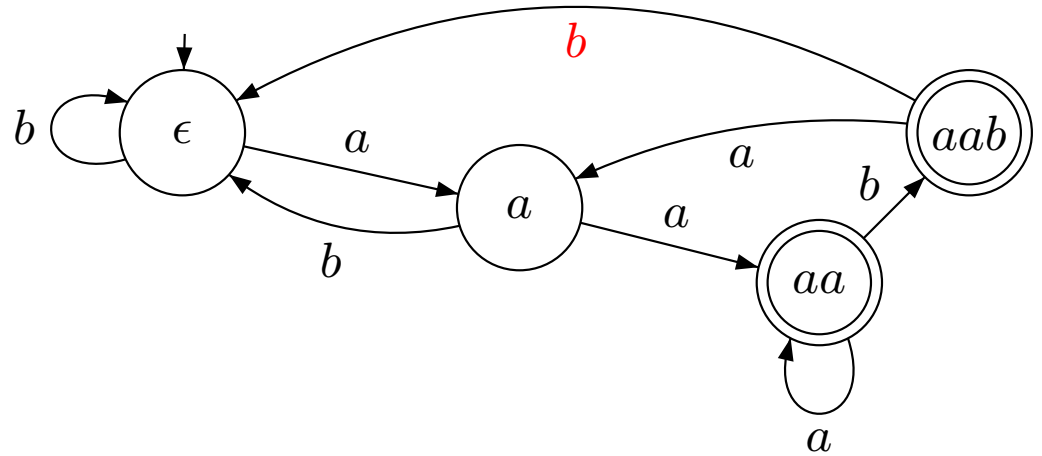
$$\mathcal{U} = \{aab, aa\} \quad \delta(aab, a) = \text{Border}(aab.a) = a$$



# Example

$$\mathcal{U} = \{aab, aa\}$$

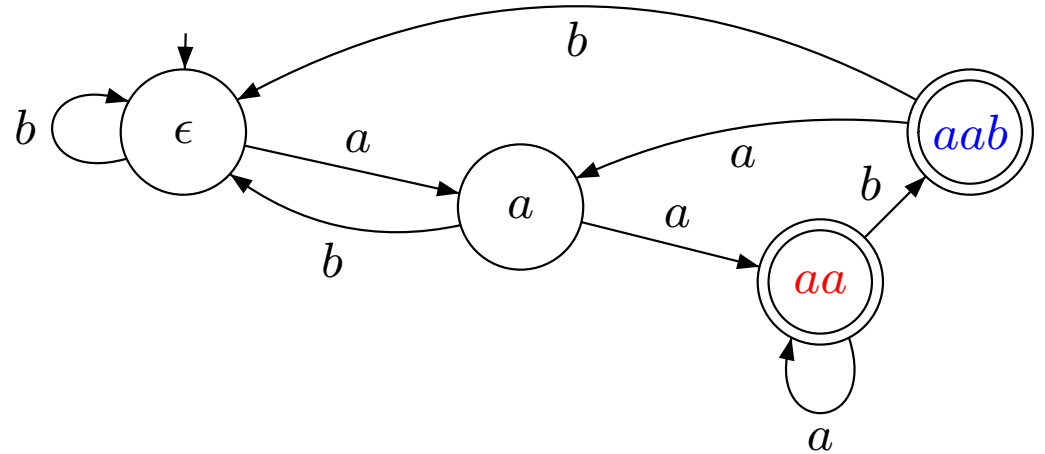
$$\delta(aab, b) = \text{Border}(aab.b) = \epsilon$$



# Example

$$\mathcal{U} = \{aab, aa\}$$

$$\mathbb{T}(x_1, x_2) = \begin{pmatrix} b & a & 0 & 0 \\ b & 0 & ax_2 & 0 \\ 0 & 0 & ax_2 & bx_1 \\ b & a & 0 & 0 \end{pmatrix},$$

