

# DNA evolution, Automata and Clumps

Pierre Nicodème

LIPN Team CALIN, University Paris 13, Villetaneuse

# Biological Motivation

- ▶ **Promoters** are DNA sequences located upstream of the gene they regulate; regulation can be positive for enhancers or negative for repressors.
- ▶ The promoters contain binding sites for regulatory proteins such as **Transcription Factors (TFs)** that are **short stretches of DNA**.

# Biological Motivation

- ▶ **Promoters** are DNA sequences located upstream of the gene they regulate; regulation can be positive for enhancers or negative for repressors.
- ▶ The promoters contain binding sites for regulatory proteins such as **Transcription Factors (TFs)** that are **short stretches of DNA**.
- ▶ **Waiting time: how long** it takes for a **Transcription Factor** to **appear** in a **promoter** under a **probabilistic model of evolution** helps understanding the **overall evolution of promoters** within species and between species?

# From infinitesimal to discrete evolution model

- ▶  $Q(t)dt$  evolution matrix for **infinitesimal time**
- ▶  $P(t)$  evolution matrix **from time**  $x$  **and time**  $x + t$

$$P(t) = e^{Q(t)} \quad (\text{Karlin-Taylor 1975})$$

- ▶  $P(1) = (\pi_{\alpha \rightarrow \beta})$  evolution matrix for **one generation (20 years)**,  $\alpha, \beta \in \{A, C, G, T\}$

# Initial $\nu(\alpha)$ and Substitution Probabilities $\pi_{\alpha \rightarrow \beta}$

$\alpha$	$\nu(\alpha)$
A	0.23889
C	0.26242
G	0.25865
T	0.24004



substitution  
probability  $\pi_{\alpha \rightarrow \beta}$   
for one generation  
(20 years)

A		A	0.9999999763
A		C	$4.54999994943 \times 10^{-9}$
A		G	$1.57499995613 \times 10^{-8}$
A		T	$3.40000001733 \times 10^{-9}$
C		A	$6.14999993408 \times 10^{-9}$
C		C	0.99999996495
C		G	$7.14999984731 \times 10^{-9}$
C		T	$2.17499993935 \times 10^{-8}$
G		A	$2.17499993935 \times 10^{-8}$
G		C	$7.14999984731 \times 10^{-9}$
G		G	0.99999996495
G		T	$6.14999993408 \times 10^{-9}$
T		A	$3.40000001733 \times 10^{-9}$
T		C	$1.57499995613 \times 10^{-8}$
T		G	$4.54999994943 \times 10^{-9}$
T		T	0.9999999763

# Probability of occurrence of a $k$ -mer at time 1

- ▶  $S_n(0)$  random DNA sequence of length  $n$  at time 0
- ▶  $S_n(1)$  sequence obtained from  $S_n(0)$  by evolution at time 1
- ▶  $b$  a  $k$ -mer (word of length  $k$  over  $\mathcal{A} = \{A, C, G, T\}$ )
- ▶  $\mathfrak{P}_n(b)$  probability that  $b$ 
  - ▶ occurs at time 1
  - ▶ while not occurring at time 0

$$\mathfrak{P}_n(b) = \mathbb{P}(b \in S_n(1) \mid b \notin S_n(0))$$

# Probability of occurrence of a $k$ -mer at time 1

- ▶  $S_n(0)$  random DNA sequence of length  $n$  at time 0
- ▶  $S_n(1)$  sequence obtained from  $S_n(0)$  by evolution at time 1
- ▶  $b$  a  $k$ -mer (word of length  $k$  over  $\mathcal{A} = \{A, C, G, T\}$ )
- ▶  $\mathfrak{P}_n(b)$  probability that  $b$ 
  - ▶ occurs at time 1
  - ▶ while not occurring at time 0

$$\mathfrak{P}_n(b) = \mathbb{P}(b \in S_n(1) \mid b \notin S_n(0))$$

## Expectation of the Waiting time $\mathfrak{E}_n(b)$

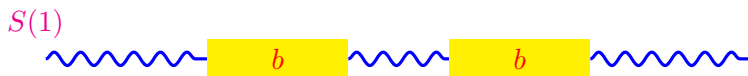
- ▶  $\mathfrak{E}_n(b) \approx \frac{1}{\mathfrak{P}_n(b)}$  (geometric distribution – BehVin2010)

# Plan of the talk

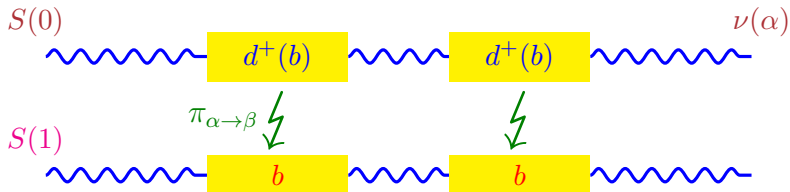
## Different computations of $\mathfrak{P}_n$

1. Behrens-Vingron (2010)
  - ▶ Approach **neglecting words correlation**.
  - ▶ **Efficient computation** of  $\mathfrak{P}_n$  with respect to this assumption.
2. Behrens-Nicaud-P.N. (2012)
  - ▶ **Rigorous and efficient approach by automata**.
  - ▶ Approach **hiding the quasi-linear behaviour** of  $\mathfrak{P}_n$
3. P.N. (NCMA2012)
  - ▶ Approach by **clump analysis**, either by **combinatorics of words** or by **automata**.
  - ▶ **Proof** of the **quasi-linear behaviour** of  $\mathfrak{P}_n$
4. P.N. (ALEA 2013)
  - ▶ **Explicit formula** for  $\mathfrak{P}_n$

# Behrens-Vingron 2010

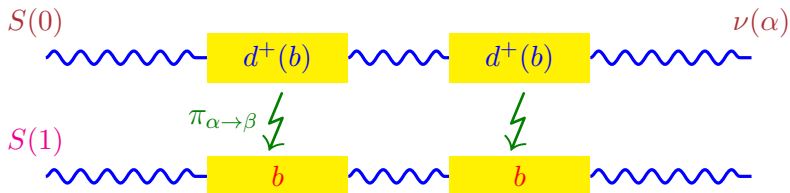


# Behrens-Vingron 2010



- $d^+(b)$  neighbors of  $b$  by substitution

# Behrens-Vingron 2010

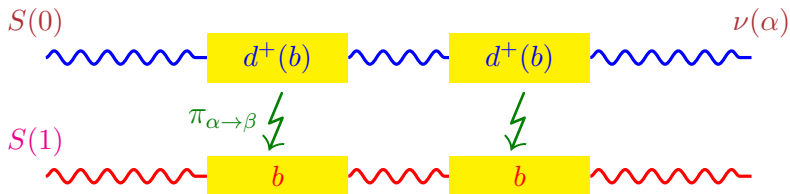


- $d^+(b)$  neighbors of  $b$  by substitution

$$\left\{ \begin{array}{l} \mathfrak{P}_n \approx \sum_{i=1}^{\lfloor n/k \rfloor} (-1)^{i+1} \binom{n-i(k-1)}{i} \Phi^i \\ \Phi = \sum_{(a_1, \dots, a_k) \in \mathcal{A}^k \setminus \{b_1, \dots, b_k\}} \nu(a_1) \times \dots \times \nu(a_k) \cdot \prod_{j=1}^k \pi_{a_j \rightarrow b_j}(1) \end{array} \right.$$

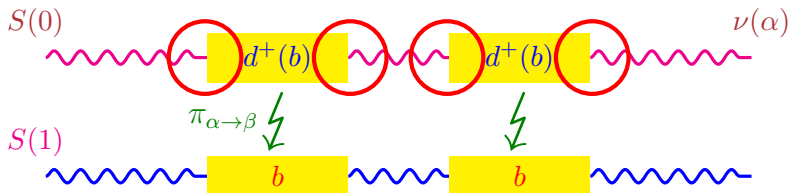
# Approximations of Behrens-Vingron 2010

- occurrences of  $b$  in  $S(1)$  **do not overlap**



# Approximations of Behrens-Vingron 2010

- ▶ occurrences of  $b$  in  $S(1)$  **do not overlap**
- ▶ possible **unwanted occurrences** of  $b$  at **junctions** in  $S(0)$

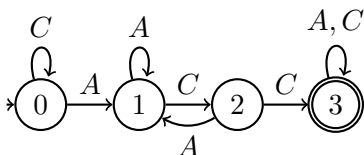


# Behrens-Nicaud-P.N. 2012

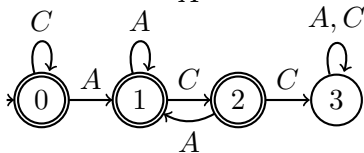
Construct an **automaton**

- ▶ on the **alphabet**  $\Sigma = \mathcal{A} \times \mathcal{A}$  with  $\mathcal{A} = \{A, C, G, T\}$
- ▶ **recognizing sequences**  $S(b) = S(0) \otimes S(1)$
- ▶ **such that**
  1.  $b \notin S(0)$
  2.  $b \in S(1)$

# Using the Knuth-Morris-Pratt automaton



$$\mathcal{M}_{\text{ACC}} = \{Q, \delta, s = 0, \textcolor{red}{F}\}$$

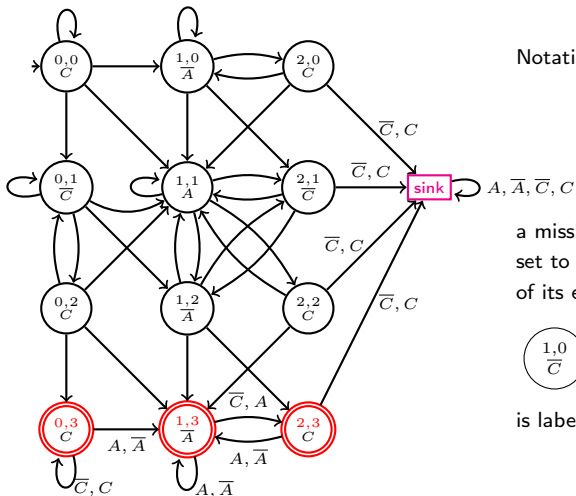


$$\overline{\mathcal{M}}_{\text{ACC}} = \{Q, \delta, s = 0, \textcolor{red}{Q} \setminus \textcolor{red}{F}\}$$

$$\begin{cases} \mathcal{M}_b = (Q = \{0, \dots, k\}, \delta_b, 0, \{\textcolor{red}{k}\}) \\ \overline{\mathcal{M}}_b = (Q = \{0, \dots, k\}, \delta_b, 0, \{\textcolor{red}{0}, \dots, \textcolor{red}{k-1}\}) \\ \mathcal{N}_b = \overline{\mathcal{M}}_b \otimes \mathcal{M}_b = (Q \times Q, \Delta, q'_0 = (0, 0), \textcolor{red}{F}' = \{0, \dots, k-1\} \times \{\textcolor{red}{k}\}) \end{cases}$$

$$\Delta((r, s), (\alpha, \beta)) = (\delta_b(r, \alpha), \delta_b(s, \beta))$$

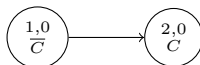
The automaton  $\mathcal{N}_{\text{ACC}} = \overline{\mathcal{M}}_{\text{ACC}} \otimes \mathcal{M}_{\text{ACC}}$  with matrix  $\mathbb{P}$



