

DNA evolution, Automata and Clumps

Pierre Nicodème

CNRS - LIPN Team CALIN, University Paris-North, 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

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

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

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

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

α	$\nu(\alpha)$
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)

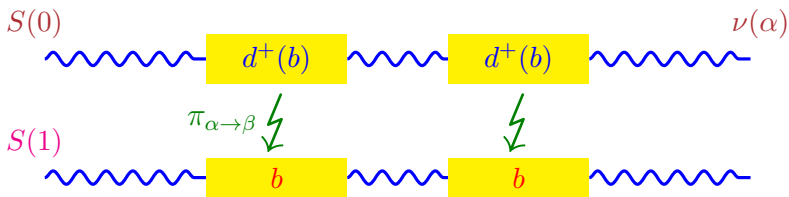
Plan of the talk

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)
 - ▶ **Non-efficient** approach by **clump analysis**, either by **combinatorics of words** or by **automata**.
 - ▶ **Proof by singularity analysis** of the **quasi-linear behaviour** of \mathfrak{P}_n

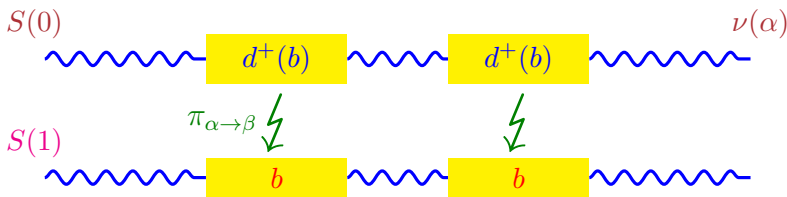
Behrens-Vingron 2010

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



Behrens-Vingron 2010

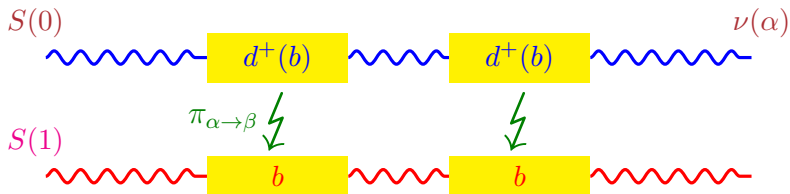
- $d^+(b)$ neighbors of b by substitution



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

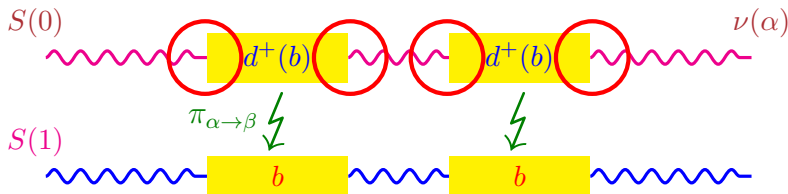
Approximations of Behrens-Vingron 2010

- 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 $S(0)$

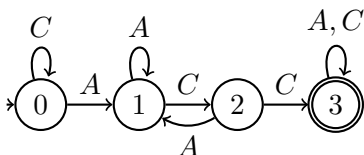


Behrens-Nicaud-P.N. 2012

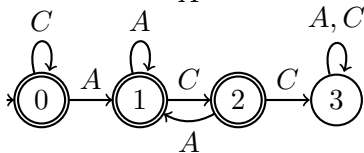
Construct an **automaton**

- ▶ 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



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

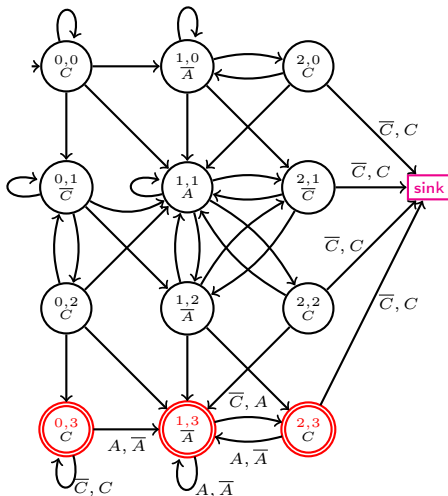


$$\overline{\mathcal{M}}_{\text{ACC}} = \{Q, \delta, s = 0, \textcolor{red}{Q} \setminus \textcolor{red}{F}\}$$

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

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

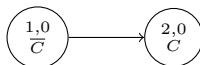
The automaton $\mathcal{N}_{\text{ACC}} = \overline{\mathcal{M}}_{\text{ACC}} \otimes \mathcal{M}_{\text{ACC}}$ with matrix \mathbb{P}



Notations for the transitions:

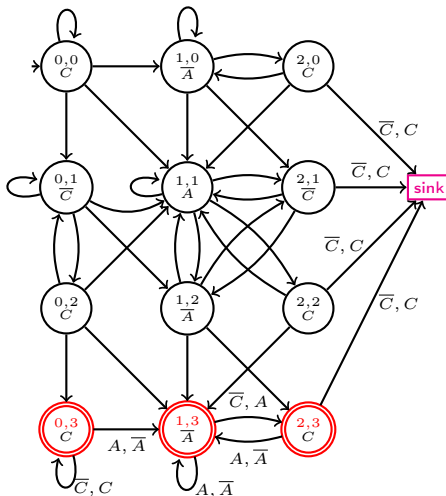
$$\begin{cases} A = \begin{pmatrix} A \\ A \end{pmatrix}, & C = \begin{pmatrix} C \\ C \end{pmatrix} \\ \overline{A} = \begin{pmatrix} A \\ C \end{pmatrix}, & \overline{C} = \begin{pmatrix} C \\ A \end{pmatrix} \end{cases}$$

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



is labelled by C

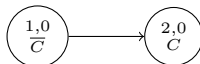
The automaton $\mathcal{N}_{\text{ACC}} = \overline{\mathcal{M}}_{\text{ACC}} \otimes \mathcal{M}_{\text{ACC}}$ with matrix \mathbb{P}



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		$\frac{E_{BNN}(T_{1000})}{E_{BV}(T_{1000})}$
	$E_{BNN}(T_{1000})/10^6$	Rank	$E_{BV}(T_{1000})/10^6$	Rank	
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
...

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

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 $r \geq 2$
 $p_r \leq n \times \pi^r \leq 2000 \times 10^{-18} < 2 \cdot 10^{-6} \times \pi$
- ▶ q_s : **probability** that s **1-neighbours** simultaneously mutate to b with $s \geq 2$
 $q_s \leq n \times \pi^s \leq 2000 \times 10^{-18} < 2 \cdot 10^{-6} \times \pi$

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 $r \geq 2$
$$p_r \leq n \times \pi^r \leq 2000 \times 10^{-18} < 2 \cdot 10^{-6} \times \pi$$
- ▶ q_s : **probability** that s **1-neighbours** simultaneously mutate to b with $s \geq 2$
$$q_s \leq n \times \pi^s \leq 2000 \times 10^{-18} < 2 \cdot 10^{-6} \times \pi$$

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

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

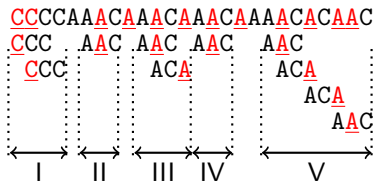
- ▶ i putative-hit positions $A \rightarrow C$
- ▶ and j putative-hit positions $C \rightarrow A$

