

Word and clump statistics and an application to DNA evolution

P. Nicodème

CNRS - Team Calin - LIPN - University Paris13

00100

00100
111100

00100
111100
1011111100

00100
111100
10111111100
01100

```
00100
111100
10111111100
01100
0100
```

```
00100
111100
10111111100
01100
0100
11100
```

00100
111100
10111111100
01100
0100
11100
11100

00100
111100
10111111100
01100
0100
11100
11100
010100

00100
111100
10111111100
01100
0100
11100
11100
010100
0110100

00100
111100
10111111100
01100
0100
11100
11100
010100
0110100
00100

```
00100
111100
10111111100
01100
0100
11100
11100
010100
0110100
00100
1011111101011111100
```

```
00100
111100
10111111100
01100
0100
11100
11100
010100
0110100
00100
1011111101011111100
00010100
```

```
00100
111100
10111111100
01100
0100
11100
11100
010100
0110100
00100
1011111101011111100
00010100
0101100
```

00100
111100
10111111100
01100
0100
11100
11100
010100
0110100
00100
1011111101011111100
00010100
0101100

100010100111

```
00100
111100
10111111100
01100
0100
11100
11100
010100
0110100
00100
1011111101011111100
00010100
0101100
```

```
100010100111
1100010100111
```



```
00100
111100
10111111100
01100
0100
11100
11100
010100
0110100
00100
1011111101011111100
00010100
0101100
```

```
100010100111
1100010100111
1000010111
```

```
00100
111100
10111111100
01100
0100
11100
11100
010100
0110100
00100
1011111101011111100
00010100
0101100
```

```
100010100111
1100010100111
1000010111
0011000011000000111
```

```
00100
111100
10111111100
01100
0100
11100
11100
010100
0110100
00100
1011111101011111100
00010100
0101100
```

```
100010100111
1100010100111
1000010111
0011000011000000111
010101010101000111
```

```
00100
111100
10111111100
01100
0100
11100
11100
010100
0110100
00100
1011111101011111100
00010100
0101100
```

```
100010100111
1100010100111
1000010111
0011000011000000111
010101010101000111
110110110011011010100101100110000011010111
```

```
00100
111100
10111111100
01100
0100
11100
11100
010100
0110100
00100
1011111101011111100
00010100
0101100
```

```
100010100111
1100010100111
1000010111
0011000011000000111
010101010101000111
110110110011011010100101100110000011010111
0111
```

```
00100
111100
10111111100
01100
0100
11100
11100
010100
0110100
00100
1011111101011111100
00010100
0101100
```

```
100010100111
1100010100111
1000010111
0011000011000000111
010101010101000111
110110110011011010100101100110000011010111
0111
111
```

```
00100
111100
10111111100
01100
0100
11100
11100
010100
0110100
00100
1011111101011111100
00010100
0101100
```

```
100010100111
1100010100111
1000010111
0011000011000000111
010101010101000111
110110110011011010100101100110000011010111
0111
111
001100010111
```

```
00100
111100
10111111100
01100
0100
11100
11100
010100
0110100
00100
1011111101011111100
00010100
0101100
```

```
100010100111
1100010100111
1000010111
0011000011000000111
010101010101000111
110110110011011010100101100110000011010111
0111
111
001100010111
0100010010010100000100111
```



```
00100
111100
10111111100
01100
0100
11100
11100
010100
0110100
00100
1011111101011111100
00010100
0101100
```

```
100010100111
1100010100111
1000010111
0011000011000000111
010101010101000111
110110110011011010100101100110000011010111
0111
111
001100010111
0100010010010100000100111
00111
```

```
00100
111100
10111111100
01100
0100
11100
11100
010100
0110100
00100
1011111101011111100
00010100
0101100
```

```
100010100111
1100010100111
1000010111
0011000011000000111
010101010101000111
110110110011011010100101100110000011010111
0111
111
001100010111
0100010010010100000100111
00111
000101100101010000111
```

```
00100
111100
10111111100
01100
0100
11100
11100
010100
0110100
00100
1011111101011111100
00010100
0101100
```

```
100010100111
1100010100111
1000010111
0011000011000000111
010101010101000111
110110110011011010100101100110000011010111
0111
111
001100010111
0100010010010100000100111
00111
000101100101010000111
10101010100010101100101010000010000110111
```

```
00100
111100
10111111100
01100
0100
11100
11100
010100
0110100
00100
1011111101011111100
00010100
0101100
```

```
100010100111
1100010100111
1000010111
0011000011000000111
010101010101000111
110110110011011010100101100110000011010111
0111
111
001100010111
0100010010010100000100111
00111
000101100101010000111
10101010100010101100101010000010000110111
```