Notations for the transitions:

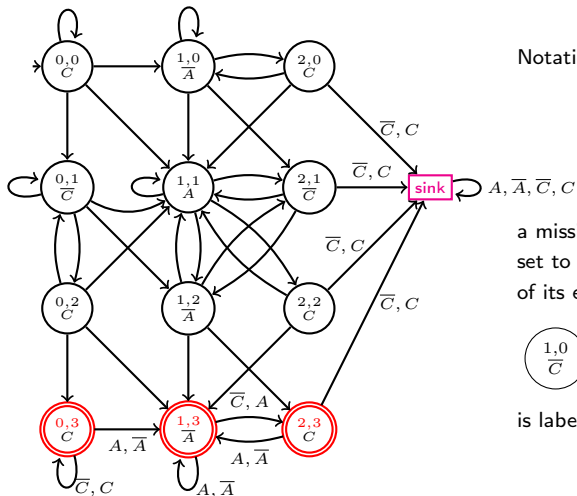
$$\begin{cases} A = \begin{pmatrix} A \\ A \end{pmatrix}, & C = \begin{pmatrix} C \\ C \end{pmatrix} \\ \overline{A} = \begin{pmatrix} A \\ C \end{pmatrix}, & \overline{C} = \begin{pmatrix} C \\ A \end{pmatrix} \end{cases}$$

a missing label of a transition is set to the letter at the bottom of its ending state



is labelled by  $C$

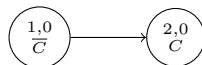
The automaton  $\mathcal{N}_{\text{ACC}} = \overline{\mathcal{M}}_{\text{ACC}} \otimes \mathcal{M}_{\text{ACC}}$  with matrix  $\mathbb{P}$



Notations for the transitions:

$$\begin{cases} A = \begin{pmatrix} A \\ A \end{pmatrix}, & C = \begin{pmatrix} C \\ C \end{pmatrix} \\ \overline{A} = \begin{pmatrix} A \\ C \end{pmatrix}, & \overline{C} = \begin{pmatrix} C \\ A \end{pmatrix} \end{cases}$$

a missing label of a transition is set to the letter at the bottom of its ending state



is labelled by  $C$

$$\mathfrak{P}_n = \mathbf{P}(S_n(1) \in \mathcal{A}^* b \mathcal{A}^* | S_n(0) \notin \mathcal{A}^* b \mathcal{A}^*) = \frac{V_{q'_0} \mathbb{P}^n V_{F'}^t}{1 - V_{q'_0} \mathbb{P}^n V_{\text{sink}}^t}$$

# Results for 5-mers of DNA

	BNN		BV		
	$E_{\text{BNN}}(T_{1000})/10^6$	Rank	$E_{\text{BV}}(T_{1000})/10^6$	Rank	$\frac{E_{\text{BNN}}(T_{1000})}{E_{\text{BV}}(T_{1000})}$
CCCCC	9,105	1021	6,304	1	1.44
GGGGG	9,570	1022	6,666	142	1.44
TTTTT	10,401	1023	7,457	993	1.39
AAAAA	10,656	1024	7,654	1024	1.39
CGCGC	7,047	699	6,446	11	1.09
TCCCC	7,076	737	6,477	17	1.09
CCCCT	7,076	738	6,477	21	1.09
GCGCG	7,127	787	6,518	31	1.09
CTCTC	7,263	883	6,679	148	1.09
...	...	...	...	...	...

$$\left\{ \begin{array}{l} 4\% \text{ of the 5-mers} \\ 0.2\% \text{ of the 7-mers} \\ 0.002\% \text{ of the 10-mers} \end{array} \right\} \quad \text{verify } \frac{E_{\text{BNN}}(T_{1000})}{E_{\text{BV}}(T_{1000})} > 1.05\%$$

# Numerical remarks

- ▶ **length** of promoters  $n \in [500 - 2000]$
- ▶ **Mutation probability**  $\pi = \max(\pi_{\alpha \rightarrow \beta}) \approx 10^{-9}$

# Numerical remarks

- ▶ **length** of promoters  $n \in [500 - 2000]$
- ▶ **Mutation probability**  $\pi = \max(\pi_{\alpha \rightarrow \beta}) \approx 10^{-9}$

We have

- ▶  $p_r$ : **probability of Mutation** to  $b$  from a  $r$ -neighbour of  $b$  with  $r \geq 2$   
$$p_r \leq n \times \pi^r \leq 2000 \times 10^{-18} < 2.10^{-6} \times \pi$$
- ▶  $q_s$ : **probability** that  $s$  **1-neighbors** simultaneously mutate to  $b$  with  $s \geq 2$   
$$q_s \leq n \times \pi^s \leq 2000 \times 10^{-18} < 2.10^{-6} \times \pi$$

## Numerical remarks

- ▶ **length** of promoters  $n \in [500 - 2000]$
- ▶ **Mutation probability**  $\pi = \max(\pi_{\alpha \rightarrow \beta}) \approx 10^{-9}$

We have

- ▶  $p_r$ : **probability of Mutation** to  $b$  from a  $r$ -neighbour of  $b$  with  $r \geq 2$   
$$p_r \leq n \times \pi^r \leq 2000 \times 10^{-18} < 2.10^{-6} \times \pi$$
- ▶  $q_s$ : **probability** that  $s$  **1-neighbors** simultaneously mutate to  $b$  with  $s \geq 2$   
$$q_s \leq n \times \pi^s \leq 2000 \times 10^{-18} < 2.10^{-6} \times \pi$$

Therefore assuming a **single mutation** in the promoter is **numerically sound**

# Clump approach

Assuming **a single mutation**

## Putative-hit positions.

- ▶ Given a **sequence**  $S(0)$  **not containing a  $k$ -mer**  $b$ ,
- ▶ a **putative-hit position** is any position of  $S(0)$  that can **lead by a mutation to an occurrence of  $b$  in  $S(1)$** ,
- ▶ where we assume that a **single** mutation has occurred.

$S(0) = \text{CCCAACAC}, \quad b = \text{ACC} \quad \rightsquigarrow \quad \underline{S}(0) = \underline{\text{C}}\text{CC}\underline{\text{A}}\text{ACAC},$

putative-hit positions **underlined** in  $\underline{S}(0)$ .

## Putative-hit positions.

- ▶ Given a **sequence**  $S(0)$  **not containing a  $k$ -mer**  $b$ ,
- ▶ a **putative-hit position** is any position of  $S(0)$  that can **lead by a mutation to an occurrence of  $b$  in  $S(1)$** ,
- ▶ where we assume that a **single** mutation has occurred.

$$S(0) = \text{CCCAACAC}, \quad b = \text{ACC} \quad \rightsquigarrow \quad \underline{S}(0) = \underline{\text{C}}\text{CC}\underline{\text{A}}\text{ACAC},$$

putative-hit positions **underlined** in  $\underline{S}(0)$ .

In a random sequence of length  $n$  with  $\mathcal{A} = \{\text{A}, \text{C}\}$ , let

- ▶  $H_{\text{A} \rightarrow \text{C}}^{(n)}$  number of putative-hit-positions  $\text{A} \rightarrow \text{C}$ ,
- ▶  $H_{\text{C} \rightarrow \text{A}}^{(n)}$  number of putative-hit-positions  $\text{C} \rightarrow \text{A}$ ,

Then

$$\mathfrak{P}_n \approx \mathbf{E}(H_{\text{A} \rightarrow \text{C}}^{(n)}) \times \pi_{\text{A} \rightarrow \text{C}} + \mathbf{E}(H_{\text{C} \rightarrow \text{A}}^{(n)}) \times \pi_{\text{C} \rightarrow \text{A}}$$

# Computing via generating functions

## Aim:

Compute

$$F_b(z, t_{A \rightarrow C}, t_{C \rightarrow A}) = \sum_{n \geq 0} \sum_{0 \leq i \leq n - |b|} \sum_{0 \leq j \leq n - |b|} f_{n,i,j} t_{A \rightarrow C}^i t_{C \rightarrow A}^j z^n$$

where  $f_{n,i,j}$  is the probability that a sequence  $S_n(0)$  **with no**  $b$ , of length  $n$ , contains

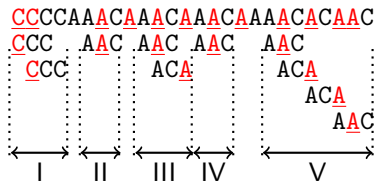
- ▶  $i$  putative-hit positions  $A \rightarrow C$
- ▶ and  $j$  putative-hit positions  $C \rightarrow A$

We have

$$\mathfrak{P}_n = [z^n] \left( \pi_{A \rightarrow C} \frac{\partial F(z, t_{A \rightarrow C}, 1)}{\partial t_{A \rightarrow C}} \Big|_{t_{A \rightarrow C}=1} + \pi_{C \rightarrow A} \frac{\partial F(z, 1, t_{C \rightarrow A})}{\partial t_{C \rightarrow A}} \Big|_{t_{C \rightarrow A}=1} \right)$$

# Putative-Hit-Positions and clump analysis

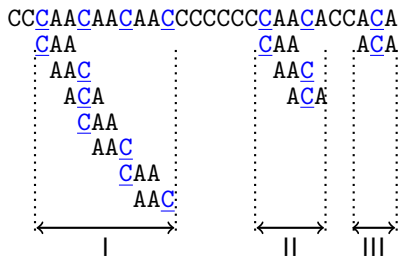
$$\mathcal{A} = \{\underline{A}, \underline{C}\} \quad b = \text{ACC} \longrightarrow d(\text{ACC}, 1) = \{\underline{C}\underline{C}\underline{C}, \underline{A}\underline{A}\underline{C}, \underline{A}\underline{C}\underline{A}\}$$



- (left)  $b = \text{ACC}$  - in clump I, when the right extension of a clump adds a new putative-hit position, this position is not necessarily in the extension, but possibly backwards left

# Putative-Hit-Positions and clump analysis

$$\mathcal{A} = \{A, C\} \quad b' = AAA \longrightarrow d(AAA, 1) = \{\underline{C}AA, A\underline{C}A, AA\underline{C}\}$$



- (right)  $b' = AAA$  - clump I contains 7 occurrences of  $d(AAA)$ , but only 4 putative-hit positions for  $b' = AAA$ . The number of word occurrences is not the relevant statistics for counting putative-hit positions

# Automata and Clumps

# Clumps of the set of words $\mathcal{U} = \{aaba, baab\}$

$\mathcal{E}_{w_1, w_2}$  correlation set from  $w_1$  to  $w_2$

$$\begin{array}{ll} \mathcal{E}_{aaba, aaba} = \{ba, aba\} & \mathcal{E}_{baab, baab} = \{aab\} \\ \mathcal{E}_{aaba, baab} = \{b\} & \mathcal{E}_{baab, aaba} = \{aa\} \end{array}$$