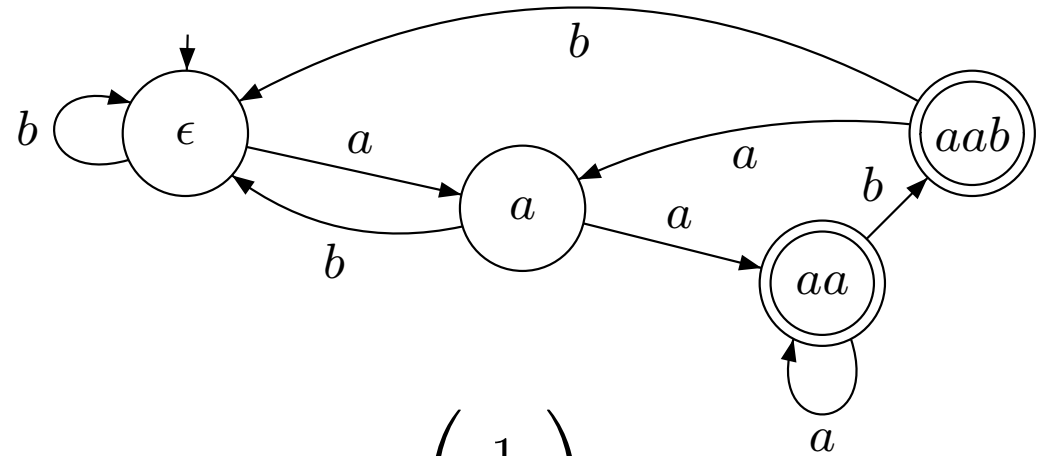


$x_1, x_2$  marks for  $aab, aa$

# Example

$$\mathcal{U} = \{aab, aa\}$$

$$\mathbb{T}(x_1, x_2) = \begin{pmatrix} b & a & 0 & 0 \\ b & 0 & ax_2 & 0 \\ 0 & 0 & ax_2 & bx_1 \\ b & a & 0 & 0 \end{pmatrix},$$



$$F(a, b, x_1, x_2) = (1, 0, 0, 0)(\mathbb{I} - \mathbb{T}(a, b, x_1, x_2))^{-1} \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}$$

$$= \frac{1 - a(x_2 - 1)}{1 - ax_2 - b + ab(x_2 - 1) - a^2bx_2(x_1 - 1)^2}.$$

# Inclusion-Exclusion: one word

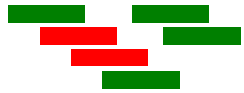
word *aaa*  $p(x)$ : unknown p.g.f of counts of *aaa*

*bbbbbaaaaaaaaaabbbbb*

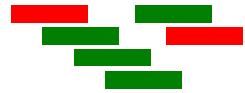


each occurrence is marked or not (flip-flop)

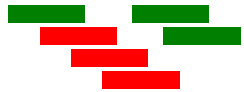
*bbbbbaaaaaaaaaabbbbb*



*bbbbbaaaaaaaaaabbbbb*



*bbbbbaaaaaaaaaabbbbb*



*bbbbbaaaaaaaaaabbbbb*



$$\begin{array}{c} \text{---} \\ \text{---} \end{array} \rightsquigarrow \begin{cases} \text{---} \\ \text{---} \end{cases} \quad x \rightsquigarrow \begin{cases} 1 \\ +x \end{cases} \quad \begin{array}{l} p(x) \rightsquigarrow p(1+x) = \phi(x) \\ \rightsquigarrow p(x) = \phi(x-1) \end{array}$$

computing **easier**  $\phi(t)$  and substituting  $t \rightsquigarrow x-1$

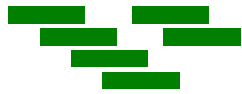
give **harder**  $p(x)$



# One word - Clusters

word *aaa*

*bbbbbaaaaaaaaaabbbb*



*bbbbbaaaaaaaaaabbbb*



*bbbbbaaaaaaaaaabbbb*



*bbbbbaaaaaaaaaabbbb*



*bbbbbaaaaaaaaaabbbb*



clusters  $\mathfrak{C}$

$$\mathfrak{C} = w + \mathfrak{C}.(\mathcal{C}_{w,w} - \epsilon) \implies \mathfrak{C}(z, x) = \frac{x\pi_w z^{|w|}}{1 - x(C(z) - 1)}$$