We have

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

Putative-Hit-Positions and clump analysis

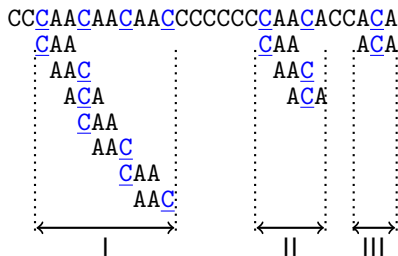
$$\mathcal{A} = \{\underline{A}, \underline{C}\} \quad b = \text{ACC} \longrightarrow d(\text{ACC}, 1) = \{\underline{C}\underline{C}\underline{C}, \underline{A}\underline{A}\underline{C}, \underline{A}\underline{C}\underline{A}\}$$



- (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

$$\mathcal{A} = \{\mathbf{A}, \mathbf{C}\} \quad b' = \mathbf{AAA} \longrightarrow d(\mathbf{AAA}, 1) = \{\mathbf{CAA}, \mathbf{ACA}, \mathbf{AAC}\}$$



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

| - Automaton approach

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

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

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

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

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

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

Algorithm

1. Build the set of strings

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

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

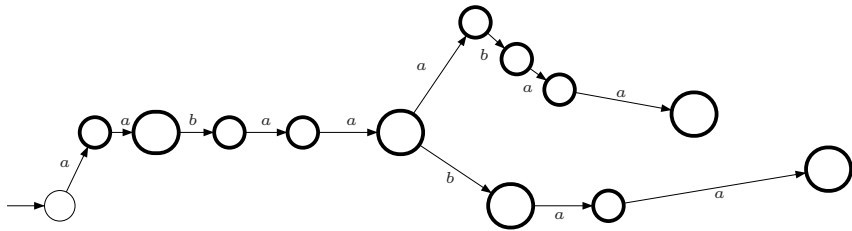
1. Build the set of strings

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

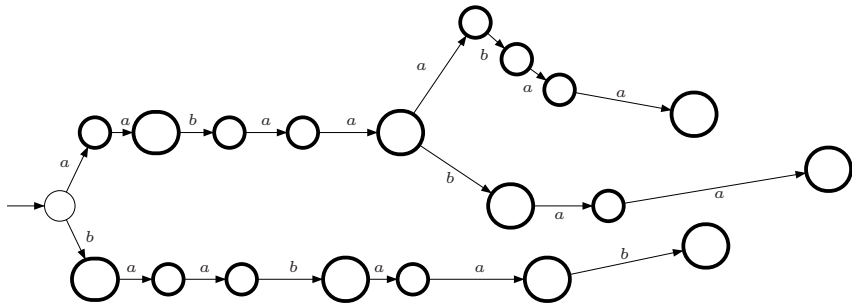
2. Build a trie \mathcal{T} on X
3. Build a **Aho-Corasick like automaton** upon \mathcal{T} . For each node ν of \mathcal{T} with “access word” v , use the transition function δ

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

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

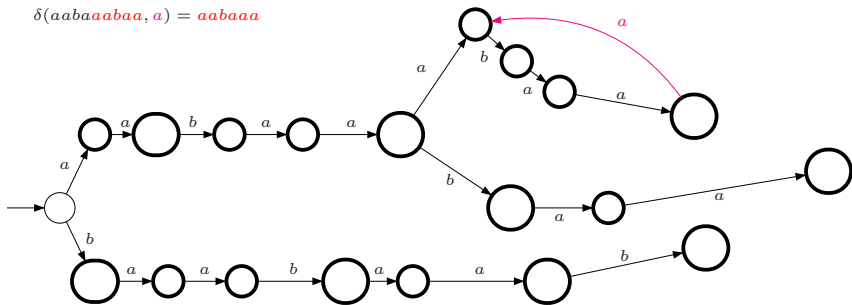


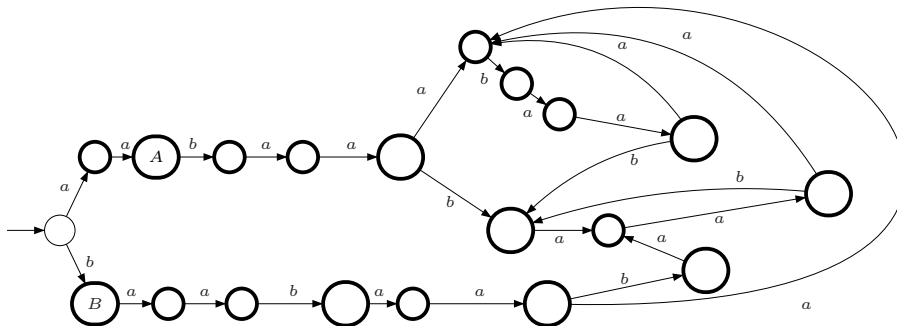
$X = \{a b a a, a a b a b a a, a a b a a b a a, a a b a a b, b a a b, b a a b a a b\}$



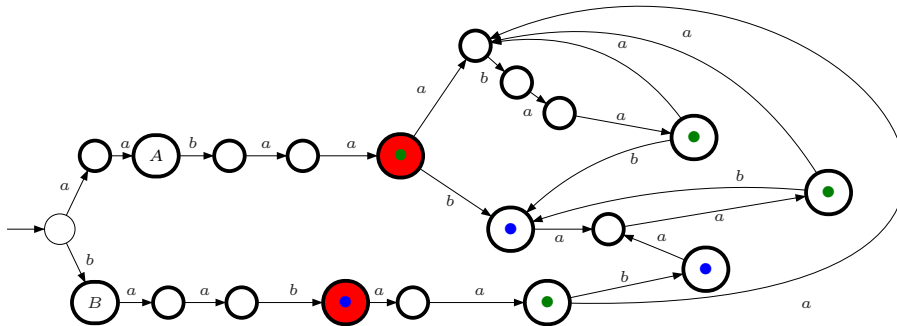
$X = \{a b a a, a a b a b a a, a a b a a b a a, a a b a a b, b a a b, b a a b a a b\}$

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



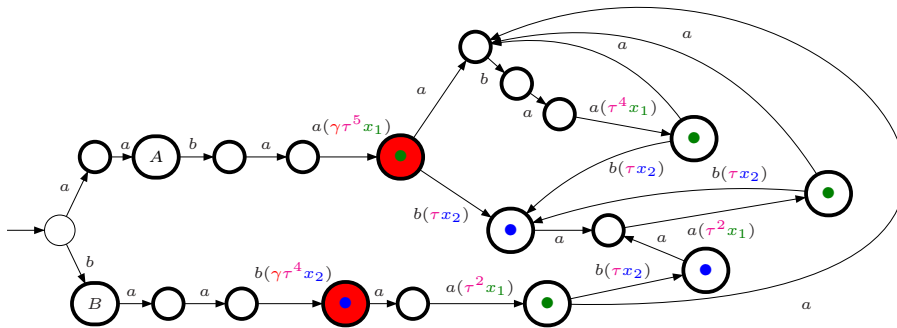


An automaton for $\mathcal{V} = \{v_1 = aabaa, v_2 = baab\}$. All transitions labeled by a and b ending respectively on state A and B are omitted.



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.