# Clumps of the set of words $\mathcal{U} = \{aabaa, baab\}$

$\mathcal{E}_{w_1, w_2}$  correlation set from  $w_1$  to  $w_2$

$$\begin{array}{ll} \mathcal{E}_{aabaa, aabaa} = \{baa, abaa\} & \mathcal{E}_{baab, baab} = \{aab\} \\ \mathcal{E}_{aabaa, baab} = \{b\} & \mathcal{E}_{baab, aabaa} = \{aa\} \end{array}$$

## Algorithm

1. Build the set of strings

$$\begin{aligned} X = & \{aabaa.(\epsilon + \mathcal{E}_{aabaa, aabaa})\} \cup \{aabaa.\mathcal{E}_{aabaa, baab}\} \\ & \cup \{baab.(\epsilon + \mathcal{E}_{baab, baab})\} \cup \{baab.\mathcal{E}_{baab, aabaa}\} \end{aligned}$$

# Clumps of the set of words $\mathcal{U} = \{aabaa, baab\}$

$\mathcal{E}_{w_1, w_2}$  correlation set from  $w_1$  to  $w_2$

$$\begin{array}{ll} \mathcal{E}_{aabaa, aabaa} = \{baa, abaa\} & \mathcal{E}_{baab, baab} = \{aab\} \\ \mathcal{E}_{aabaa, baab} = \{b\} & \mathcal{E}_{baab, aabaa} = \{aa\} \end{array}$$

## Algorithm

1. Build the set of strings

$$\begin{aligned} X = & \{aabaa.(\epsilon + \mathcal{E}_{aabaa, aabaa})\} \cup \{aabaa.\mathcal{E}_{aabaa, baab}\} \\ & \cup \{baab.(\epsilon + \mathcal{E}_{baab, baab})\} \cup \{baab.\mathcal{E}_{baab, aabaa}\} \end{aligned}$$

2. Build a trie  $\mathcal{T}$  on  $X$

# Clumps of the set of words $\mathcal{U} = \{aabaa, baab\}$

$\mathcal{E}_{w_1, w_2}$  correlation set from  $w_1$  to  $w_2$

$$\begin{array}{ll} \mathcal{E}_{aabaa, aabaa} = \{baa, abaa\} & \mathcal{E}_{baab, baab} = \{aab\} \\ \mathcal{E}_{aabaa, baab} = \{b\} & \mathcal{E}_{baab, aabaa} = \{aa\} \end{array}$$

## Algorithm

1. Build the set of strings

$$\begin{aligned} X = & \{aabaa.(\epsilon + \mathcal{E}_{aabaa, aabaa})\} \cup \{aabaa.\mathcal{E}_{aabaa, baab}\} \\ & \cup \{baab.(\epsilon + \mathcal{E}_{baab, baab})\} \cup \{baab.\mathcal{E}_{baab, aabaa}\} \end{aligned}$$

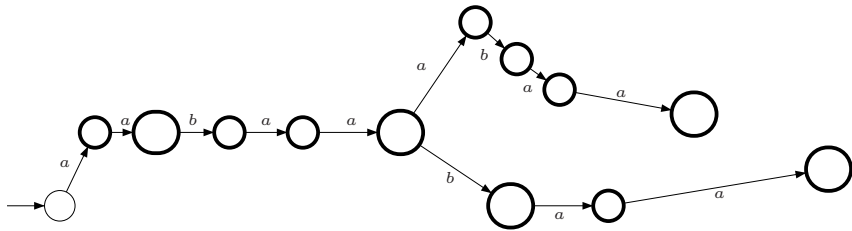
2. Build a trie  $\mathcal{T}$  on  $X$

3. Build a **Aho-Corasick like automaton** upon  $\mathcal{T}$ . For each node  $\nu$  of  $\mathcal{T}$  with “access word”  $v$ , use the transition function  $\delta$

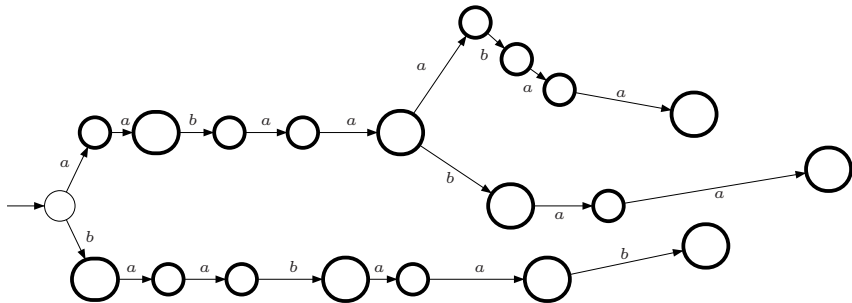
$\delta(\nu, \ell) =$  node accessed by the **longest prefix** in  $X$  that is **suffix** of  $v.\ell$

(Bassino-Clément-Fayolle-P.N. 2008)

$X = \{a b a a, a b a a b a a, a b a a a b a a, a b a a b\}$

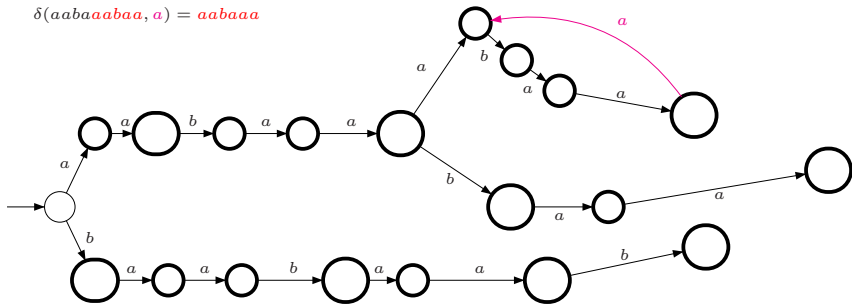


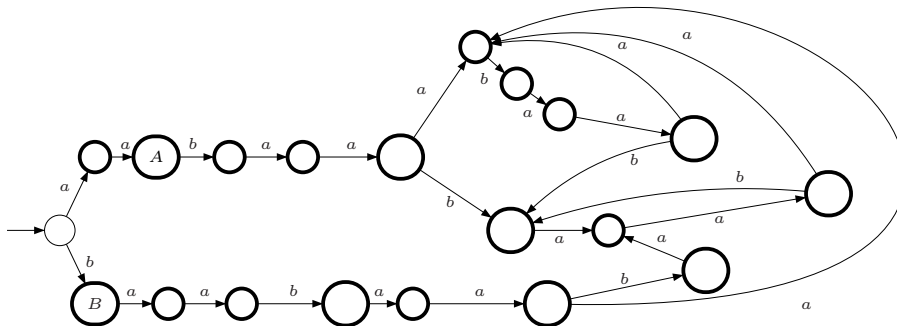
$X = \{a b a a, a a b a b a a, a a b a a b a a, a a b a a b, b a a b, b a a b a a b\}$



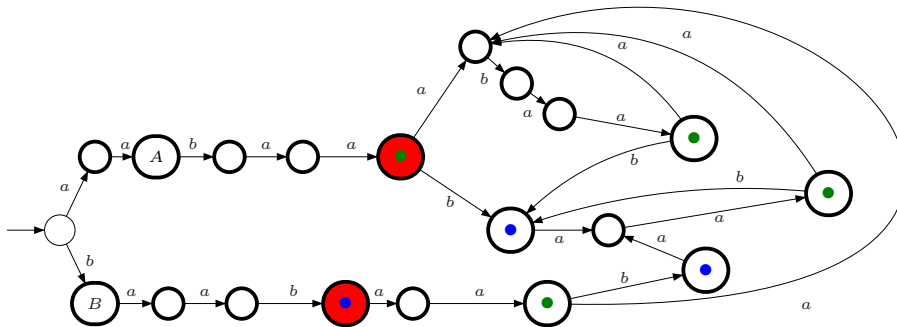
$X = \{a b a a, a a b a b a a, a a b a a b a a, a a b a a b, b a a b, b a a b a a b\}$

$\delta(a a b a a a b a a, a) = a a b a a a$



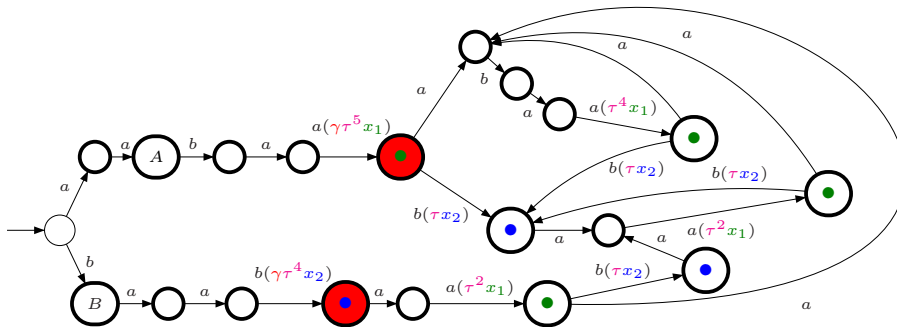


An automaton for  $\mathcal{V} = \{v_1 = aabaa, v_2 = baab\}$ . All transitions labeled by  $a$  and  $b$  ending respectively on state  $A$  and  $B$  are omitted.



An automaton for  $\mathcal{V} = \{v_1 = aabaa, v_2 = baab\}$ . All transitions labeled by  $a$  and  $b$  ending respectively on state  $A$  and  $B$  are omitted.

- ●, ● → the corresponding prefix (or state) ends with some occurrence of aabaa, baab.
- red states → states where we have entered a new clump



An automaton for  $\mathcal{V} = \{v_1 = aabaa, v_2 = baab\}$ . All transitions labeled by  $a$  and  $b$  ending respectively on state  $A$  and  $B$  are omitted.

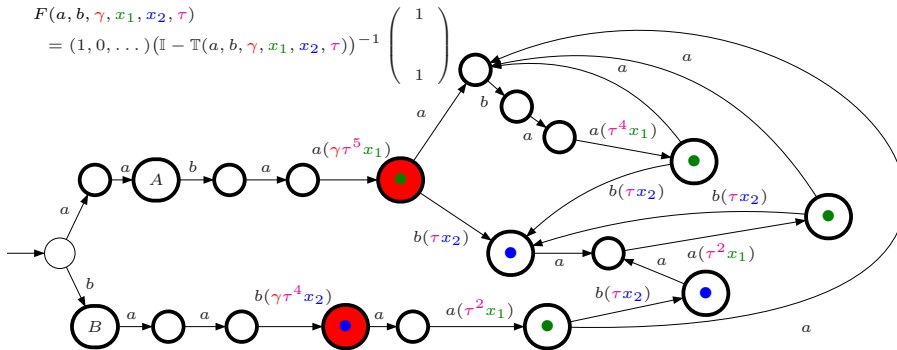
- ▶  $\bullet, \bullet \rightarrow$  the corresponding prefix (or state) ends with some occurrence of  $aabaa, baab$ .
- ▶ **red states**  $\rightarrow$  states where we have entered a **new clump**

Formal weights on transitions

- ▶  $\gamma \rightarrow$  the **number of clumps**
- ▶  $\tau \rightarrow$  total **length of clumps**
- ▶  $x_1, x_2 \rightarrow$  occurrences of  $aabaa, baab$

$$F(a, b, \gamma, x_1, x_2, \tau)$$

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



An automaton for  $\mathcal{V} = \{v_1 = aabaa, v_2 = baab\}$ . All transitions labeled by  $a$  and  $b$  ending respectively on state  $A$  and  $B$  are omitted.