1 1 1000011100

$\frac{1}{2}$ $\frac{1}{2}$ 1000011100
0001101000010110

1	1	1000011100
2	2	0001101000010110
3	3	1001110101

1	1	1000011100
2	2	0001101000010110
3	3	1001110101
4		111011110

1	1	1000011100
2	2	0001101000010110
3	3	1001110101
4		111011110
5		111011111

1	1	1000011100
2	2	0001101000010110
3	3	1001110101
4		111011110
5		111011111
6	4	001000101010

1	1	1000011100
2	2	0001101000010110
3	3	1001110101
4		111011110
5		111011111
6	4	001000101010
7	5	11001011011

1	1	1000011100
2	2	0001101000010110
3	3	1001110101
4		111011110
5		111011111
6	4	001000101010
7	5	11001011011
8	6	1001100010

1	1	1000011100
2	2	0001101000010110
3	3	1001110101
4		111011110
5		111011111
6	4	001000101010
7	5	11001011011
8	6	1001100010
9		10111010011

1	1	1000011100
2	2	0001101000010110
3	3	1001110101
4		111011110
5		111011111
6	4	001000101010
7	5	11001011011
8	6	1001100010
9		10111010011
10	7	0011000001010

1	1	1000011100
2	2	0001101000010110
3	3	1001110101
4		111011110
5		111011111
6	4	001000101010
7	5	11001011011
8	6	1001100010
9		10111010011
10	7	0011000001010
11		0111110010

1	1	1000011100
2	2	0001101000010110
3	3	1001110101
4		111011110
5		111011111
6	4	001000101010
7	5	11001011011
8	6	1001100010
9		10111010011
10	7	0011000001010
11		0111110010
12		11010111100010

1	1	1000011100
2	2	0001101000010110
3	3	1001110101
4		111011110
5		111011111
6	4	001000101010
7	5	11001011011
8	6	1001100010
9		10111010011
10	7	0011000001010
11		0111110010
12		11010111100010
13		111010011

1	1	1000011100
2	2	0001101000010110
3	3	1001110101
4		111011110
5		111011111
6	4	001000101010
7	5	11001011011
8	6	1001100010
9		10111010011
10	7	0011000001010
11		0111110010
12		11010111100010
13		111010011
14	8	101001011000

1	1	1000011100
2	2	0001101000010110
3	3	1001110101
4		111011110
5		111011111
6	4	001000101010
7	5	11001011011
8	6	1001100010
9		10111010011
10	7	0011000001010
11		0111110010
12		11010111100010
13		111010011
14	8	101001011000
15	9	001000011011

1	1	1000011100
2	2	0001101000010110
3	3	1001110101
4		111011110
5		111011111
6	4	001000101010
7	5	11001011011
8	6	1001100010
9		10111010011
10	7	0011000001010
11		0111110010
12		11010111100010
13		111010011
14	8	101001011000
15	9	001000011011
16	10	001000011110

1	1	1000011100
2	2	0001101000010110
3	3	1001110101
4		111011110
5		111011111
6	4	001000101010
7	5	11001011011
8	6	1001100010
9		10111010011
10	7	0011000001010
11		0111110010
12		11010111100010
13		111010011
14	8	101001011000
15	9	001000011011
16	10	001000011110
17	11	000011000111111

1	1	1000011100
2	2	0001101000010110
3	3	1001110101
4		111011110
5		111011111
6	4	001000101010
7	5	11001011011
8	6	1001100010
9		10111010011
10	7	0011000001010
11		0111110010
12		11010111100010
13		111010011
14	8	101001011000
15	9	001000011011
16	10	001000011110
17	11	000011000111111
18	12	1001010011