$$\mathcal{T} = \text{Seq}(\mathcal{A} + \mathfrak{C}) \implies \Phi(z, x) = \frac{1}{1 - z - \mathfrak{C}(z, x)}$$

$$F(z, x) = \Phi(z, x - 1)$$

# Three words - Clusters (Goulden-Jackson)

$$\mathcal{U} = \{aba, bab, aa\}$$

bbbbbabababaabbbb



bbbbbabababaabbbb



bbbbbabababaabbbb



clusters  $\mathfrak{C}_{i,j}$  begin with  $w_i$  and finish with  $w_j$

$$\mathfrak{C}_{i,j} = w_i \mathcal{C}_{w_i, w_j} + \sum_{1 \leq k \leq 3} \mathfrak{C}_{i,k} \cdot (\mathcal{C}_{w_k, w_j} - \delta_{kj} \epsilon)$$

$$\mathfrak{C} = (w_1 \bullet, w_2 \bullet, w_3 \bullet) \left( \mathbf{I} - \begin{pmatrix} \mathcal{C}_{w_1, w_1} \bullet - \epsilon & \mathcal{C}_{w_1, w_2} \bullet & \mathcal{C}_{w_1, w_3} \bullet \\ \mathcal{C}_{w_2, w_1} \bullet & \mathcal{C}_{w_2, w_2} \bullet - \epsilon & \mathcal{C}_{w_2, w_3} \bullet \\ \mathcal{C}_{w_3, w_1} \bullet & \mathcal{C}_{w_3, w_2} \bullet & \mathcal{C}_{w_3, w_3} \bullet - \epsilon \end{pmatrix} \right)^{-1} \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$$

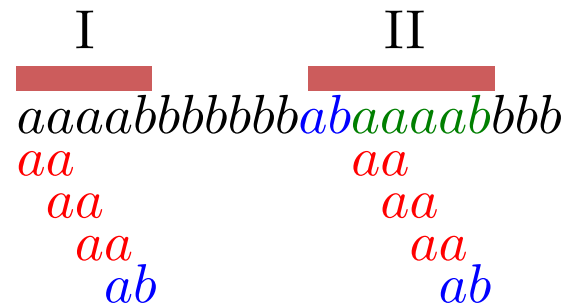
$$\mathcal{T} = \text{Seq}(\mathcal{A} + \mathfrak{C}) \implies \Phi(z, \mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3) = \frac{1}{1 - z - \mathfrak{C}(z, \mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3)}$$

$$F(z, \mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3) = \Phi(z, \mathbf{x}_1 - 1, \mathbf{x}_2 - 1, \mathbf{x}_3 - 1) = \frac{1}{1 - z - \mathfrak{C}(z, \mathbf{x}_1 - 1, \mathbf{x}_2 - 1, \mathbf{x}_3 - 1)}$$

# General Case: Non Reduced Set of Words

$$\mathcal{U} = \{aa, ab, baaaab\}$$

Consider clusters I and II

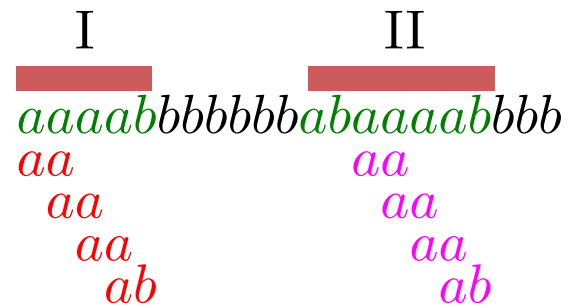


- **Cluster I**: as in the one word or reduced case: “flip-flop” counting
- **Cluster II**: there are supplementary matches: “sticky” counting, attached to a right extension of one word to another word (here  $ab$  to  $baaaab$ ).