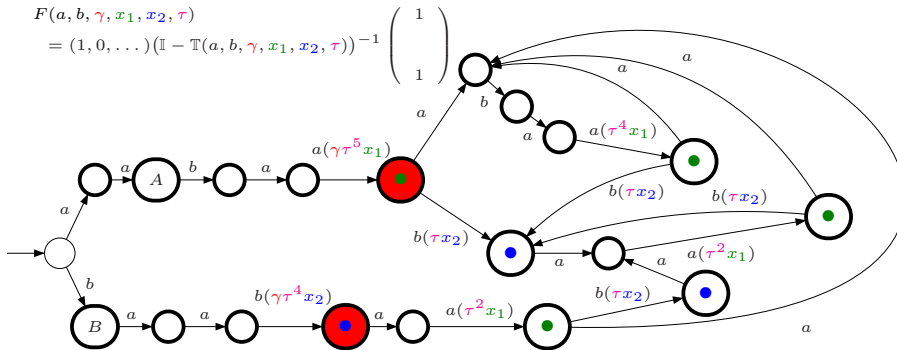
- ▶ $\bullet, \bullet \rightarrow$ the corresponding prefix (or state) ends with some occurrence of $aabaa, baab$.
- ▶ **red states** \rightarrow states where we have entered a **new clump**

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

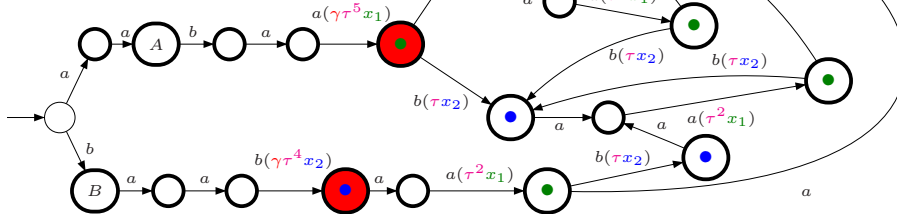
- ▶ $\gamma \rightarrow$ the **number of clumps**
- ▶ $\tau \rightarrow$ total **length of clumps**
- ▶ $x_1, x_2 \rightarrow$ occurrences of $aabaa, baab$

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

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

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

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



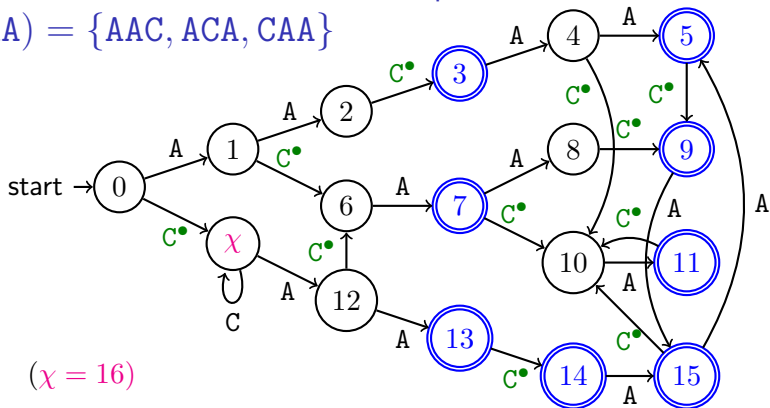
An automaton for $\mathcal{V} = \{v_1 = aabaa, v_2 = baab\}$. All transitions labeled by a and b ending respectively on state A and B are omitted.

- ▶ $\bullet, \bullet \rightarrow$ 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

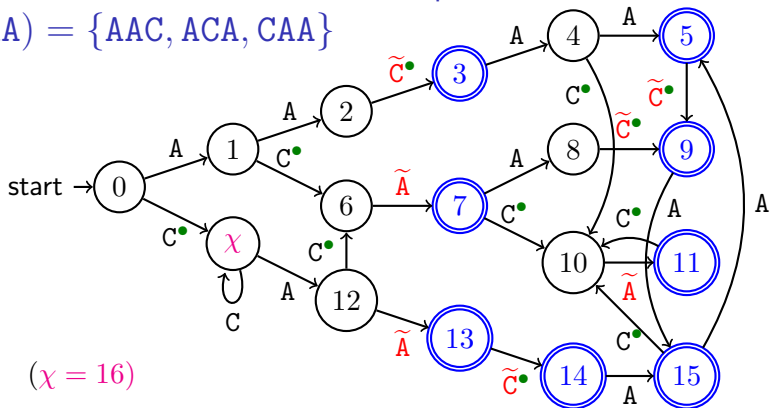
- ▶ $\gamma \rightarrow$ the **number of clumps**
- ▶ $\tau \rightarrow$ total **length of clumps**
- ▶ $x_1, x_2 \rightarrow$ occurrences of $aabaa, baab$

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



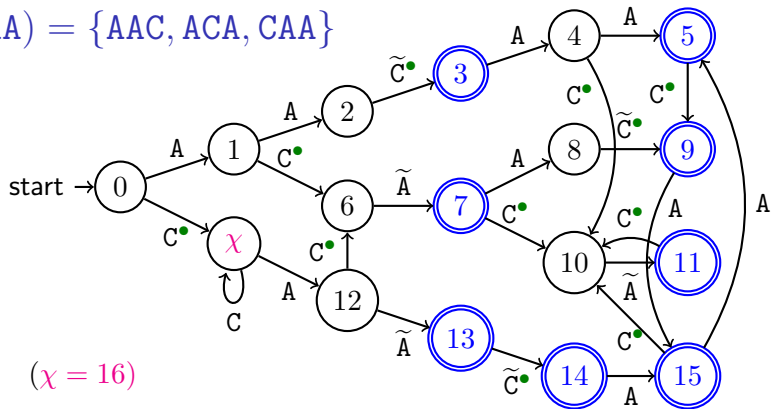
- ▶ **Double circles** signals an occurrence of a word of $d(aaa)$.
- ▶ **Avoiding** AAA leads to **missing transitions** A
- ▶ The **missing transitions** C **point to the state** χ .
- ▶ **• characters** mark **putative-hit-positions**

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



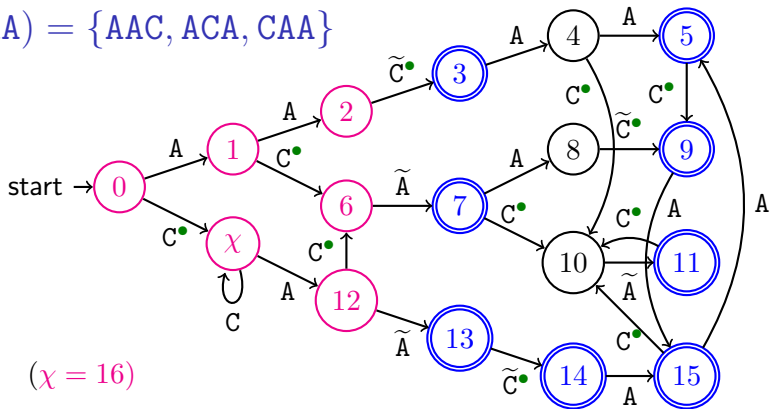
- ▶ **Double circles** signals an occurrence of a word of $d(\text{aaa})$.
- ▶ **Avoiding** AAA leads to **missing transitions** A
- ▶ The **missing transitions** C **point to the state** χ .
- ▶ **• characters** mark **putative-hit-positions**
- ▶ Transitions covered by tildes (\tilde{A}, \tilde{C})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

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



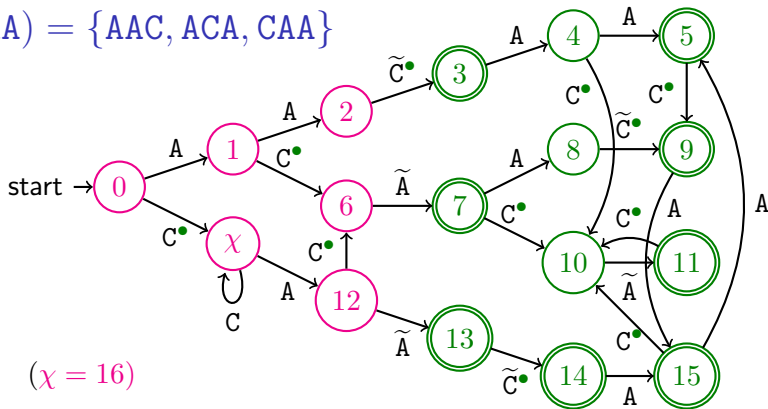
- $O = \{q, \delta(0, w) = q, w \in X\}$, (**occurrence** of a word of $d(aaa)$).