1	1	1000011100
2	2	0001101000010110
3	3	1001110101
4		111011110
5		111011111
6	4	001000101010
7	5	11001011011
8	6	1001100010
9		10111010011
10	7	0011000001010
11		0111110010
12		11010111100010
13		111010011
14	8	101001011000
15	9	001000011011
16	10	001000011110
17	11	000011000111111
18	12	1001010011
19	13	00010101001110001

1	1	1000011100
2	2	0001101000010110
3	3	1001110101
4		111011110
5		111011111
6	4	001000101010
7	5	11001011011
8	6	1001100010
9		10111010011
10	7	0011000001010
11		0111110010
12		11010111100010
13		111010011
14	8	101001011000
15	9	001000011011
16	10	001000011110
17	11	000011000111111
18	12	1001010011
19	13	00010101001110001
20		111101011

1	1	1000011100
2	2	0001101000010110
3	3	1001110101
4		111011110
5		111011111
6	4	001000101010
7	5	11001011011
8	6	1001100010
9		10111010011
10	7	0011000001010
11		0111110010
12		11010111100010
13		111010011
14	8	101001011000
15	9	001000011011
16	10	001000011110
17	11	000011000111111
18	12	1001010011
19	13	00010101001110001
20		111101011
21	14	10101001001000

1	1	1000011100
2	2	0001101000010110
3	3	1001110101
4		111011110
5		111011111
6	4	001000101010
7	5	11001011011
8	6	1001100010
9		10111010011
10	7	0011000001010
11		0111110010
12		11010111100010
13		111010011
14	8	101001011000
15	9	001000011011
16	10	001000011110
17	11	000011000111111
18	12	1001010011
19	13	00010101001110001
20		111101011
21	14	10101001001000
22	15	101001111100

1	1	1000011100
2	2	0001101000010110
3	3	1001110101
4		111011110
5		111011111
6	4	001000101010
7	5	11001011011
8	6	1001100010
9		10111010011
10	7	0011000001010
11		0111110010
12		11010111100010
13		111010011
14	8	101001011000
15	9	001000011011
16	10	001000011110
17	11	000011000111111
18	12	1001010011
19	13	00010101001110001
20		111101011
21	14	10101001001000
22	15	101001111100
23		0111010111

1	1	1000011100
2	2	0001101000010110
3	3	1001110101
4		111011110
5		111011111
6	4	001000101010
7	5	11001011011
8	6	1001100010
9		10111010011
10	7	0011000001010
11		0111110010
12		11010111100010
13		111010011
14	8	101001011000
15	9	001000011011
16	10	001000011110
17	11	000011000111111
18	12	1001010011
19	13	00010101001110001
20		111101011
21	14	10101001001000
22	15	101001111100
23		0111010111
24	16	0001001010110

1	1	1000011100
2	2	0001101000010110
3	3	1001110101
4		111011110
5		111011111
6	4	001000101010
7	5	11001011011
8	6	1001100010
9		10111010011
10	7	0011000001010
11		0111110010
12		11010111100010
13		111010011
14	8	101001011000
15	9	001000011011
16	10	001000011110
17	11	000011000111111
18	12	1001010011
19	13	00010101001110001
20		111101011
21	14	10101001001000
22	15	101001111100
23		0111010111
24	16	0001001010110
25		111110111

1	1	1000011100
2	2	0001101000010110
3	3	1001110101
4		111011110
5		111011111
6	4	001000101010
7	5	11001011011
8	6	1001100010
9		10111010011
10	7	0011000001010
11		0111110010
12		11010111100010
13		111010011
14	8	101001011000
15	9	001000011011
16	10	001000011110
17	11	000011000111111
18	12	1001010011
19	13	00010101001110001
20		111101011
21	14	10101001001000
22	15	101001111100
23		0111010111
24	16	0001001010110
25		111110111
26		0110111000011

1	1	1000011100
2	2	0001101000010110
3	3	1001110101
4		111011110
5		111011111
6	4	001000101010
7	5	11001011011
8	6	1001100010
9		10111010011
10	7	0011000001010
11		0111110010
12		11010111100010
13		111010011
14	8	101001011000
15	9	001000011011
16	10	001000011110
17	11	000011000111111
18	12	1001010011
19	13	00010101001110001
20		111101011
21	14	10101001001000
22	15	101001111100
23		0111010111
24	16	0001001010110
25		111110111
26		0110111000011
27	17	1000000000