- ▶  $\bullet, \bullet \rightarrow$  the corresponding prefix (or state) ends with some occurrence of  $aabaa, baab$ .
- ▶ **red states**  $\rightarrow$  states where we have entered a **new clump**

Formal weights on transitions

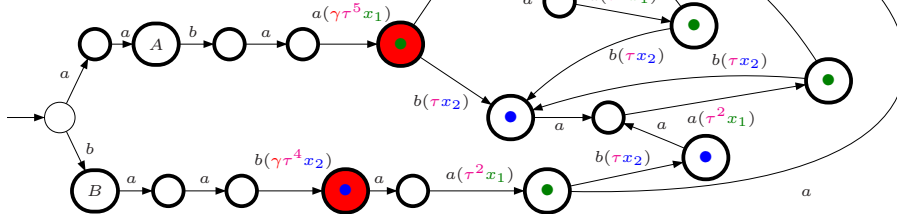
- ▶  $\gamma \rightarrow$  the **number of clumps**
- ▶  $\tau \rightarrow$  total **length of clumps**
- ▶  $x_1, x_2 \rightarrow$  occurrences of  $aabaa, baab$

$$F(a, b, \gamma, x_1, x_2, \tau)$$

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

$$a \rightsquigarrow \pi_a z, b \rightsquigarrow \pi_b z$$

$$[z^n] F(\pi_a z, \pi_b z, \dots) \rightarrow () \mathbb{T}^n(\pi_a, \pi_b, \dots)()$$



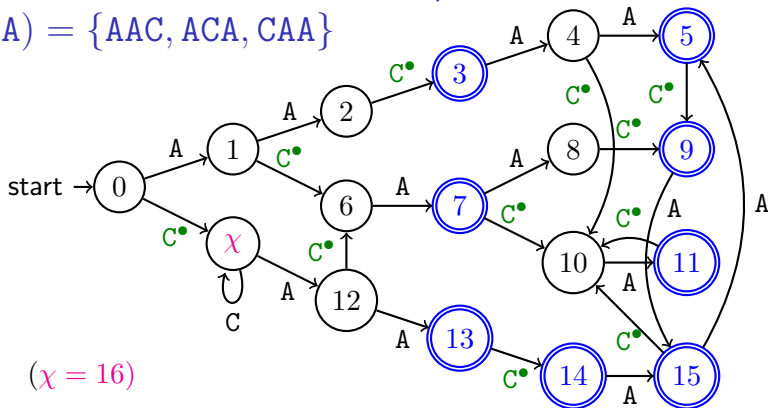
An automaton for  $\mathcal{V} = \{v_1 = aabaa, v_2 = baab\}$ . All transitions labeled by  $a$  and  $b$  ending respectively on state  $A$  and  $B$  are omitted.

- ▶  $\bullet, \bullet \rightarrow$  the corresponding prefix (or state) ends with some occurrence of  $aabaa, baab$ .
- ▶ **red states**  $\rightarrow$  states where we have entered a **new clump**

Formal weights on transitions

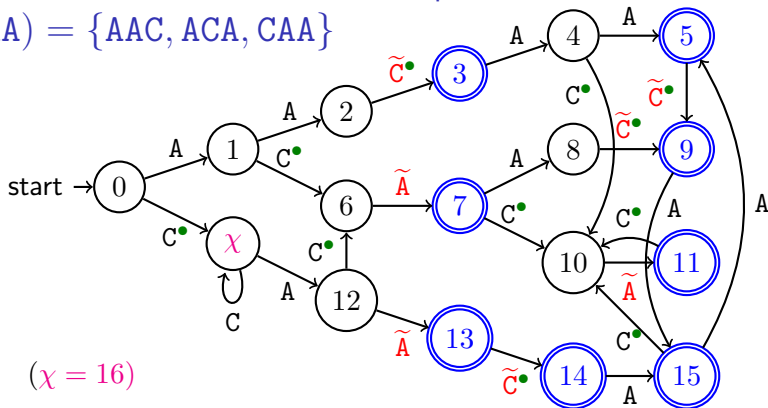
- ▶  $\gamma \rightarrow$  the **number of clumps**
- ▶  $\tau \rightarrow$  total **length of clumps**
- ▶  $x_1, x_2 \rightarrow$  occurrences of  $aabaa, baab$

# Automaton for constrained clumps of $d(\text{AAA}) = \{\text{AAC}, \text{ACA}, \text{CAA}\}$



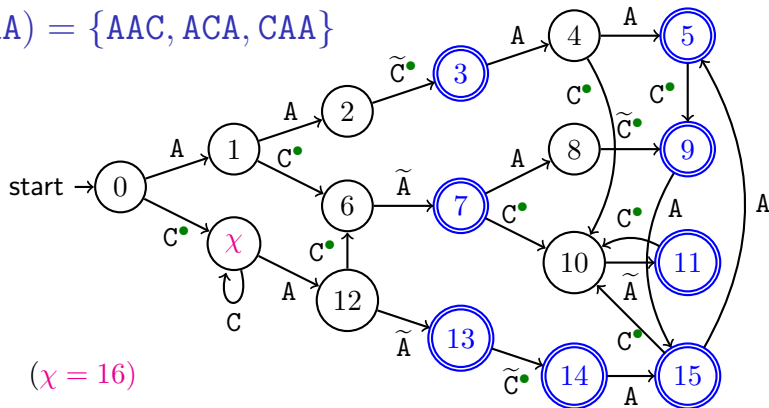
- ▶ **Double circles** signals an occurrence of a word of  $d(\text{aaa})$ .
- ▶ **Avoiding** AAA leads to **missing transitions** A
- ▶ The **missing transitions** C **point to the state**  $\chi$ .
- ▶ **• characters** mark **putative-hit-positions**

# Automaton for constrained clumps of $d(\text{AAA}) = \{\text{AAC}, \text{ACA}, \text{CAA}\}$



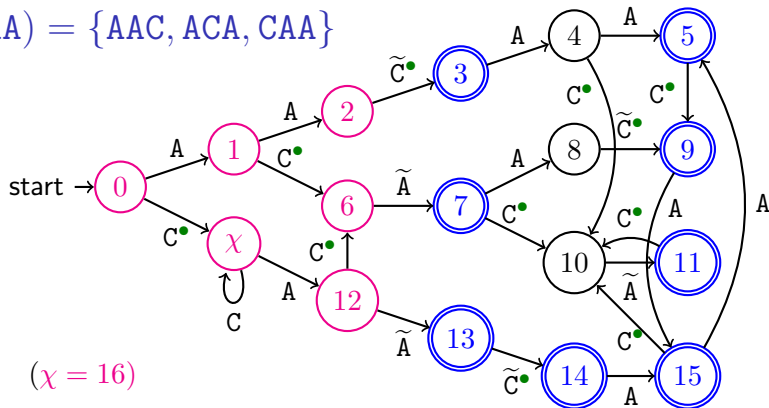
- ▶ **Double circles** signals an occurrence of a word of  $d(\text{aaa})$ .
- ▶ **Avoiding** AAA leads to **missing transitions** A
- ▶ The **missing transitions** C **point to the state**  $\chi$ .
- ▶ **• characters** mark **putative-hit-positions**
- ▶ Transitions covered by tildes ( $\tilde{A}, \tilde{C}$ ) emits a signal counting a putative-hit position.

# Automaton for constrained clumps of $d(AAA) = \{AAC, ACA, CAA\}$



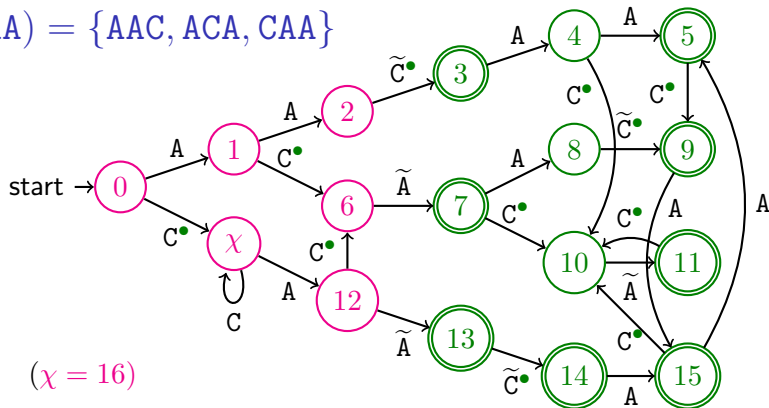
- $O = \{q, \delta(0, w) = q, w \in X\}$ , (**occurrence** of a word of  $d(aaa)$ ).

# Automaton for constrained clumps of $d(\text{AAA}) = \{\text{AAC}, \text{ACA}, \text{CAA}\}$



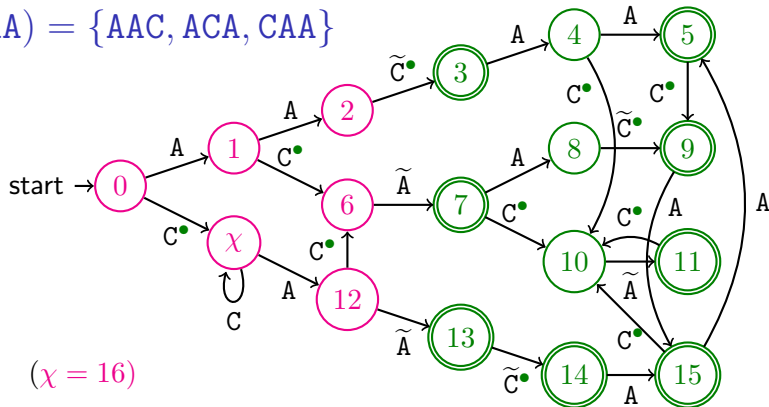
- ▶  $O = \{q, \delta(0, w) = q, w \in X\}$ , (**occurrence** of a word of  $d(\text{aaa})$ ).
- ▶  $\overline{E} = \{q, \delta(0, w) = q, w \in \widehat{\text{Pref}}(d(b))\}$ , with  $\widehat{\text{Pref}}(d(b))$  set of **strict prefixes** of words of  $d(b)$ .