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



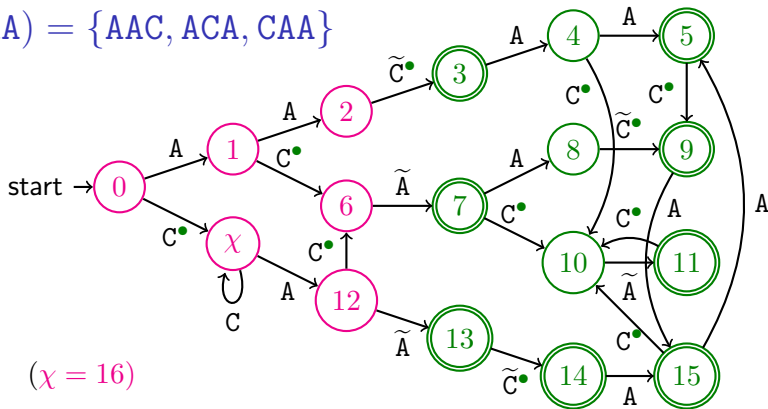
- ▶ $O = \{q, \delta(0, w) = q, w \in X\}$, (**occurrence** of a word of $d(\text{aaa})$).
- ▶ $\overline{E} = \{q, \delta(0, w) = q, w \in \widehat{\text{Pref}}(d(b))\}$, with $\widehat{\text{Pref}}(d(b))$ set of **strict prefixes** of words of $d(b)$.

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



- ▶ $O = \{q, \delta(0, w) = q, w \in X\}$, (**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 $E = Q \setminus \overline{E}$

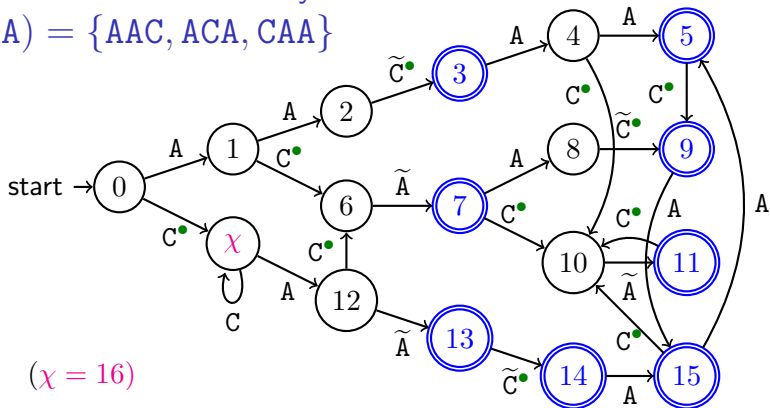
Automaton for constrained clumps of $d(\text{AAA}) = \{\text{AAC}, \text{ACA}, \text{CAA}\}$



- ▶ $O = \{q, \delta(0, w) = q, w \in X\}$, (**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 θ

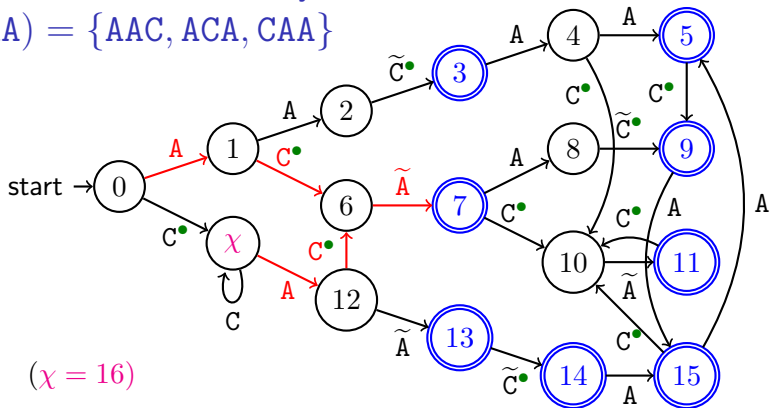
$$d(AAA) = \{AAC, ACA, CAA\}$$



$(\chi = 16)$

Definition of an auxiliary function θ

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

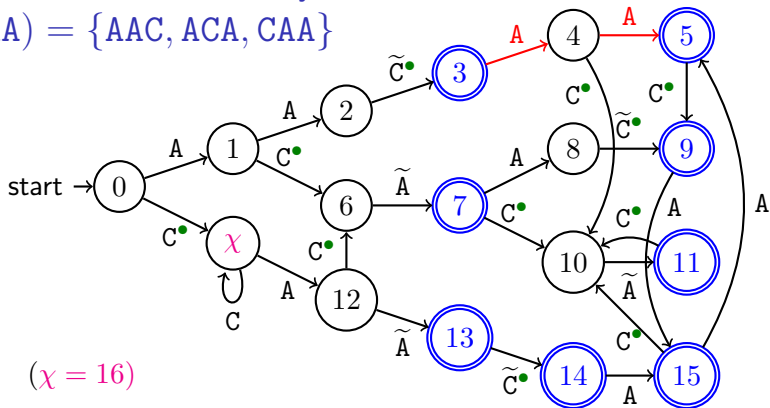


($\chi = 16$)

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

Definition of an auxiliary function θ

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

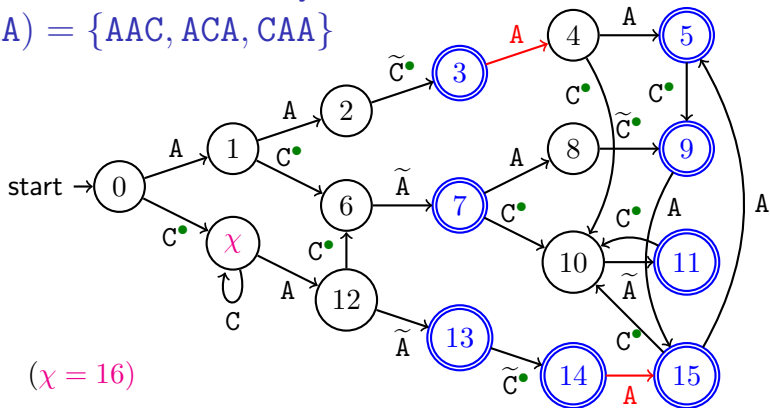


($\chi = 16$)