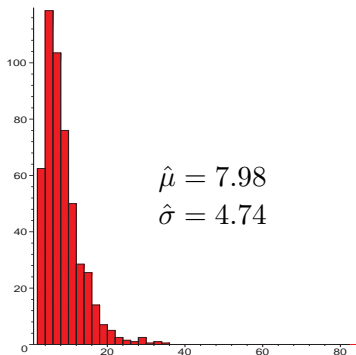
1	1	1000011100
2	2	0001101000010110
3	3	1001110101
4		111011110
5		111011111
6	4	001000101010
7	5	11001011011
8	6	1001100010
9		10111010011
10	7	0011000001010
11		0111110010
12		11010111100010
13		111010011
14	8	101001011000
15	9	001000011011
16	10	001000011110
17	11	000011000111111
18	12	1001010011
19	13	00010101001110001
20		111101011
21	14	10101001001000
22	15	101001111100
23		0111010111
24	16	0001001010110
25		111110111
26		0110111000011
27	17	100000000
28		00111111011

1	1	1000011100
2	2	0001101000010110
3	3	1001110101
4		111011110
5		111011111
6	4	001000101010
7	5	11001011011
8	6	1001100010
9		10111010011
10	7	0011000001010
11		0111110010
12		11010111100010
13		111010011
14	8	101001011000
15	9	001000011011
16	10	001000011110
17	11	000011000111111
18	12	1001010011
19	13	000101010011110001
20		111101011
21	14	10101001001000
22	15	101001111100
23		0111010111
24	16	0001001010110
25		111110111
26		0110111000011
27	17	1000000000
28		00111111011
29	18	11000011001

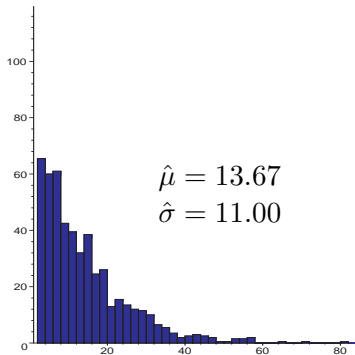
1	1	1000011100
2	2	0001101000010110
3	3	1001110101
4		111011110
5		111011111
6	4	001000101010
7	5	11001011011
8	6	1001100010
9		10111010011
10	7	0011000001010
11		0111110010
12		11010111100010
13		111010011
14	8	101001011000
15	9	001000011011
16	10	001000011110
17	11	000011000111111
18	12	1001010011
19	13	000101010011110001
20		111101011
21	14	10101001001000
22	15	101001111100
23		0111010111
24	16	0001001010110
25		111110111
26		0110111000011
27	17	1000000000
28		00111111011
29	18	11000011001
30		0111001000

1	1	1000011100
2	2	0001101000010110
3	3	1001110101
4		111011110
5		111011111
6	4	001000101010
7	5	11001011011
8	6	1001100010
9		10111010011
10	7	0011000001010
11		0111110010
12		11010111100010
13		111010011
14	8	101001011000
15	9	001000011011
16	10	001000011110
17	11	000011000111111
18	12	1001010011
19	13	000101010011110001
20		111101011
21	14	10101001001000
22	15	101001111100
23		0111010111
24	16	0001001010110
25		111110111
26		0110111000011
27	17	1000000000
28		00111111011
29	18	11000011001
30		0111001000

Waiting time (1000 random texts)

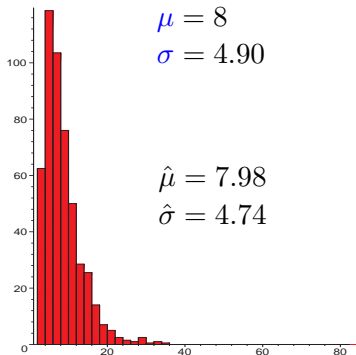


100

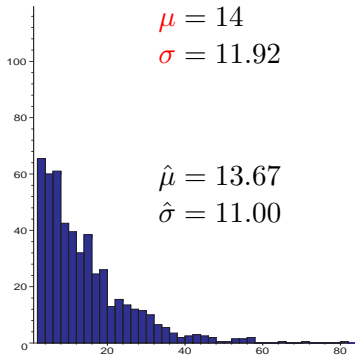


111

Waiting time (1000 random texts)



100



111

What is going on?