Main idea: consider the “reduced” backbone by counting flip-flop occurrences, add the sticky occurrences during the right extensions from one word to another (generalization of correlation)

# “Flip-Flop” versus “Sticky”

$$\mathcal{U} = \{u_1 = aa, u_2 = ab, u_3 = baaaaab\}$$



- Cluster I “flip-flop”
- Cluster II contains “sticky” words

# Inclusion-Exclusion: Non-Reduced Case

$$\mathcal{U} = \{u_1 = aa, u_2 = ab, u_3 = baaaab\}$$

aaaabbbbbbabaaaabbbb  
 aa aa  
 aa ab aa  
 aa aa  
 ab ab  
 baaaab

I II  
  
 aaaabbbbbbabaaaabbbb  
 aa aa  
 aa ab aa  
 aa aa  
 ab ab  
 baaaab

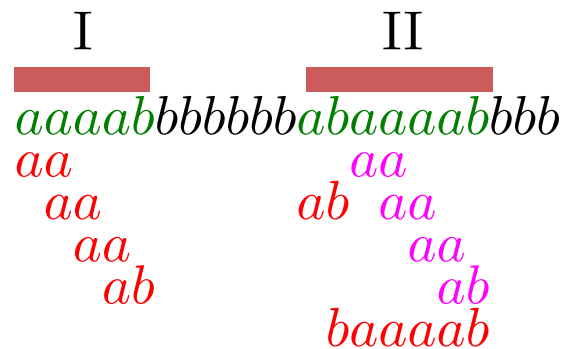
I II  
  
 aaaabbbbbbabaaaabbbb  
 aa aa  
 aa ab aa  
 aa aa  
 ab ab  
 baaaab

I II  
  
 aaaabbbbbbabaaaabbbb  
 aa aa  
 aa ab aa  
 aa aa  
 ab ab  
 baaaab

1. **flip-flop** all occurrences giving **clusters**
2. **forget factor** occurrences
3. add **sticky** occurrences

# Counting “Flip-Flop” versus “Sticky”

$$\mathcal{U} = \{u_1 = aa, u_2 = ab, u_3 = baaaaab\}$$



– Cluster I “flip-flop”:  $t_i \rightsquigarrow x_i - 1$

– Cluster II

1. “flip-flop”:  $t_i \rightsquigarrow x_i - 1$

2. “sticky”:  $v_i \rightsquigarrow x_i$

Remark: counting the sticky occurrences by  $v_i = 1 + t_i$  and doing

$t_i \rightsquigarrow x_i - 1$  is correct

# Right Extension Sets and Matrices

**Right Extension Set** of a pair of words  $(h_1, h_2)$

$$\mathcal{E}_{h_1, h_2} = \{ e \mid \text{there exists } e' \in \mathcal{A}^+ \text{ such that } h_1 e = e' h_2 \text{ with } 0 < |e| < |h_2| \}.$$

if  $h_1 \neq h_2$  have no factor relation,  $\mathcal{E}_{h_1, h_2} = \mathcal{C}_{h_1, h_2}$  but  $\mathcal{E}_{h, h} = \mathcal{C}_h - \epsilon$

**Right Extension Matrix** of a vector of words  $\mathbf{u} = (u_1, \dots, u_r)$

$$\mathcal{E}_{\mathbf{u}} = (\mathcal{E}_{u_i, u_j})_{1 \leq i, j \leq r}.$$

## Examples

$$\mathbf{u}_1 = (aba, ab) \Rightarrow \mathcal{E}_{\mathbf{u}_1} = \begin{pmatrix} ba & b \\ \emptyset & \emptyset \end{pmatrix} \quad \mathcal{E}_{ab, aba} = \emptyset \quad \begin{cases} aba = |aba \\ e' = \epsilon \notin \mathcal{A}^+ \end{cases}$$

$$\mathbf{u}_2 = (aaaa, aaa) \Rightarrow \mathcal{E}_{\mathbf{u}_2} = \begin{pmatrix} a+a^2+a^3 & a+a^2 \\ a^2+a^3 & a+a^2 \end{pmatrix} \quad \begin{cases} a \notin \mathcal{E}_{aaa, aaaa} & aaa.a = |aaaa \\ aa \in \mathcal{E}_{aaa, aaaa} & aaa.aa = a.aaaa \end{cases}$$