# Automaton for constrained clumps of $d(AAA) = \{AAC, ACA, CAA\}$



- ▶  $O = \{q, \delta(0, w) = q, w \in X\}$ , (**occurrence** of a word of  $d(aaa)$ ).
- ▶  $\overline{E} = \{q, \delta(0, w) = q, w \in \widehat{\text{Pref}}(d(b))\}$ , with  $\widehat{\text{Pref}}(d(b))$  set of **strict prefixes** of words of  $d(b)$ .
- ▶ **Clump-Core** of the automaton  $E = Q \setminus \overline{E}$

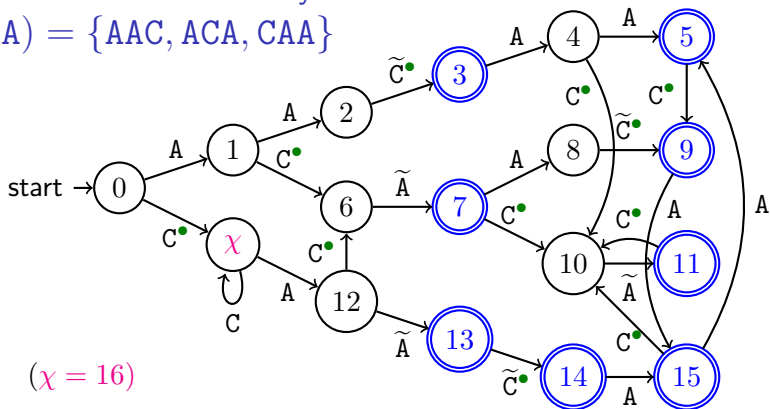
# Automaton for constrained clumps of $d(\text{AAA}) = \{\text{AAC}, \text{ACA}, \text{CAA}\}$



- ▶  $O = \{q, \delta(0, w) = q, w \in X\}$ , (**occurrence** of a word of  $d(\text{aaa})$ ).
- ▶  $\overline{E} = \{q, \delta(0, w) = q, w \in \widehat{\text{Pref}}(d(b))\}$ , with  $\widehat{\text{Pref}}(d(b))$  set of **strict prefixes** of words of  $d(b)$ .
- ▶ **Clump-Core** of the automaton  $E = Q \setminus \overline{E}$
- ▶ **Markov property:**  $\forall q \in E, |\{w \in \mathcal{A}^{[b]}; \delta(x, w) = q\}| = 1$

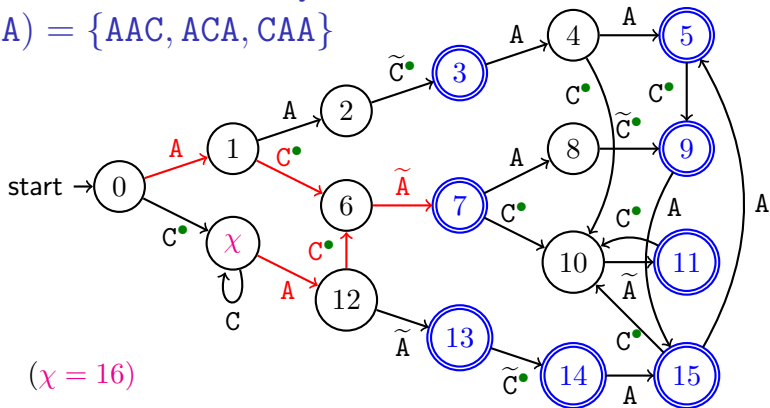
# Definition of an auxiliary function $\theta$

$$d(AAA) = \{AAC, ACA, CAA\}$$



# Definition of an auxiliary function $\theta$

$$d(AAA) = \{AAC, ACA, CAA\}$$

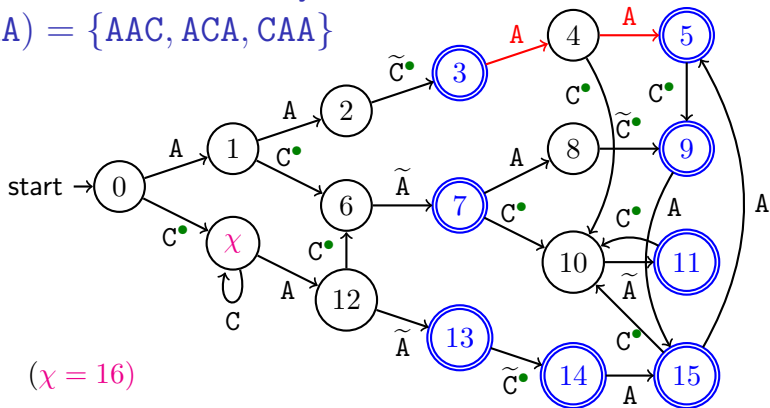


( $\chi = 16$ )

$$\theta(7) = \text{ACA}$$

# Definition of an auxiliary function $\theta$

$$d(AAA) = \{AAC, ACA, CAA\}$$

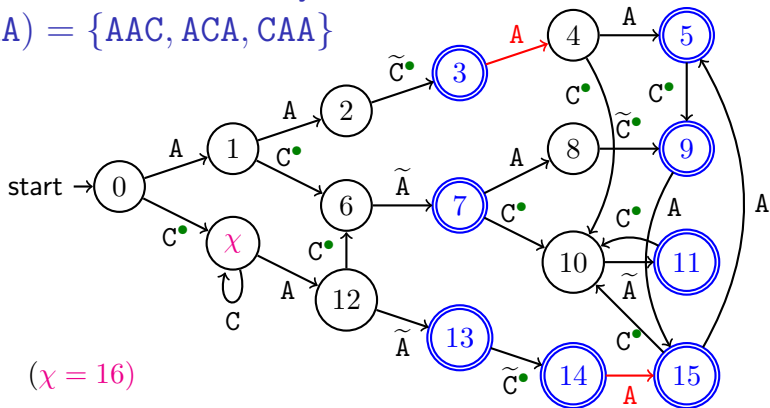


( $\chi = 16$ )

$$\theta(7) = \text{ACA} , \quad \theta(5) = \text{AA}$$

# Definition of an auxiliary function $\theta$

$$d(AAA) = \{AAC, ACA, CAA\}$$



( $\chi = 16$ )

$$\theta(7) = \text{ACA} , \quad \theta(5) = \text{AA} , \quad \theta(4) = \theta(15) = \text{A}$$

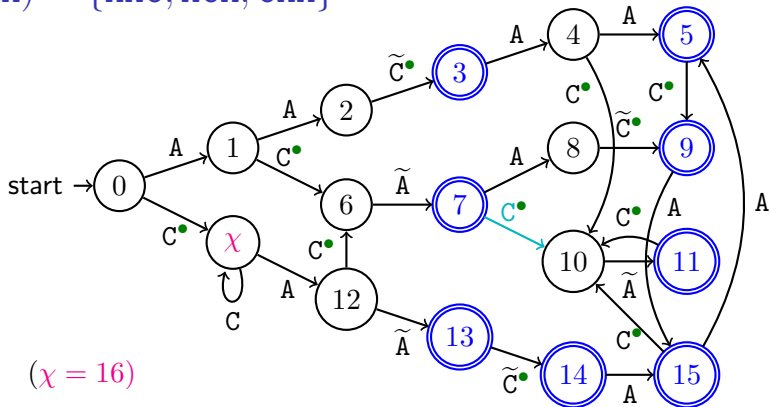
## Formal definition of $\theta$

For each state  $o \in O$  (recognizing an **occurrence** of  $d(b)$ ),

$$\theta(o) = \left\{ \begin{array}{l} \text{ } w \text{ with } |w| \leq |b|, \text{ of maximal length,} \\ \text{verifying} \left| \begin{array}{l} (a) \text{ there exists } q \text{ such that } \delta(q, w) = o, \\ (b) \text{ there is no } u \in \widehat{\text{Pref}}(w) \\ \text{such that } \delta(q, u) \in O \end{array} \right. \end{array} \right.$$

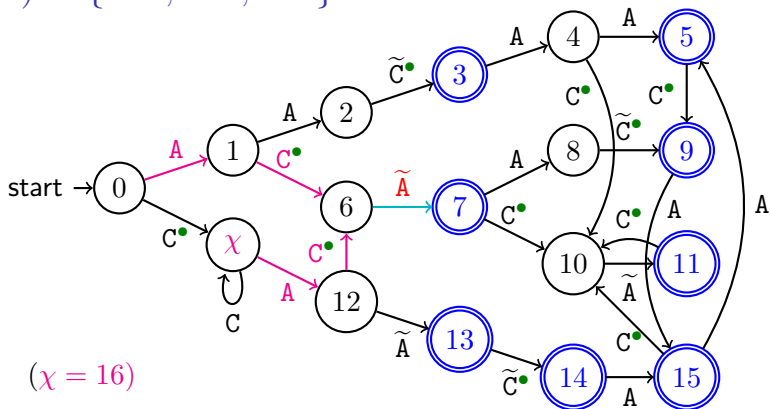
By the **Markov** property,  $\theta(o)$  defines a **unique word**

Adjacency matrix  $\mathbb{H}(t) = (h_{ij}(t))$

$$d(\text{AAA}) = \{\text{AAC}, \text{ACA}, \text{CAA}\}$$


►  $h_{7,10}(t_{\mathbf{C} \rightarrow \mathbf{A}}) = \nu_{\mathbf{C}}$

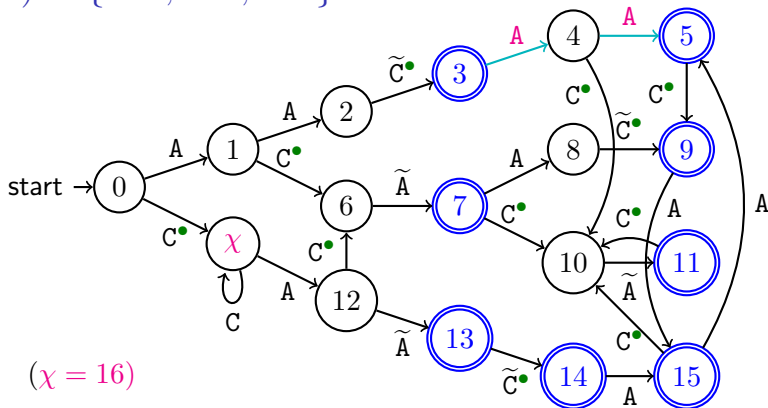
Adjacency matrix  $\mathbb{H}(t) = (h_{ij}(t))$   
 $d(\text{AAA}) = \{\text{AAC}, \text{ACA}, \text{CAA}\}$



- $h_{7,10}(t_{\text{C} \rightarrow \text{A}}) = \nu_{\text{C}}$
- $h_{6,7}(t_{\text{C} \rightarrow \text{A}}) = \nu_{\text{A}} t_{\text{C} \rightarrow \text{A}}$

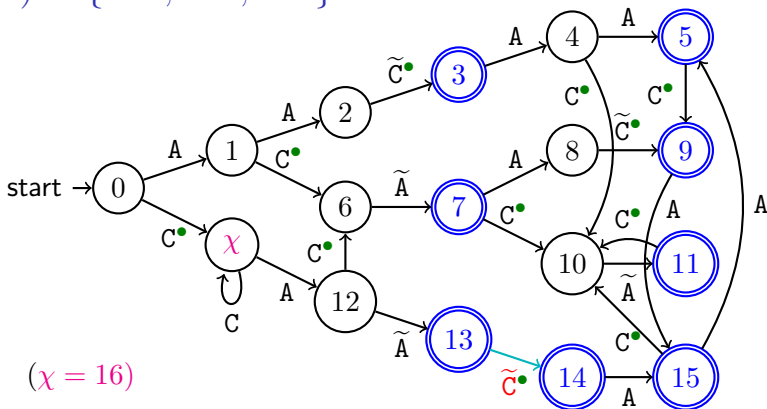
$$(\theta(7) = \text{A} \text{C}^{\bullet} \text{A})$$

Adjacency matrix  $\mathbb{H}(t) = (h_{ij}(t))$   
 $d(\text{AAA}) = \{\text{AAC}, \text{ACA}, \text{CAA}\}$