- ▶ Probability of appearance at a given position

$$\mathbf{P}(100) = \frac{1}{8}$$

$$\mathbf{P}(111) = \frac{1}{8}$$

What is going on?

- ▶ Probability of appearance at a given position

$$P(100) = \frac{1}{8}$$

$$P(111) = \frac{1}{8}$$

- ▶ BUT the 111 occur often by CLUMPS

...0111110

111

111...

...011110

111...

- ▶ while the 100 NEVER OVERLAP

What is going on?

- ▶ Probability of appearance at a given position

$$P(100) = \frac{1}{8}$$

$$P(111) = \frac{1}{8}$$

- ▶ BUT the 111 occur often by CLUMPS

...0111110

111

111...

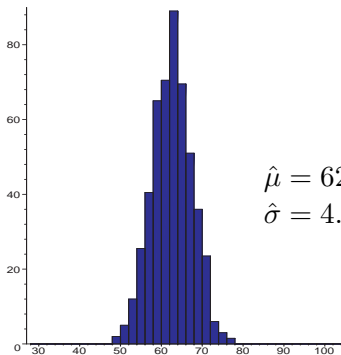
...011110

111...

- ▶ while the 100 NEVER OVERLAP

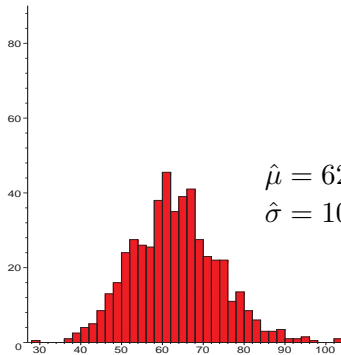
Expected waiting time: 111 – 14 100 – 7

Number of occurrences 1000 texts of size 500



$$\hat{\mu} = 62.38$$
$$\hat{\sigma} = 4.90$$

100



$$\hat{\mu} = 62.07$$
$$\hat{\sigma} = 10.79$$

111

Part I

Statistics of reduced patterns

1 word - Bernoulli model

\mathcal{A} alphabet, w considered word

Polynomial of autocorrelation

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

$w = ababa$

$ababa$

$ababa|$

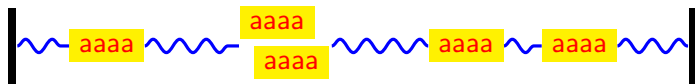
$ababa$

$ababa$

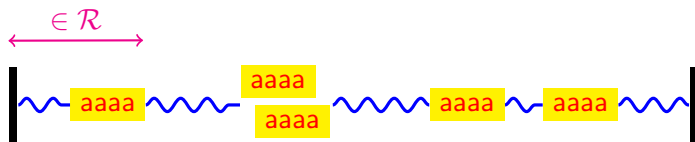
$$\mathcal{C}_{ababa} = \{\epsilon, ba, baba\}$$

$$C_{ababa}(z) = \sum_{v \in \mathcal{C}_{ababa}} \mathbf{P}(v)z^{|v|} = 1 + \omega_a\omega_b z^2 + \omega_a^2\omega_b^2 z^4$$

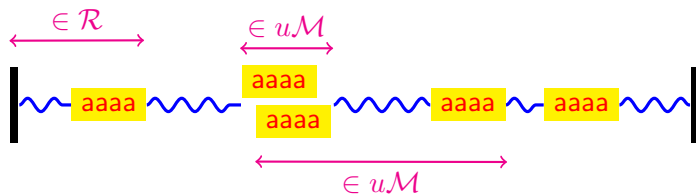
Guibas-Odlyzko decomposition for a word u



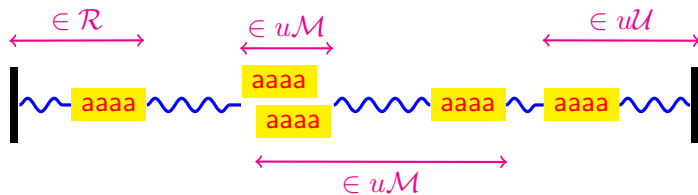
Guibas-Odlyzko decomposition for a word u