# Counting Sticky Words

$$\mathcal{U} = \{u_1 = aa, u_2 = baaaabaaaab\} \quad \mathcal{E}_{u_2, u_2} = \{aaaab, aaaabaaaab\}$$

baaaaabaaaabaaaab  
baaaaabaaaabaaaab

$$N_{2,1}(6) = 9 - 6 = 3$$

baaaaabaaaabaaaab  
baaaaabaaaabaaaab

$$N_{2,1}(11) = 9 - 3 = 6$$

$$N_{i,j}(k) = |u_i|_j - |u_i[1 \dots |u_i| - k]|_j.$$

$$\langle \mathcal{E}_{u_2, u_2} \rangle_2 = \pi_a^4 \pi_b z^5 (t_1 + 1)^3 t_2 + \pi_a^8 \pi_b^2 z^{10} (t_1 + 1)^6 t_2$$



# Formal Setting

$N_{i,j}(k)$  counts the number of occurrences of  $u_j$  in  $u_i$  ending in the last  $k$  positions

$$N_{i,j}(k) = |u_i|_j - |u_i[1 \dots |u_i| - k]|_j.$$

$\langle s \rangle_i$  **formal weight** of a **suffix** of word  $u_i$

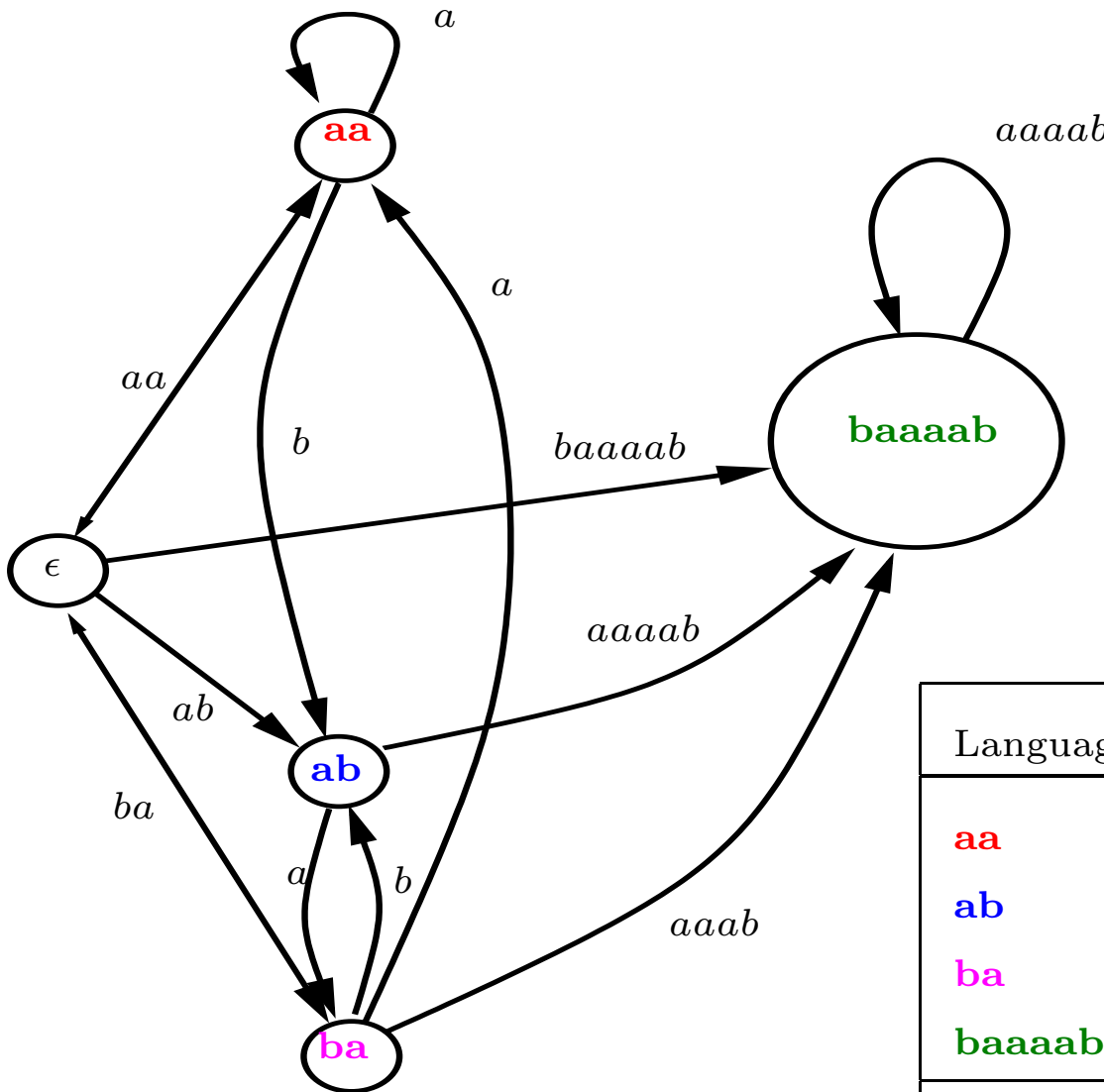
$$\langle s \rangle_i = \pi(s) z^{|s|} t_i \prod_{m \neq i} (t_m + 1)^{N_{i,m}(|s|)}.$$