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

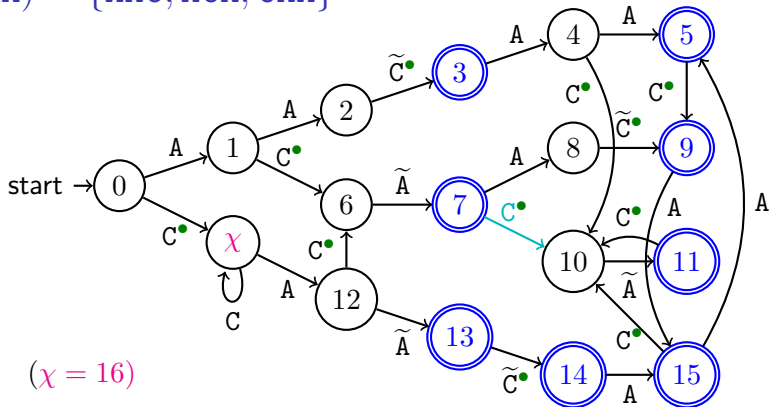
Definition of an auxiliary function θ

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



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

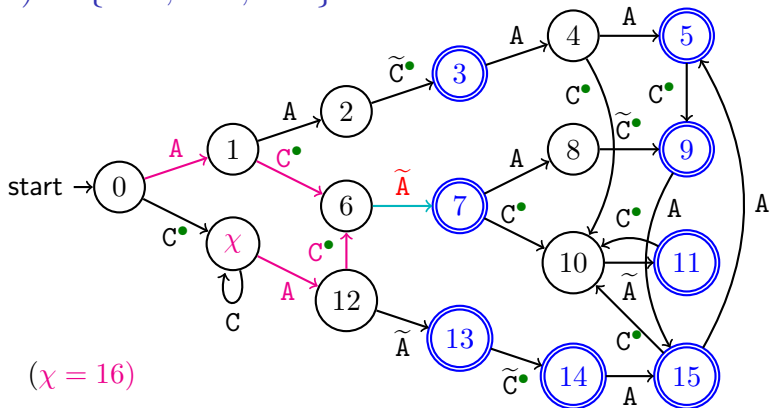
Adjacency matrix $\mathbb{H}(t) = (h_{ij}(t))$
 $d(\text{AAA}) = \{\text{AAC}, \text{ACA}, \text{CAA}\}$



($\chi = 16$)

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

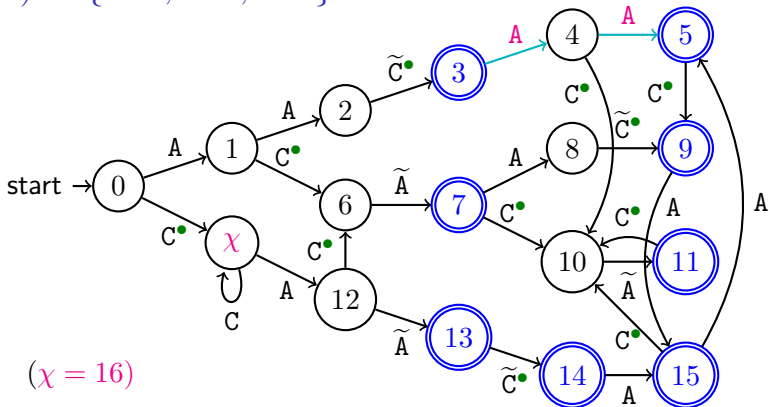
Adjacency matrix $\mathbb{H}(t) = (h_{ij}(t))$
 $d(\text{AAA}) = \{\text{AAC}, \text{ACA}, \text{CAA}\}$



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

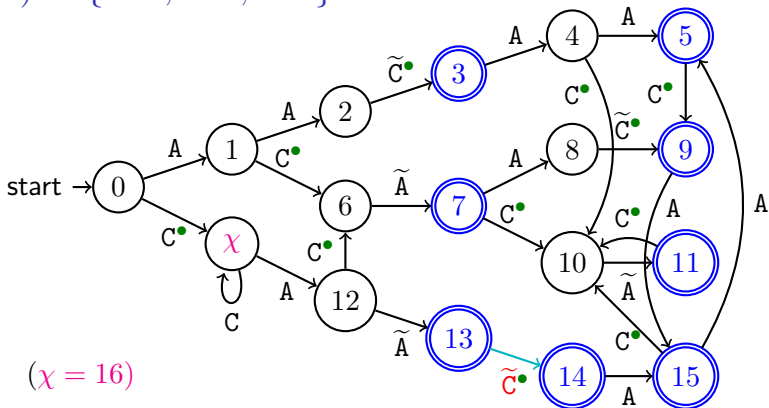
$$(\theta(7) = \text{A} \text{C}^{\bullet} \text{A})$$

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 $\frac{f(z)}{g(z)} \rightarrow \text{gfun}[\text{diffeqtores}] \rightarrow \text{recurrence}$ 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 ,

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 ,

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

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 ,

- ▶ $\widehat{f}_n^{(b)} = \mathbf{P}(S_n(0)) = \mathbf{P}(\text{not going into sink})$
- ▶ η_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 $\rightarrow \text{gfun}[\text{rectoproc}] \rightarrow \text{procedure Proc}(n)=[z^n] \frac{f(z)}{g(z)}$

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

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

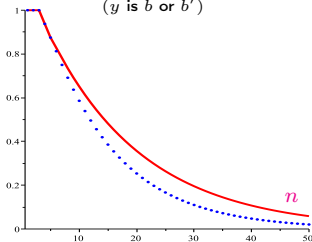
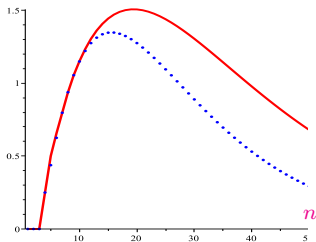
where, $P(z, t)$ and $Q(z, t)$ are **polynoms**, and,
in a **random** sequence $S_n(0)$ of length n with **no occurrence** of b ,

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

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

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



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

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

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

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

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

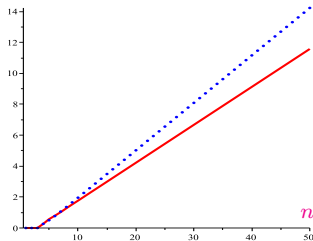
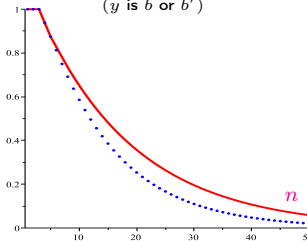
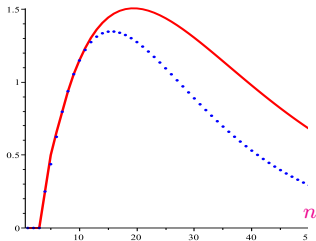
$$t = t_{\mathbf{A} \rightarrow \mathbf{C}} = t_{\mathbf{C} \rightarrow \mathbf{A}}$$

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

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

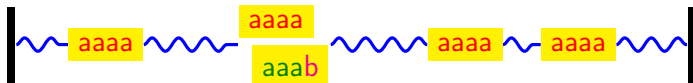
$$\mathfrak{P}_n \approx [z^n] \sum_{\alpha \neq \beta \in \{A, C, G, T\}} \frac{\partial F_b(z, 1, \dots, 1, \pi_{\alpha \rightarrow \beta} t_{\alpha \rightarrow \beta}, 1, \dots)}{\partial t_{\alpha \rightarrow \beta}} \Big|_{t_{\alpha \rightarrow \beta} = 1} / \widehat{f}_n^{(b)}$$