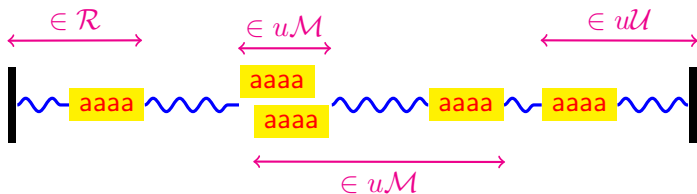
Guibas-Odlyzko decomposition for a word u



Guibas-Odlyzko decomposition for a word u



Guibas-Odlyzko decomposition for a word u



$$A^* = \mathcal{N} \cup \mathcal{R} \cdot \mathcal{M}^* \cdot \mathcal{U} \quad \mathcal{L}_\bullet = \mathcal{N} \cup \mathcal{R} \bullet \cdot (\mathcal{M} \bullet)^* \cdot \mathcal{U}$$

need to compute $\mathcal{N}, \mathcal{R}, \mathcal{M}, \mathcal{U}$

The languages \mathcal{R} , \mathcal{M} and \mathcal{U}

First $\mathcal{R} = \{ t = u.w \text{ and } \nexists r, s, t = r.w.s \}$

$aaaaaababa \in \mathcal{R}$, $bbbbbabababa \notin \mathcal{R}$

The languages \mathcal{R} , \mathcal{M} and \mathcal{U}

First $\mathcal{R} = \{ t = u.w \text{ and } \nexists r, s, t = r.w.s \}$

$aaaaaababa \in \mathcal{R}$, $bbbbbabababa \notin \mathcal{R}$

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

$ababa$
 $aaaaababa \in \mathcal{M}$ $ababa$
 $babbbbbbbbababa \notin \mathcal{M}$ $ababa$
 $ba \in \mathcal{M}$

The languages \mathcal{R} , \mathcal{M} and \mathcal{U}

First $\mathcal{R} = \{ t = u.w \text{ and } \exists r, s, t = r.w.s \}$

$aaaaaababa \in \mathcal{R}$, $bbbbbabababa \notin \mathcal{R}$

Minimal $\mathcal{M} = \{ t, w.t = u.w \text{ and } \exists r, s, w.t = r.w.s \}$

$ababa aaaaaababa \in \mathcal{M}$ $ababa bbbbbbbbababa \notin \mathcal{M}$ $ababa ba \in \mathcal{M}$

Ultimate $\mathcal{U} = \{ t, \exists r, s, w.t = r.w.s \}$

$ababa aabbbabbbbbbb \in \mathcal{U}$ $ababa bbbbbbbbbbbbb \notin \mathcal{U}$

Not $\mathcal{N} = \overline{\mathcal{A}^*.w.\mathcal{A}^*} = \{ t, \exists r, s, t = r.w.s \}$

Régnier-Szpankowski - Equations over the languages

- ▶ (I) $\mathcal{A}^* = \mathcal{U} + \mathcal{M}\mathcal{A}^* \iff w.\mathcal{A}^* = w.\mathcal{U} + w.\mathcal{M}\mathcal{A}^*$
– for any word beginning with w
1. 1 single occurrence of $w \Rightarrow w.\mathcal{U}$
 2. several occurrences of $w \Rightarrow w.\mathcal{M}.\mathcal{A}^*$

Régnier-Szpankowski - Equations over the languages

- ▶ (I) $\mathcal{A}^* = \mathcal{U} + \mathcal{M}\mathcal{A}^*$ $\Leftrightarrow w.\mathcal{A}^* = w.\mathcal{U} + w.\mathcal{M}\mathcal{A}^*$
 - for any word beginning with w
 - 1. 1 single occurrence of $w \Rightarrow w.\mathcal{U}$
 - 2. several occurrences of $w \Rightarrow w.\mathcal{M}.\mathcal{A}^*$