extension to a set of words  $S$  which are suffixes of  $u_i$

$$\langle S \rangle_i = \sum_{s \in S} \langle s \rangle_i.$$

$$\mathcal{E}_{i,j} \rightsquigarrow \langle \mathcal{E}_{i,j} \rangle_j$$

# Right Extension Graph



$$\mathcal{U} = \{\mathbf{aa}, \mathbf{ab}, \mathbf{ba}, \mathbf{baaaaab}\}$$

Language	G. F.
$\mathbf{aa}$	$t_1 z^2$
$\mathbf{ab}$	$t_2 z^2$
$\mathbf{ba}$	$t_3 z^2$
$\mathbf{baaaaab}$	$t_4 z^6$
$\mathcal{E}_{\mathbf{ab}, \mathbf{ba}} = \{a\}$	$t_3 z$
$\mathcal{E}_{\mathbf{ba}, \mathbf{baaaaab}} = \{aaaab\}$	$(1 + t_1)^2 (1 + t_2) t_4 z^4$
$\mathcal{E}_{\mathbf{baaaaab}, \mathbf{baaaaab}} = \{aaaaab\}$	$(1 + t_1)^3 (1 + t_2) t_4 z^5$

# Putting Things Together

$$\text{Let } \langle \mathbf{u} \rangle = (\langle u_1 \rangle_1, \dots, \langle u_r \rangle_r) \quad \text{and} \quad \langle \mathcal{E}_{\mathbf{u}} \rangle = \begin{pmatrix} \dots & \dots & \dots \\ \dots & \langle \mathcal{E}_{i,j} \rangle_j & \dots \\ \dots & \dots & \dots \end{pmatrix}$$

**Proposition I.** *The generating function  $\mathfrak{C}(z, \mathbf{t})$  of clusters built from the set  $\mathcal{U} = \{u_1, \dots, u_r\}$  is given by*

$$\mathfrak{C}(z, \mathbf{t}) = \langle \mathbf{u} \rangle \cdot \left( \mathbb{I} - \langle \mathcal{E}_{\mathbf{u}} \rangle \right)^{-1} \cdot \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix},$$

where  $\mathbf{u} = (u_1, \dots, u_r)$ ,  $\mathbf{t} = (t_1, \dots, t_r)$

**Proposition II.** *The generating function  $F(z, \mathbf{x})$  counting matches of a non-reduced set of words is*

$$F(z, \mathbf{x}) = \frac{1}{1 - z - \mathfrak{C}(z, \mathbf{x} - \mathbf{1})}$$