- ▶  $h_{7,10}(t_{\text{C} \rightarrow \text{A}}) = \nu_{\text{C}}$
- ▶  $h_{6,7}(t_{\text{C} \rightarrow \text{A}}) = \nu_{\text{A}} t_{\text{C} \rightarrow \text{A}}$  ( $\theta(7) = \text{AC}^{\bullet}\text{A}$ )
- ▶  $h_{3,4}(t_{\text{C} \rightarrow \text{A}}) = h_{4,5}(t_{\text{C} \rightarrow \text{A}}) = \nu_{\text{A}}$  ( $\theta(5) = \text{AA}$ )

Adjacency matrix  $\mathbb{H}(t) = (h_{ij}(t))$   
 $d(\text{AAA}) = \{\text{AAC}, \text{ACA}, \text{CAA}\}$



- ▶  $h_{7,10}(t_{C \rightarrow A}) = \nu_C$
- ▶  $h_{6,7}(t_{C \rightarrow A}) = \nu_A t_{C \rightarrow A}$
- ▶  $h_{3,4}(t_{C \rightarrow A}) = h_{4,5}(t_{C \rightarrow A}) = \nu_A$
- ▶  $h_{13,14}(t_{C \rightarrow A}) = \nu_C t_{C \rightarrow A}$

$$(\theta(7) = A C^\bullet A)$$

$$(\theta(5) = A A)$$

$$(\theta(14) = C^\bullet)$$

## Formal definition of the adjacency matrix $\mathbb{H}(t)$

(a)  $h_{ij}(t) = 0$  if there is no transition from  $i$  to  $j$

(b) With  $\delta(i, \alpha) = j$ ,

$$h_{i,j}(t) = \begin{cases} \nu(\alpha) & \text{if } \begin{cases} j \notin O, \\ j \in O \text{ and } \theta(j) \text{ contains no putative-hit position} \end{cases} \\ \nu(\alpha) \times t & \text{elsewhere} \end{cases}$$

## Formal definition of the adjacency matrix $\mathbb{H}(t)$

(a)  $h_{ij}(t) = 0$  if there is no transition from  $i$  to  $j$

(b) With  $\delta(i, \alpha) = j$ ,

$$h_{i,j}(t) = \begin{cases} \nu(\alpha) & \text{if } \begin{cases} j \notin O, \\ j \in O \text{ and } \theta(j) \text{ contains no putative-hit position} \end{cases} \\ \nu(\alpha) \times t & \text{elsewhere} \end{cases}$$

## From matrix to generating function

$$\begin{aligned} F_b(z, t) &= (1, 0, \dots, 0) \times (\mathbb{I} + z\mathbb{H}(t) + \dots + z^n\mathbb{H}^n(t) + \dots) \times \mathbf{1}^t \\ &= (1, 0, \dots, 0) \times (\mathbb{I} - z\mathbb{H}(t))^{-1} \times \mathbf{1}^t. \end{aligned}$$

Entries of  $(\mathbb{I} - z\mathbb{H}(t))^{-1}$  rational functions in  $z$  and  $t$

# Rational functions and gfun

**rational** function  $\frac{f(z)}{g(z)}$   $\rightarrow$  **gfun**[diffeqtorec]  $\rightarrow$  **recurrence** equations

**recurrence** equations  $\rightarrow$  **gfun**[rectoproc]  $\rightarrow$  **procedure** **Proc**( $n$ )= $[z^n]\frac{f(z)}{g(z)}$

# Rational functions, Taylor coefficient of order $n$ and gfun

**rational** function  $\frac{f(z)}{g(z)} \rightarrow \text{gfun}[\text{diffeqtorec}] \rightarrow \text{recurrence equations}$

**recurrence equations**  $\rightarrow \text{gfun}[\text{rectoproc}] \rightarrow \text{procedure Proc}(n)=[z^n] \frac{f(z)}{g(z)}$

$$F_b(z, t) = \frac{P(z, t)}{Q(z, t)} \quad \text{and} \quad F_b(z, 1) = \sum_{n \geq 0} \widehat{f}_n^{(b)} z^n = \frac{P(z, 1)}{Q(z, 1)},$$

$$E(z) = \sum_n \eta_n z^n = \left. \frac{\partial}{\partial t} F_b(z, t) \right|_{t=1} = \frac{P'_t(z, 1)}{Q(z, 1)} - \frac{P(z, 1)Q'_t(z, 1)}{Q^2(z, 1)}$$

# Rational functions, Taylor coefficient of order $n$ and gfun

**rational** function  $\frac{f(z)}{g(z)} \rightarrow \text{gfun}[\text{diffeqtorec}] \rightarrow \text{recurrence}$  equations

**recurrence** equations  $\rightarrow \text{gfun}[\text{rectoproc}] \rightarrow \text{procedure Proc}(n)=[z^n] \frac{f(z)}{g(z)}$

$$F_b(z, t) = \frac{P(z, t)}{Q(z, t)} \quad \text{and} \quad F_b(z, 1) = \sum_{n \geq 0} \widehat{f}_n^{(b)} z^n = \frac{P(z, 1)}{Q(z, 1)},$$

$$E(z) = \sum_n \eta_n z^n = \left. \frac{\partial}{\partial t} F_b(z, t) \right|_{t=1} = \frac{P'_t(z, 1)}{Q(z, 1)} - \frac{P(z, 1)Q'_t(z, 1)}{Q^2(z, 1)}$$

where,  $P(z, t)$  and  $Q(z, t)$  are **polynoms**, and,  
in a **random** sequence  $S_n(0)$  of length  $n$  with **no occurrence** of  $b$ ,

# Rational functions, Taylor coefficient of order $n$ and gfun

**rational** function  $\frac{f(z)}{g(z)} \rightarrow \text{gfun}[\text{diffeqtores}] \rightarrow \text{recurrence}$  equations

**recurrence** equations  $\rightarrow \text{gfun}[\text{rectoproc}] \rightarrow \text{procedure Proc}(n)=[z^n] \frac{f(z)}{g(z)}$

$$F_b(z, t) = \frac{P(z, t)}{Q(z, t)} \quad \text{and} \quad F_b(z, 1) = \sum_{n \geq 0} \widehat{f}_n^{(b)} z^n = \frac{P(z, 1)}{Q(z, 1)},$$

$$E(z) = \sum_n \eta_n z^n = \left. \frac{\partial}{\partial t} F_b(z, t) \right|_{t=1} = \frac{P'_t(z, 1)}{Q(z, 1)} - \frac{P(z, 1)Q'_t(z, 1)}{Q^2(z, 1)}$$

where,  $P(z, t)$  and  $Q(z, t)$  are **polynoms**, and,  
in a **random** sequence  $S_n(0)$  of length  $n$  with **no occurrence** of  $b$ ,

$$\blacktriangleright \widehat{f}_n^{(b)} = \mathbf{P}(S_n(0)) = \mathbf{P}(\text{not going into sink})$$

# Rational functions, Taylor coefficient of order $n$ and gfun

**rational** function  $\frac{f(z)}{g(z)} \rightarrow \text{gfun}[\text{diffeqtores}] \rightarrow \text{recurrence equations}$

**recurrence equations**  $\rightarrow \text{gfun}[\text{rectoproc}] \rightarrow \text{procedure Proc}(n)=[z^n] \frac{f(z)}{g(z)}$

$$F_b(z, t) = \frac{P(z, t)}{Q(z, t)} \quad \text{and} \quad F_b(z, 1) = \sum_{n \geq 0} \widehat{f}_n^{(b)} z^n = \frac{P(z, 1)}{Q(z, 1)},$$

$$E(z) = \sum_n \eta_n z^n = \left. \frac{\partial}{\partial t} F_b(z, t) \right|_{t=1} = \frac{P'_t(z, 1)}{Q(z, 1)} - \frac{P(z, 1)Q'_t(z, 1)}{Q^2(z, 1)}$$

where,  $P(z, t)$  and  $Q(z, t)$  are **polynoms**, and,  
in a **random** sequence  $S_n(0)$  of length  $n$  with **no occurrence** of  $b$ ,

- ▶  $\widehat{f}_n^{(b)} = \mathbf{P}(S_n(0)) = \mathbf{P}(\text{not going into sink})$
- ▶  $\eta_n$  is the **unconditionned probability** of the expectation of the count of putative-hit positions

# Rational functions, Taylor coefficient of order $n$ and gfun

**rational** function  $\frac{f(z)}{g(z)} \rightarrow \text{gfun}[\text{diffeqtores}] \rightarrow \text{recurrence}$  equations

**recurrence** equations  $\rightarrow \text{gfun}[\text{rectoproc}] \rightarrow \text{procedure Proc}(n)=[z^n] \frac{f(z)}{g(z)}$

$$F_b(z, t) = \frac{P(z, t)}{Q(z, t)} \quad \text{and} \quad F_b(z, 1) = \sum_{n \geq 0} \widehat{f}_n^{(b)} z^n = \frac{P(z, 1)}{Q(z, 1)},$$

$$E(z) = \sum_n \eta_n z^n = \left. \frac{\partial}{\partial t} F_b(z, t) \right|_{t=1} = \frac{P'_t(z, 1)}{Q(z, 1)} - \frac{P(z, 1)Q'_t(z, 1)}{Q^2(z, 1)}$$

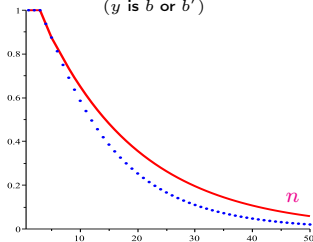
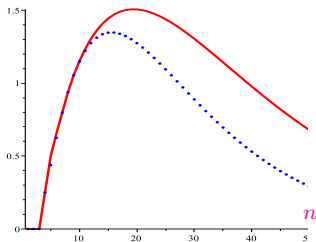
where,  $P(z, t)$  and  $Q(z, t)$  are **polynomials**, and,  
in a **random** sequence  $S_n(0)$  of length  $n$  with **no occurrence** of  $b$ ,

- ▶  $\widehat{f}_n^{(b)} = \mathbf{P}(S_n(0)) = \mathbf{P}(\text{not going into sink})$
- ▶  $\eta_n$  is the **unconditionned probability** of the expectation of the count of putative-hit positions
- ▶ **Conditionned expectation:**  $\widetilde{\eta}_n = \eta_n / \widehat{f}_n^{(b)}$