- ▶ (II) $\mathcal{A}^*w = \mathcal{R}.\mathcal{C} + \mathcal{R}.\mathcal{A}^*.w$ (remark $\epsilon \in \mathcal{C}$)
 - for any word finishing by w
 - 1. first occurrence w overlaps last one
or single occurrence of $w \Rightarrow \mathcal{R}.\mathcal{C}$
 - 2. first occurrence of w does not overlap last $\Rightarrow \mathcal{R}.\mathcal{A}^*.w$

Equations over the languages (continued)

- (III) $\left\{ \begin{array}{l} \mathcal{M}^+ = \mathcal{A}^*.w + \mathcal{C} - \epsilon \\ \Leftrightarrow w.\mathcal{M}^+ = w.\mathcal{A}^*.w + w.(\mathcal{C} - \epsilon) \end{array} \right.$
1. $w.\mathcal{M}^+$, \Rightarrow words beginning and finishing with w
 2. $w.\mathcal{A}^*.w$, \Rightarrow first and last occurrences do not overlap
 3. $w.(\mathcal{C} - \epsilon)$ \Rightarrow first and last occurrences overlap

Equations over the languages (continued)

- ▶ (III) $\left\{ \begin{array}{l} \mathcal{M}^+ = \mathcal{A}^*.w + \mathcal{C} - \epsilon \\ \Leftrightarrow w.\mathcal{M}^+ = w.\mathcal{A}^*.w + w.(\mathcal{C} - \epsilon) \end{array} \right.$
1. $w.\mathcal{M}^+$, \Rightarrow words beginning and finishing with w
 2. $w.\mathcal{A}^*.w$, \Rightarrow first and last occurrences do not overlap
 3. $w.(\mathcal{C} - \epsilon)$ \Rightarrow first and last occurrences overlap
- ▶ (IV) $\mathcal{N}.\mathcal{A} = \mathcal{R} + \mathcal{N} - \epsilon$
- concatenate a letter to a word of $N \Rightarrow$
1. creates a match $\Rightarrow \mathcal{R}$
 2. does not create a match $\Rightarrow \mathcal{N}$
 3. (empty word ϵ forbidden)

From Languages to Generating Functions

- (I) $\mathcal{A}^* = \mathcal{U} + \mathcal{M}\mathcal{A}^*$ $\rightsquigarrow \frac{1}{1-z} = U(z) + \frac{M(z)}{1-z}$
- (II) $\mathcal{A}^*w = \mathcal{R}\mathcal{A} + \mathcal{R}\mathcal{A}^*.w$ $\rightsquigarrow \frac{\omega_w z^{|w|}}{1-z} = F(z) \left(C(z) + \frac{\omega_w z^{|w|}}{1-z} \right)$
- (III) $\mathcal{M}^+ = \mathcal{A}^*.w + \mathcal{A} - \epsilon$ $\rightsquigarrow \frac{M(z)}{1-M(z)} = \frac{\omega_w z^{|w|}}{1-z} + C(z) - 1$
- (IV) $\mathcal{N}\mathcal{A} = \mathcal{R} + \mathcal{N} - \epsilon$ $\rightsquigarrow zN(z) = F(z) + N(z) - 1$

Solving the system

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

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

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

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

$$\mathcal{L}_\bullet = \mathcal{N} \cup \mathcal{R}_\bullet \cdot (\mathcal{M}_\bullet)^* \cdot \mathcal{U}$$

$$L(z, u) = N(z) + \frac{R(z)uU(z)}{1 - uM(z)} = \frac{1}{1 - z - \frac{(u-1)\omega_w z^{|w|}}{1 - (u-1)(C(z) - 1)}}$$

Moments

X_n random variable counting the number of occurrences of the word w in a **random text** of length \mathbf{n}

$$\blacktriangleright \mathbf{E}(X_n) = [z^n] \frac{\partial L(z, u)}{\partial u} \Big|_{u=1} = (\mathbf{n} - |w| + 1)\omega_w$$

$$\begin{aligned} \blacktriangleright \mathbf{Var}(X_n) &= [z^n] \frac{\partial}{\partial u} u \frac{\partial L(z, u)}{\partial u} \Big|_{u=1} - \mathbf{E}^2(X_n) \\ &= \mathbf{n} \times \omega_w (2\mathcal{C}(1) - 1 - (2|w| - 1)\omega_w) + O(1) \end{aligned}$$