# Examples

$$\mathcal{U} = \{u\}$$

$$\mathfrak{G}(z, t) = \frac{t\langle u \rangle}{1 - t\langle \mathcal{E}_u \rangle} = \frac{t\pi(u)z^{|u|}}{1 - t(C(z) - 1)}$$

$$\mathcal{U} = \{u_1, u_2\}$$

$$\mathfrak{G}(z, t_1, t_2)$$

$$= \frac{t_1\langle u_1 \rangle_1 + t_2\langle u_2 \rangle_2 - t_1t_2(\langle u_1 \rangle_1[\langle \mathcal{E}_{2,2} \rangle_2 - \langle \mathcal{E}_{1,2} \rangle_2] + \langle u_2 \rangle_2[\langle \mathcal{E}_{1,1} \rangle_1 - \langle \mathcal{E}_{2,1} \rangle_1])}{1 - t_2\langle \mathcal{E}_{2,2} \rangle_2 - t_1\langle \mathcal{E}_{1,1} \rangle_1 + t_1t_2(\langle \mathcal{E}_{1,1} \rangle_1\langle \mathcal{E}_{2,2} \rangle_2 - \langle \mathcal{E}_{2,1} \rangle_1\langle \mathcal{E}_{1,2} \rangle_2)}$$

# Algorithmic computation

INIT( $\mathcal{A}_U$ )

```
1  for  $i \leftarrow 1$  to  $r$  do
2       $f_i(u_i) \leftarrow 1$ 
3  for  $w \in \text{Pref}(U)$  by a postorder traversal of the tree do
4      for  $i \leftarrow 1$  to  $r$  do
5          for  $\alpha \in \mathcal{A}$  such that  $w \cdot \alpha \in \text{Pref}(u_i)$  do
6               $f_i(w) \leftarrow \pi(\alpha) \cdot f_i(w \cdot \alpha) \prod_{j \neq i} (1 + t_j)^{\llbracket u_j \text{ suffix of } w \cdot \alpha \rrbracket}$ 
7  return  $(f_i)_{1 \leq i \leq r}$ 
```

BUILD-EXTENSION-MATRIX( $\mathcal{A}_U$ )

```
1  ▷ Initialize the matrix  $(\mathcal{E}_{i,j})_{1 \leq i,j \leq r}$ 
2  for  $i \leftarrow 1$  to  $r$  do
3      for  $j \leftarrow 1$  to  $r$  do
4           $\mathcal{E}_{i,j} \leftarrow 0$ 
5  ▷ Compute the maps  $(f_i(w))$  for  $i = 1..r$  and  $w \in \text{Pref}(U)$ 
6   $(f_i)_{1 \leq i \leq r} \leftarrow \text{INIT}(\mathcal{A}_U)$ 
7  ▷ Main loop
8  for  $i \leftarrow 1$  to  $r$  do
9       $v \leftarrow u_i$ 
10     do for  $j \leftarrow 1$  to  $r$  do
11          $\mathcal{E}_{i,j} \leftarrow \mathcal{E}_{i,j} + f_j(v)$ 
12          $v \leftarrow \text{Border}(v)$ 
13     while  $v \neq \epsilon$ 
14  return  $E$ 
```

Time complexity of the main loop  $O(s \times r^2)$ , where  $r$  is the number of words and  $s$  is the length of the longest suffix chain

(sequence  $(u_1 = u, u_2 = \text{Border}(u_1), u_3 = \text{Border}(u_2), \dots, u_s = \text{Border}(u_{s-1}) = \epsilon)$ )