- ▶ 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**

II - Formal Languages Approach

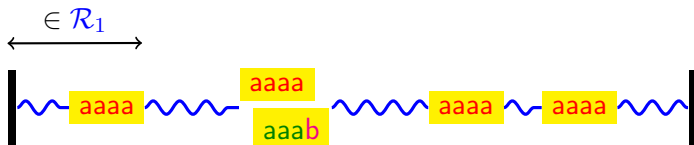
Guibas-Odlyzko decomposition - occurrences of a word u

$$U = (\text{aaaa}, \text{aaab}) \quad \begin{cases} u_1 = \text{aaaa} \\ u_2 = \text{aaab} \end{cases}$$



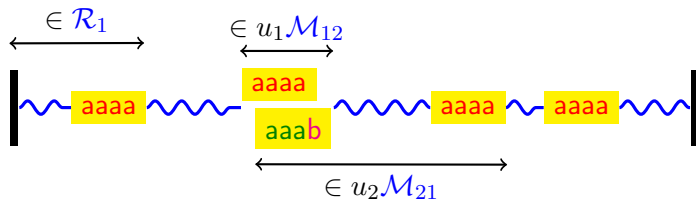
Guibas-Odlyzko decomposition - occurrences of a word u

$$U = (\text{aaaa}, \text{aaab}) \quad \begin{cases} u_1 = \text{aaaa} \\ u_2 = \text{aaab} \end{cases}$$



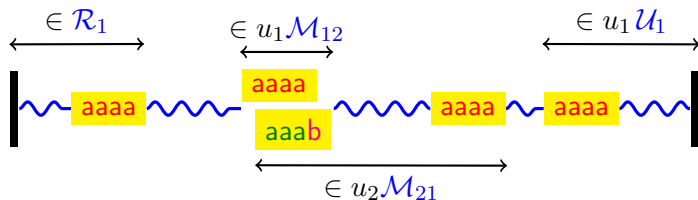
Guibas-Odlyzko decomposition - occurrences of a word u

$$U = (\text{aaaa}, \text{aaab}) \quad \begin{cases} u_1 = \text{aaaa} \\ u_2 = \text{aaab} \end{cases}$$



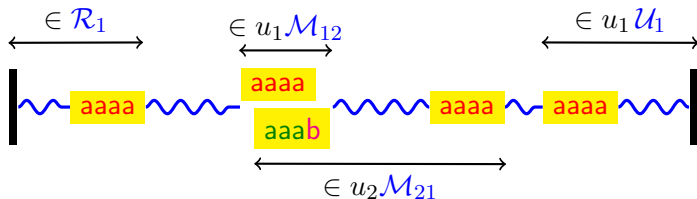
Guibas-Odlyzko decomposition - occurrences of a word u

$$U = (\text{aaaa}, \text{aaab}) \quad \begin{cases} u_1 = \text{aaaa} \\ u_2 = \text{aaab} \end{cases}$$



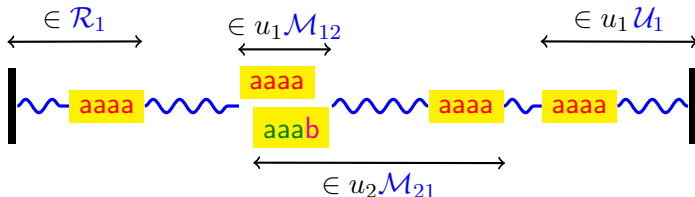
Guibas-Odlyzko decomposition - occurrences of a word u

$$U = (\text{aaaa}, \text{aaab}) \quad \begin{cases} u_1 = \text{aaaa} \\ u_2 = \text{aaab} \end{cases}$$



Guibas-Odlyzko decomposition - occurrences of a word u

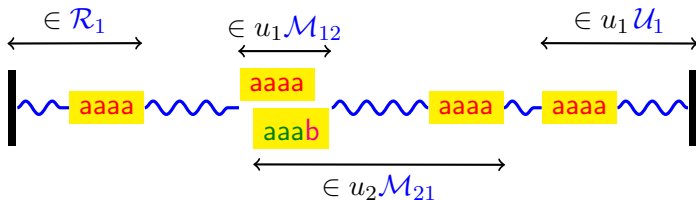
$$U = (\text{aaaa}, \text{aaab}) \quad \begin{cases} u_1 = \text{aaaa} \\ u_2 = \text{aaab} \end{cases}$$



- ▶ The “**Right**” language \mathcal{R}_i associated to the word u_i is the set of words $\mathcal{R}_i = \{r \mid r = e \cdot u_i \text{ and there is no } v \in U \text{ such that } r = xvy \text{ with } |y| > 0\}$.
- ▶ The “**Minimal**” language \mathcal{M}_{ij} leading from a word u_i to a word u_j is the set of words $\mathcal{M}_{ij} = \{m \mid u_i \cdot m = e \cdot u_j \text{ and there is no } v \in U \text{ such that } u_i \cdot m = xvy \text{ with } |x| > 0, |y| > 0\}$.
- ▶ The “**Ultimate**” language \mathcal{U}_i of words following the last occurrence of the word u_i (such that this occurrence is the last occurrence of U in the text) is the set of words $\mathcal{U}_i = \{u \mid \text{there is no } v \in U \text{ such that } u_i \cdot u = xvy \text{ with } |x| > 0\}$.
- ▶ 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

$$U = (\text{aaaa}, \text{aaab}) \quad \begin{cases} u_1 = \text{aaaa} \\ u_2 = \text{aaab} \end{cases}$$



$$F(z, x_1, x_2) = \mathcal{N}(z) + (\mathcal{R}_1(z)x_1, \mathcal{R}_2(z)x_2) \begin{pmatrix} \mathcal{M}_{11}(z)x_1 & \mathcal{M}_{12}(z)x_2 \\ \mathcal{M}_{21}(z)x_1 & \mathcal{M}_{22}(z)x_2 \end{pmatrix} \begin{pmatrix} \mathcal{U}_1(z) \\ \mathcal{U}_2(z) \end{pmatrix}$$

Computing the languages

- ▶ \mathcal{C}_{u_1, u_2} **correlation set** of two words u_1 and u_2
 $\mathcal{C}_{u_1, u_2} = \{ e \mid \exists e' \in \mathcal{A}^+, u_1 e = e' u_2 \text{ with } |e| < |u_2| \}.$
- ▶ $\mathcal{C}_u = \mathcal{C}_{u, u}$ **autocorrelation set** ($\epsilon \in \mathcal{C}_u$)

Computing the languages

- ▶ \mathcal{C}_{u_1, u_2} **correlation set** of two words u_1 and u_2
 $\mathcal{C}_{u_1, u_2} = \{ e \mid \exists e' \in \mathcal{A}^+, u_1 e = e' u_2 \text{ with } |e| < |u_2| \}.$
- ▶ $\mathcal{C}_u = \mathcal{C}_{u, u}$ **autocorrelation set** ($\epsilon \in \mathcal{C}_u$)

▶ Régnier-Szpankowski Equations

$$\bigcup_{k \geq 1} (\mathbb{M}^k)_{i,j} = \mathcal{A}^* \cdot u_j + \mathcal{C}_{ij} - \delta_{ij} \epsilon, \quad \mathcal{U}_i \cdot \mathcal{A} = \bigcup_j \mathcal{M}_{ij} + \mathcal{U}_i - \epsilon,$$

$$\mathcal{A} \cdot \mathcal{R}_j - (\mathcal{R}_j - u_j) = \bigcup_i u_i \mathcal{M}_{ij}, \quad \mathcal{N} \cdot u_j = \mathcal{R}_j + \bigcup_i \mathcal{R}_i (\mathcal{C}_{ij} - \delta_{ij} \epsilon),$$