# An unexpected behaviour

$$\eta_n = \mathbf{E}(H_n^{(A \rightarrow C)}) + \mathbf{E}(H_n^{(C \rightarrow A)})$$

$$\widehat{f}_n^{(y)} = \mathbf{P}(|S_n(0)|_y = 0) \\ (y \text{ is } b \text{ or } b')$$



$$b = \text{ACAC} \quad b' = \text{AACC}$$

$$\nu(A) = \nu(C) = \frac{1}{2}$$

$$\mathbf{E}(H_n^{(A \rightarrow C)}) + \mathbf{E}(H_n^{(C \rightarrow A)}) = \left. \frac{\partial F_b(z, t)}{\partial t} \right|_{t=1}$$

$$\pi_{A \rightarrow C} = \pi_{C \rightarrow A}$$

$$t = t_{A \rightarrow C} = t_{C \rightarrow A}$$

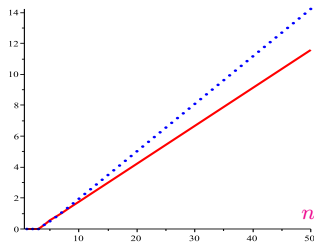
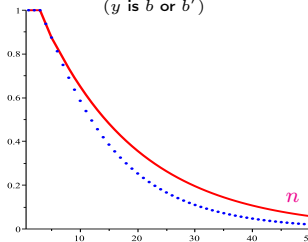
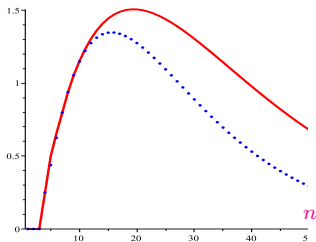
$$\pi_{A \rightarrow A} = \pi_{C \rightarrow C}$$

# An unexpected behaviour

$$\eta_n = \mathbf{E}(H_n^{(A \rightarrow C)}) + \mathbf{E}(H_n^{(C \rightarrow A)})$$

$$\hat{f}_n^{(y)} = \mathbf{P}(|S_n(0)|_y = 0) \\ (y \text{ is } b \text{ or } b')$$

$$\tilde{\eta}_n = \eta_n / \hat{f}_n^{(y)}$$



$$b = \text{ACAC} \quad b' = \text{AACC}$$

$$\nu(A) = \nu(C) = \frac{1}{2}$$

$$\mathbf{E}(H_n^{(A \rightarrow C)}) + \mathbf{E}(H_n^{(C \rightarrow A)}) = \left. \frac{\partial F_b(z, t)}{\partial t} \right|_{t=1}$$

$$\pi_{A \rightarrow C} = \pi_{C \rightarrow A}$$

$$t = t_{A \rightarrow C} = t_{C \rightarrow A}$$

$$\pi_{A \rightarrow A} = \pi_{C \rightarrow C}$$

## A proof by singularity analysis

$$F_b(z, t) = \frac{P(z, t)}{Q(z, t)} \quad P(z, t) \text{ and } Q(z, t) \text{ polynomials}$$

$$F_b(z, 1) = \sum_{n \geq 0} \hat{f}_n^{(b)} z^n = \frac{P(z, 1)}{Q(z, 1)}$$

$\hat{f}_n^{(b)}$  probability that  $S_n(0)$  has **no occurrence** of  $b$ .

$$E(z) = \sum_{n \geq 0} \mathbf{E}(H_n) z^n = \frac{P'_x(z, 1)}{Q(z, 1)} - \frac{P(z, 1)Q'_x(z, 1)}{Q^2(z, 1)}$$

## A proof by singularity analysis

$$F_b(z, t) = \frac{P(z, t)}{Q(z, t)} \quad P(z, t) \text{ and } Q(z, t) \text{ polynomials}$$

$$F_b(z, 1) = \sum_{n \geq 0} \hat{f}_n^{(b)} z^n = \frac{P(z, 1)}{Q(z, 1)}$$

$\hat{f}_n^{(b)}$  probability that  $S_n(0)$  has **no occurrence** of  $b$ .

$$E(z) = \sum_{n \geq 0} \mathbf{E}(H_n) z^n = \frac{P'_x(z, 1)}{Q(z, 1)} - \frac{P(z, 1)Q'_x(z, 1)}{Q^2(z, 1)}$$

The **dominant singularity**  $\tau$  is the smallest positive solution of  $Q(z, 1) = 0$ . Use suitable Cauchy integrals

$$\begin{cases} \hat{f}_n^{(b)} = \psi \times \tau^{-(n-1)} (1 + \mathcal{O}(B^n)), & (B < 1) \\ \mathbf{E}(H_n) = [z^n] E(z) = \tau^{-n} (\phi_1 \times n + \phi_2) \times (1 + \mathcal{O}(B^n)) \end{cases}$$

$$\implies \mathbf{E}(\tilde{H}_n) = \frac{\mathbf{E}(H_n)}{\hat{f}_n^{(b)}} = (c_1 \times n + c_2) \times (1 + \mathcal{O}(B^n)), \quad (B < 1).$$

## General case

Compute  $F_b(z, t_{A \rightarrow C}, t_{A \rightarrow G}, t_{A \rightarrow T}, t_{C \rightarrow A}, \dots, t_{T \rightarrow C}, t_{T \rightarrow G})$

$$\widehat{f}_n^{(b)} = [z^n] F_b(z, 1, 1, \dots, 1, 1)$$

$$\mathfrak{P}_n \approx [z^n] \sum_{\alpha \neq \beta \in \{A, C, G, T\}} \frac{\partial F_b(z, 1, \dots, 1, \pi_{\alpha \rightarrow \beta} t_{\alpha \rightarrow \beta}, 1, \dots)}{\partial t_{\alpha \rightarrow \beta}} \Big|_{t_{\alpha \rightarrow \beta} = 1} / \widehat{f}_n^{(b)}$$

- ▶ The **dominant singularities** of **all the terms of the sum** are **equal** to the **dominant singularity** of  $F_b(z, 1, 1, \dots, 1, 1)$
- ▶  $\mathfrak{P}_n$  behaves **quasi-linearly**

# Formal Languages Approach

Assuming **a single mutation**

# Guibas-Odlyzko decomposition - occurrences of a word $u$

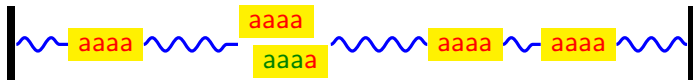
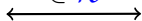
$$b = \text{aaaa}$$



# Guibas-Odlyzko decomposition - occurrences of a word $u$

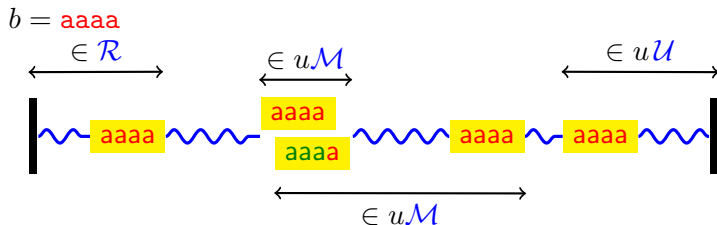
$b = \text{aaaa}$

$\in \mathcal{R}$

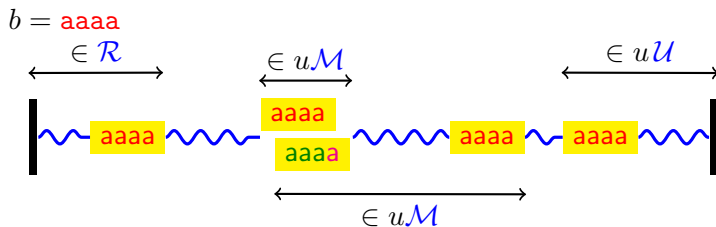




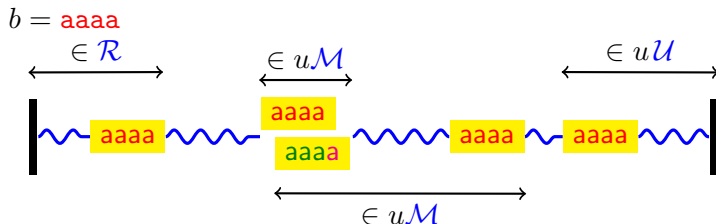
# Guibas-Odlyzko decomposition - occurrences of a word $u$



# Guibas-Odlyzko decomposition - occurrences of a word $u$

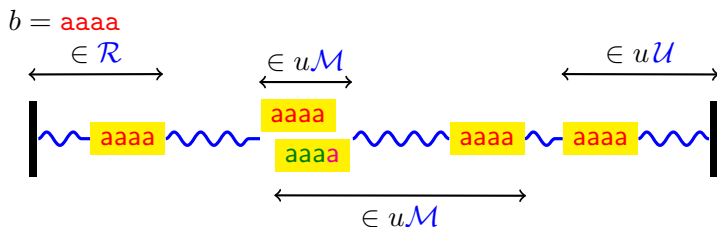


# Guibas-Odlyzko decomposition - occurrences of a word $u$



- ▶ The “**Right**” language  $\mathcal{R}$  associated to the word  $u$  is the set of words  $\mathcal{R} = \{r \mid r = e \cdot u \text{ and there is no } v \in U \text{ such that } r = xvy \text{ with } |y| > 0\}$ .
- ▶ The “**Minimal**” language  $\mathcal{M}$  leading from a word  $u$  to a word  $u$  is the set of words  $\mathcal{M} = \{m \mid u \cdot m = e \cdot u \text{ and there is no } v \in U \text{ such that } u \cdot m = xvy \text{ with } |x| > 0, |y| > 0\}$ .
- ▶ The “**Ultimate**” language  $\mathcal{U}$  of words following the last occurrence of the word  $u$  (such that this occurrence is the last occurrence of  $U$  in the text) is the set of words  $\mathcal{U} = \{u \mid \text{there is no } v \in U \text{ such that } u \cdot u = xvy \text{ with } |x| > 0\}$ .
- ▶ The “**Not**” language  $\mathcal{N}$  is the set of words with no occurrences of  $U$ ,  $\mathcal{N} = \{n \mid \text{there is no } v \in U \text{ such that } n = xvy\}$ .

# Guibas-Odlyzko decomposition - occurrences of a word $u$



$$\left. \begin{aligned} R(z) &= \frac{\mathbf{P}(b)z^{|b|}}{D(z)}, & M(z) &= 1 - \frac{1-z}{D(z)}, \\ U(z) &= \frac{1}{D(z)}, & Z(z) &= \frac{C(z)}{D(z)}, \end{aligned} \right| \text{ with } D(z) = (1-z)C(z) + \mathbf{P}(b)z^{|b|},$$

$\mathcal{C}$  **autocorrelation set** of the word  $b$

$$\mathcal{C} = \{w; \quad b.w = u.b, \quad 0 \leq |w| < |b|\}$$

$$C(z) = \sum_{w \in \mathcal{C}} \mathbf{P}(w)z^{|w|}$$

One mutation  $\rightsquigarrow$  One short clump in  $S_n(1)$

- ▶  $b$  without autocorrelation: **exactly one occurrence** of  $b$  in  $S_n(1)$ .
- ▶  $b$  with autocorrelation: **exactly one short clump**  $b.c$  with  $c \in \mathcal{C}$  in  $S_n(1)$ .
- ▶ the **positions of mutation** within a short clump are constrained

$X \neq A \quad b = AAAAA$        $\begin{array}{ccccccc} & & & 12345678 \\ & & & XXX\textcolor{blue}{A}\textcolor{red}{A}\textcolor{blue}{A}\textcolor{red}{A}\textcolor{blue}{A}\textcolor{red}{A}XXX \end{array}$

**exactly one mutation** occurred at **position 4** or **5**

## The right generating function $F(z)$

- ▶ We need **avoiding**  $b$  outside of the short clump in  $S_n(1)$   
 $b.c$  **short clump** with  $c \in \mathcal{C}$   $\rightsquigarrow$  **G.F.**  $G_c(z)$
- ▶  $G_c(z) = R(z) \times \mathbf{P}(c)z^{|c|} \times U(z) = \frac{\mathbf{P}(b.c)z^{|b.c|}}{D^2(z)}$
- ▶  $N(z) = \frac{C(z)}{D(z)}$  (zero occurrence of  $b$  in  $S_n(0)$ )