▶ Automaton Computation

$$\mathcal{R}_i = \bigotimes_{1 \leq r \leq k} \overline{\mathcal{A}^* u_r \mathcal{A}^*} \cdot \mathcal{A} \bigotimes \mathcal{A}^* u_i$$

$$u_i \mathcal{M}_{ij} = u_i \mathcal{A}^* \bigotimes \mathcal{A}^* u_j \bigotimes_{1 \leq r \leq k} \overline{\mathcal{A} \mathcal{A}^* u_r \mathcal{A}^*} \mathcal{A}$$

$$u_j \mathcal{U}_j = u_j \mathcal{A}^* \bigotimes_{1 \leq r \leq k} \mathcal{A} \cdot \overline{\mathcal{A}^* u_r \mathcal{A}^*}$$

$$\mathcal{N} = \text{NOT} \left(\bigotimes_{1 \leq r \leq k} \overline{\mathcal{A}^* u_r \mathcal{A}^*} \right)$$

Constrained Guibas-Odlyzko languages

Example: $b = \text{AA}$, $d_\ell(b) = (\text{AC}, \text{CA})$

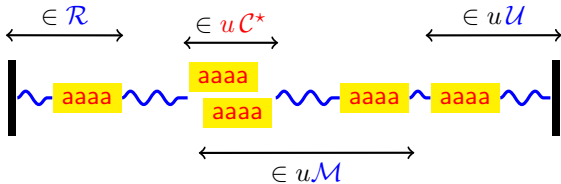
- ▶ We need **avoiding AA** in $S(0)$ and therefore **in the Right, Minimal and Ultimate languages**
- ▶ Build the **Régnier-Szpankowski languages** for the pattern $(\text{AC}, \text{CA}, \text{AA})$

$$\mathcal{L} = \mathcal{N} + (\mathcal{R}_1, \mathcal{R}_2, \mathcal{R}_3) \begin{pmatrix} \mathcal{M}_{11} & \mathcal{M}_{12} & \mathcal{M}_{13} \\ \mathcal{M}_{21} & \mathcal{M}_{22} & \mathcal{M}_{23} \\ \mathcal{M}_{31} & \mathcal{M}_{32} & \mathcal{M}_{33} \end{pmatrix} \begin{pmatrix} \mathcal{U}_1 \\ \mathcal{U}_2 \\ \mathcal{U}_3 \end{pmatrix}$$

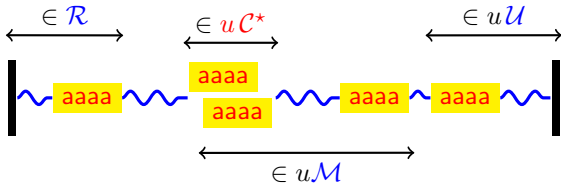
$$\widehat{\mathcal{L}} = \mathcal{N} + (\mathcal{R}_1, \mathcal{R}_2) \begin{pmatrix} \mathcal{M}_{11} & \mathcal{M}_{12} \\ \mathcal{M}_{21} & \mathcal{M}_{22} \end{pmatrix} \begin{pmatrix} \mathcal{U}_1 \\ \mathcal{U}_2 \end{pmatrix}$$

Notations: write $\widehat{\mathcal{N}}, \widehat{\mathcal{R}}_i, \widehat{\mathcal{M}}_{ij}, \widehat{\mathcal{U}}_j$ for **constrained languages**

Clump Analysis (Bassino-Clément-Fayolle-P.N. 2008)



Clump Analysis (Bassino-Clément-Fayolle-P.N. 2008)

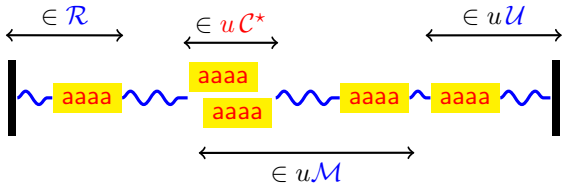


- ▶ *residual language* $\mathcal{D} = \mathcal{L}.u^-$: $\mathcal{D} = \{h, h \cdot u \in \mathcal{L}\}$
- ▶ $\mathcal{L}_2 - \mathcal{L}_1 = \mathcal{L}_2 \setminus \mathcal{L}_1 = \{h; h \in \mathcal{L}_2, h \notin \mathcal{L}_1\}$

Combinatorial decomposition (one word)

$$\mathcal{A}^* = \mathcal{N} + \mathcal{R}u^-u\mathcal{C}^*\left((\mathcal{M} - \mathcal{K})u^-u\mathcal{C}^*\right)^*\mathcal{U}$$

Clump Analysis (Bassino-Clément-Fayolle-P.N. 2008)



- ▶ *residual language* $\mathcal{D} = \mathcal{L}.u^-$: $\mathcal{D} = \{h, h \cdot u \in \mathcal{L}\}$
- ▶ $\mathcal{L}_2 - \mathcal{L}_1 = \mathcal{L}_2 \setminus \mathcal{L}_1 = \{h; h \in \mathcal{L}_2, h \notin \mathcal{L}_1\}$

Combinatorial decomposition (one word)

$$\begin{aligned}
 \mathcal{A}^* &= \mathcal{N} + \mathcal{R}u^-u\mathcal{C}^* \left((\mathcal{M} - \mathcal{K})u^-u\mathcal{C}^* \right)^* \mathcal{U} \\
 &= \mathcal{N} + \mathcal{R}u^-u\mathcal{K}^* \left((\mathcal{M} - \mathcal{K})u^-u\mathcal{K}^* \right)^* \mathcal{U} \\
 &= \mathcal{N} + \mathcal{R}u^-u\mathbf{S} \left((\mathcal{M} - \mathcal{K})u^-u\mathbf{S} \right)^* \mathcal{U}
 \end{aligned}$$

Clumps: $\mathbf{S} = u\mathcal{C}^* = u\mathcal{K}^*$

Some combinatorial properties

$$u = \textcolor{blue}{aaaaa}$$

$$\mathcal{C} - \{\epsilon\} = \{\textcolor{red}{a}, \textcolor{red}{aa}, \textcolor{red}{aaa}, \textcolor{red}{aaaa}\}$$

$$\mathcal{K} = \{\textcolor{red}{a}\}$$

$$\mathcal{M} = \{\textcolor{red}{a}, b(b + \textcolor{green}{ab} + \textcolor{green}{aab} + \textcolor{green}{aaab} + \textcolor{green}{aaaab})^* \textcolor{blue}{aaaaa}\}$$

Properties

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

Lemma.

Let $\mathcal{C}_o = \mathcal{C} - \{\epsilon\}$ be the strict autocorrelation set of a word u

- ▶ the Prefix code $\mathcal{K} = \mathcal{C}_o - \mathcal{C}_o \mathcal{A}^+$ generates **unambiguously** $\mathcal{C}^+ - \{\epsilon\}$, which implies that $\mathcal{K}^* = \mathcal{C}_o^*$

Clumps of reduced sets of words

Minimal Correlation Language: $\mathcal{K}_{ij} = (\mathcal{C}_{ij} - \mathcal{C}_{ij}\mathcal{A}^+) \cap \mathcal{M}_{ij}$

Lemma: $\mathcal{M}_{ij} - \mathcal{K}_{ij} = \mathcal{L}v_j$

Decomposition of a text by clumps:

$$\mathbb{K} = \begin{pmatrix} \mathcal{K}_{11} & \mathcal{K}_{12} \\ \mathcal{K}_{21} & \mathcal{K}_{22} \end{pmatrix}, \quad \mathbb{S} = \mathbb{K}^* \quad \mathbb{G} = \begin{pmatrix} v_1\mathbb{S}_{11} & v_1\mathbb{S}_{12} \\ v_2\mathbb{S}_{21} & v_2\mathbb{S}_{22} \end{pmatrix}$$

$$\mathcal{A}^* = \mathcal{N} + (\mathcal{R}_1v_1^{-1}, \mathcal{R}_2v_2^{-1})\mathbb{G}\left(\left((\mathcal{M}_{ij} - \mathcal{K}_{ij})v_j^{-1}\right)\mathbb{G}\right)^* \begin{pmatrix} \mathcal{U}_1 \\ \mathcal{U}_2 \end{pmatrix}$$

Constrained clumps

- ▶ **finite** code languages \mathcal{K}_{ij} **easy to compute**
- ▶ we must however **avoid** the **forbidden word** b while **extending clumps**
- ▶ $v_i, v_j \in d_\ell(b) \rightsquigarrow \hat{\mathcal{K}}_{ij} = \{h \in \mathcal{K}_{ij}; \quad |v_i.h|_b = 0\}$

sets \mathcal{K}_{ij} finite \implies computation of $\hat{\mathcal{K}}_{ij}$ by string-matching

Constrained clumps

- ▶ **finite** code languages \mathcal{K}_{ij} **easy to compute**
- ▶ we must however **avoid** the **forbidden word** b while **extending clumps**
- ▶ $v_i, v_j \in d_\ell(b) \rightsquigarrow \hat{\mathcal{K}}_{ij} = \{h \in \mathcal{K}_{ij}; \quad |v_i.h|_b = 0\}$

sets \mathcal{K}_{ij} finite \implies computation of $\hat{\mathcal{K}}_{ij}$ by string-matching

- ▶ **Decomposition by constrained clumps**

$$\hat{\mathcal{A}}_b^* = \hat{\mathcal{N}} + (\hat{\mathcal{R}}_1 v_1^-, \dots, \hat{\mathcal{R}}_r v_r^-) \hat{\mathcal{G}} \left((\hat{\mathcal{M}} - \hat{\mathcal{K}})^- \hat{\mathcal{G}} \right)^* \begin{pmatrix} \hat{u}_1 \\ \vdots \\ \hat{u}_r \end{pmatrix}$$

$$\text{with } \begin{cases} \hat{\mathcal{K}} = (\hat{\mathcal{K}}_{ij}), \\ \hat{\mathcal{S}} = \hat{\mathcal{K}}^*, \\ \hat{\mathcal{G}} = (v_i \hat{\mathcal{S}}_{ij}) \end{cases}$$

Generating function of the number of putative-hit positions

- ▶ $v_i(z, t) = \nu(v_i) t z^{|v_i|}$ for each $v_i \in d(b)$.
- ▶ for each $\hat{\mathcal{K}}_{ij}$, we can compute by string matching the number of putative-hit positions in each word of $v_i \cdot \hat{\mathcal{K}}_{ij}$.

$$\hat{\mathcal{K}}_{ij}(z, t) = \sum_{w \in \hat{\mathcal{K}}_{ij}} \nu(w) t^{\text{put-hit-pos}(v_i.w)-1} z^{|w|},$$

Generating function of the number of putative-hit positions

- ▶ $v_i(z, t) = \nu(v_i) t z^{|v_i|}$ for each $v_i \in d(b)$.
- ▶ for each $\hat{\mathcal{K}}_{ij}$, we can compute by string matching the number of putative-hit positions in each word of $v_i \cdot \hat{\mathcal{K}}_{ij}$.

$$\hat{\mathcal{K}}_{ij}(z, t) = \sum_{w \in \hat{\mathcal{K}}_{ij}} \nu(w) t^{\text{put-hit-pos}(v_i.w)-1} z^{|w|},$$

$$\hat{\mathbb{K}}(z, t) = \left(\hat{\mathcal{K}}_{ij}(z, t) \right), \quad \hat{\mathbb{S}}(z, t) = \left(\mathbb{I} - \hat{\mathbb{K}}(z, t) \right)^{-1},$$

$$\hat{\mathbb{G}}(z, t) = \left(v_i(z, t) \hat{\mathbb{S}}_{ij}(z, t) \right).$$

Generating function of the number of putative-hit positions

- ▶ $v_i(z, t) = \nu(v_i) t z^{|v_i|}$ for each $v_i \in d(b)$.
- ▶ for each $\hat{\mathcal{K}}_{ij}$, we can compute by string matching the number of putative-hit positions in each word of $v_i \cdot \hat{\mathcal{K}}_{ij}$.

$$\hat{\mathcal{K}}_{ij}(z, t) = \sum_{w \in \hat{\mathcal{K}}_{ij}} \nu(w) t^{\text{put-hit-pos}(v_i, w) - 1} z^{|w|},$$

$$\hat{\mathbb{K}}(z, t) = \left(\hat{\mathcal{K}}_{ij}(z, t) \right), \quad \hat{\mathbb{S}}(z, t) = \left(\mathbb{I} - \hat{\mathbb{K}}(z, t) \right)^{-1},$$

$$\hat{\mathbb{G}}(z, t) = \left(v_i(z, t) \hat{\mathbb{S}}_{ij}(z, t) \right).$$

$$F_b(z, t) = \hat{\mathcal{A}}_b^*(z, t)$$

$$= \hat{\mathcal{N}}(z) + (\hat{\mathcal{R}}_1 v_1^-(z), \dots, \hat{\mathcal{R}}_r v_r^-(z)) \hat{\mathbb{G}}(z, t) \left((\hat{\mathbb{M}} - \hat{\mathbb{K}})^-(z) \hat{\mathbb{G}}(z, t) \right)^* \begin{pmatrix} \hat{u}_1(z) \\ \vdots \\ \hat{u}_r(z) \end{pmatrix}$$

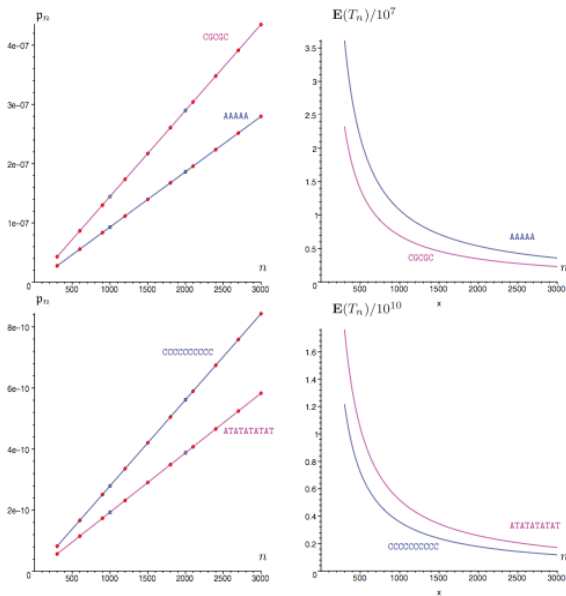


FIG. 5. Plots of the probability p_n (left) and of the expected waiting time $E(T_n)$ (right). (Top) $b = \text{AAAAA}$ (blue) and $b' = \text{CGCGC}$ (magenta). (Bottom) $b = \text{CCCCC CCCC}$ (blue) and $b' = \text{ATATATA TAT}$ (magenta). In the linear plots of the probability, the anchors values 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)