## The right generating function $F(z)$

- ▶ We need **avoiding**  $b$  outside of the short clump in  $S_n(1)$   
 $b.c$  **short clump** with  $c \in \mathcal{C}$   $\rightsquigarrow$  **G.F.**  $G_c(z)$
- ▶  $G_c(z) = R(z) \times \mathbf{P}(c)z^{|c|} \times U(z) = \frac{\mathbf{P}(b.c)z^{|b.c|}}{D^2(z)}$
- ▶  $N(z) = \frac{C(z)}{D(z)}$  (zero occurrence of  $b$  in  $S_n(0)$ )
- ▶  $D(z) = (1 - z)C(z) + \mathbf{P}(b)z^{|b|}$
- ▶  $G_c(z)$  and  $N(z)$  have the **same dominant singularity**  $\omega$

## The right generating function $F(z)$

- ▶ We need **avoiding**  $b$  outside of the short clump in  $S_n(1)$   
 $b.c$  **short clump** with  $c \in \mathcal{C}$   $\rightsquigarrow$  **G.F.**  $G_c(z)$
- ▶  $G_c(z) = R(z) \times \mathbf{P}(c)z^{|c|} \times U(z) = \frac{\mathbf{P}(b.c)z^{|b.c|}}{D^2(z)}$
- ▶  $N(z) = \frac{C(z)}{D(z)}$  (zero occurrence of  $b$  in  $S_n(0)$ )
- ▶  $D(z) = (1 - z)C(z) + \mathbf{P}(b)z^{|b|}$
- ▶  $G_c(z)$  and  $N(z)$  have the **same dominant singularity**  $\omega$
- ▶  $F(z) = \sum_{c \in \mathcal{C}} \frac{\mathbf{P}(b.c)z^{|b.c|}}{D^2(z)}$

## Asymptotics

$\mathfrak{q}_n$  approximation of  $\mathfrak{p}_n$   $\omega$  dominant singularity of  $D(z)$

$$\begin{aligned}\mathfrak{q}_n &= \frac{\mathbf{P}(b)}{C(\omega)D'(\omega)} \\ &\times \sum_{c \in \mathcal{C}} (|b| - |c|) \mathbf{P}(c) \omega^{|b \cdot c|} \\ &\times \left( (\textcolor{red}{n} - |b \cdot c| + 1) \omega^{-1} + \frac{D''(\omega)}{D'(\omega)} \right) + o(\mathbf{P}(b)).\end{aligned}$$

## An even more approximated result

$$D(z) = (1 - z)C(z) + \mathbf{P}(b)z^{|b|}$$

**by bootstrapping**  $\omega \approx 1 + \frac{\mathbf{P}(b)}{C(1) + |b|\mathbf{P}(b)} \approx 1$

Using  $\omega \approx 1$  gives

$$\mathbf{q}_n^{(\text{approx})} = \frac{\mathbf{P}(b)}{C^2(1)} \times \sum_{c \in \mathcal{C}} (|b| - |c|) \mathbf{P}(c) (\textcolor{red}{n} - |b.c| + 1)$$

# What about the approximations?

- ▶ The **distribution of letters**  $\nu'$  on  $S_n(1)$  is different of **the distribution**  $\nu$  on  $S_n(0)$
- ▶ There could have occurred **several mutations** within occurrences of  $b$  in  $S_n(1)$
- ▶ It could have occurred **parasites mutations** outside of the occurrences of  $b$  in  $S_n(1)$ .

## Distribution of letters in $S_n(1)$

$$\nu'(\alpha) = \nu(\alpha) \times p_{\alpha \rightarrow \alpha} + \sum_{\beta \neq \alpha} \nu(\beta) \times p_{\beta \rightarrow \alpha}$$

for small  $\epsilon_\alpha$ , we have  $\nu'(\alpha) = \nu(\alpha) \times (1 + \epsilon_\alpha)$ .

$$\eta = \max_{\alpha \in \mathcal{A}} \left| 1 - \frac{1}{1 - \epsilon_\alpha} \right|$$

for any language  $\mathcal{L} \in \mathcal{A}^n$ ,

$$(1 - \eta)^n \sum_{w \in \mathcal{L}} \Pr_{(\nu)}(w) \leq \sum_{w \in \mathcal{L}} \Pr_{(\nu')}(w) \leq (1 + \eta)^n \sum_{w \in \mathcal{L}} \Pr_{(\nu)}(w).$$

$$n \times \eta = o(1) \implies \sum_{w \in \mathcal{L}} \Pr_{(\nu')}(w) = \sum_{w \in \mathcal{L}} \Pr_{(\nu)}(w) \times (1 + \mathcal{O}(n\eta))$$

# Several mutations within occurrences of $b$ in $S_n(1)$

We can have **several short clumps**

1. Compute the generating function  $I(z)$  of sequences with isolated occurrences of  $b$
2. Starting from these isolated occurrences, extend the short clumps

(Bassino-Clément-Fayolle-P.N. 2008)

$$b = aaaaaa \quad \mathcal{C} - \{\epsilon\} = \{a, aa, aaa, aaaa\}$$

$$\text{Prefix code } \mathcal{K} = \mathcal{C}_o - \mathcal{C}_o \mathcal{A}^+ \quad \mathcal{C}_o = \mathcal{C} \setminus \epsilon \quad \mathcal{K} = \{a\}$$

$$\mathcal{M} = \{a, b(b + ab + aab + aaab + aaaaab)^*aaaaa\}$$

## Properties

- ▶  $\mathcal{K} \subset \mathcal{M}$
- ▶  $\mathcal{M} - \mathcal{K} = \mathcal{L}b$

G.F.  $I(z)$  of sequences with isolated occurrences of  $b$

$$\mathcal{I} = \mathcal{R}.(\mathcal{M} - \mathcal{K})^*.\mathcal{U}$$

$$K(z) = \sum_{w \in \mathcal{K}} \mathbf{P}(w)z^{|w|} \text{ is a polynomial}$$

$$\mathcal{M}(z) - K(z) = 1 - \frac{1-z}{D(z)} = \frac{p(z)\mathbf{P}(b)z^{|b|}}{D(z)} \text{ with } p(z) \text{ polynomial.}$$

$$I(z) = \frac{x\mathbf{P}(b)z^{|b|}}{D^2(z) \times \left(1 - \frac{xp(z)\mathbf{P}(b)z^{|b|}}{D(z)}\right)}$$

contribution  $\gamma_r$  of  $r$  clumps

$$\frac{\gamma_r}{\omega^n} = \frac{1}{\omega^n} [z^n] \frac{r! (\mathbf{P}(b)z^{|b|})^r p^{r-1}(z)}{((1-z)C(z) + \mathbf{P}(b)z^{|b|})^{r+1}} = \Theta(n^r)$$

## Parasite mutations

$$\begin{aligned} & \mathbf{P}(S_n(1) \text{after a parasite mutation} ) \\ & \leq \mathbf{P}(S_n(1) \text{before the mutation} ) \times \psi \end{aligned}$$

where  $\psi = \frac{\max_{\alpha, \beta \in \mathcal{A}; \alpha \neq \beta} p_{\alpha \rightarrow \beta}}{\min_{\alpha \in \mathcal{A}} p_{\alpha \rightarrow \alpha}}.$

**Theorem**[P.N. 2013]. The conditioned probability  $\mathbf{p}_n$  that a random sequence of length  $n$  that does not contain a  $k$ -mer  $b$  at time 0 evolves at time 1 to a random sequence that contains this  $k$ -mer verifies

$$\mathbf{p}_n = \mathbf{q}_n \times (1 + \mathcal{O}(n\psi)) + \mathcal{O}(n^2\psi^2)$$

where

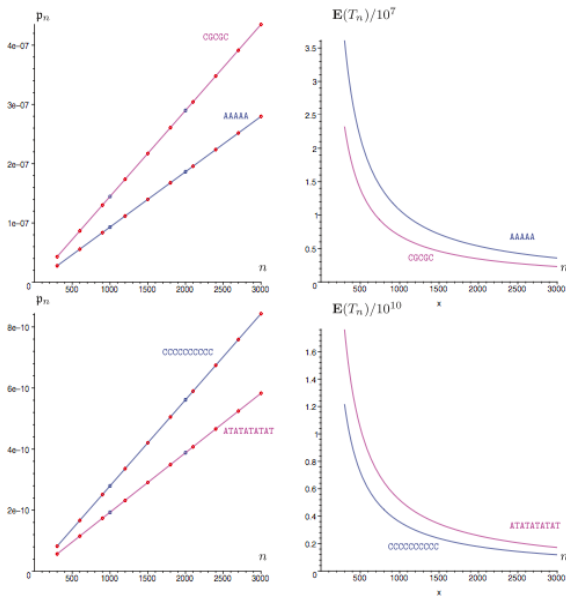
$$\begin{aligned} \mathbf{q}_n &= \frac{\mathbf{P}(b)}{C(\omega)D'(\omega)} \\ &\times \sum_{c \in \mathcal{C}} (|b| - |c|) \mathbf{P}(c) \omega^{|b \cdot c|} \\ &\times \left( (\mathbf{n} - |b \cdot c| + 1) \omega^{-1} + \frac{D''(\omega)}{D'(\omega)} \right) + o(\mathbf{P}(b)). \end{aligned}$$

$$\psi = \frac{\max_{\alpha, \beta \in \mathcal{A}; \alpha \neq \beta} p_{\alpha \rightarrow \beta}}{\min_{\alpha \in \mathcal{A}} p_{\alpha \rightarrow \alpha}}$$

# Numerical validation

$\mathcal{A} = \{A, C, G, T\}$  - uniform Bernoulli model for  $S(0)$ .

	$b = \text{AAAAA}$ and		for $\alpha \neq \beta$ , $p_{\alpha \rightarrow \beta} = 10^{-8}$	
Length $n$	$\mathfrak{p}_n \times 10^6$	$\mathfrak{h}_n \times 10^6$	$\mathfrak{q}_n \times 10^6$	$\mathfrak{q}_n^{(\text{approx})} \times 10^6$
10000	1.03335528	1.03335588	1.03335587	1.02703244
100000	10.3368481	10.3369021	10.3369021	10.2742439
10000000	1033.19278	1033.72699	1033.72698	1027.46750



**FIG. 5.** Plots of the probability  $p_n$  (left) and of the expected waiting time  $E(T_n)$  (right). (Top)  $b = \text{AAAAA}$  (blue) and  $b' = \text{CGCGC}$  (magenta). (Bottom)  $b = \text{CCCCC}$  (blue) and  $b' = \text{ATATATA}$  (magenta). In the linear plots of the probability, the anchors values for  $n = 1000$  and  $n = 2000$  (computed by automata) are represented by boxes; the straight lines are the straight lines going through the corresponding points and the circles are test values also computed by automata. The fit is perfect as expected from singularity analysis.

(from Behrens-Nicaud-P.N., JCB 19,5, 2012)