Number of occurrences

	100	111
$L(z, u)$	$\frac{1}{1 - z + (1 - u)\frac{z^3}{8}}$	$\frac{1 + \frac{1 - u}{8}(4z + 2z^2)}{1 - z - \frac{1 - u}{8}(4z - 2z^2 - z^3)}$
$\mu(z)$	$\frac{z^3}{8(1 - z)^2}$	$\frac{z^3}{8(1 - z)^2}$
$m_{(2)}(z)$	$\frac{z^3}{8(1 - z)^2} - \frac{z^3}{32(1 - z)^3}$	$\frac{8z^3 - 4z^5 - 2z^6}{32(1 - z)^3}$
μ_n	$(n - 2)/8$	$(n - 2)/8$
σ_n	$\sqrt{3n}/8$	$\sqrt{15n - 40}/8$
μ_{500}	$249/4 = 62.25$	$249/4 = 62.25$
σ_{500}	$\sqrt{1500}/8 \approx 4.84$	$\sqrt{7460}/8 \approx 10.80$

Reduced compound patterns

$W = \{w_1, w_2\}$ and w_1 (resp. w_2) is **not factor** of w_2 (resp. w_1)

Correlation of words

$$\mathbb{C}(z) = (\mathcal{C}_{i,j}(z)) \quad \mathcal{C}_{i,j} = \{h, w_i.h = u.w_j\}$$

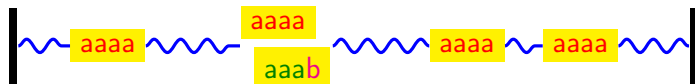
Example $W = \{aab, abaa\}$

$$\mathbb{C}(z) = \begin{pmatrix} 1 & \omega_a^2 z^2 \\ \omega_b z & 1 + \omega_a^2 \omega_b z^3 \end{pmatrix}$$

Languages **Right**, **Minimal**, **Ultimate** $\mathcal{R}_i, \mathcal{M}_{i,j}, \mathcal{U}_i$

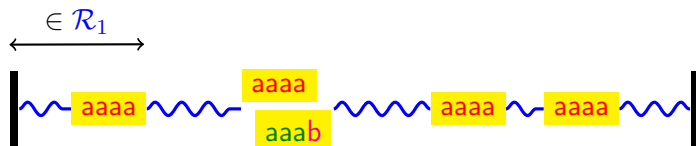
Guibas-Odlyzko decomposition for a pattern (u_1, u_2)

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



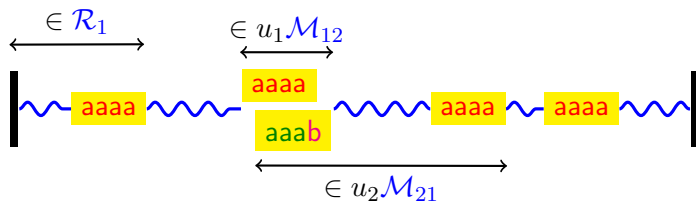
Guibas-Odlyzko decomposition for a pattern (u_1, u_2)

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



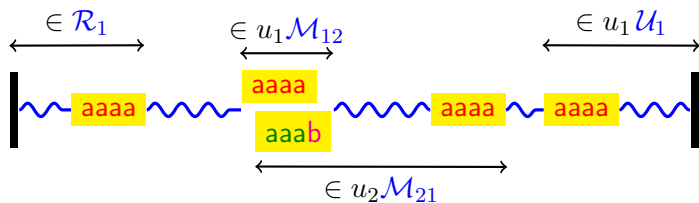
Guibas-Odlyzko decomposition for a pattern (u_1, u_2)

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



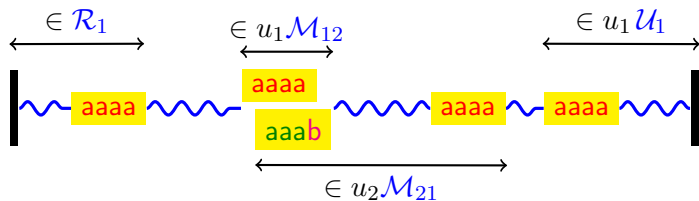
Guibas-Odlyzko decomposition for a pattern (u_1, u_2)

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



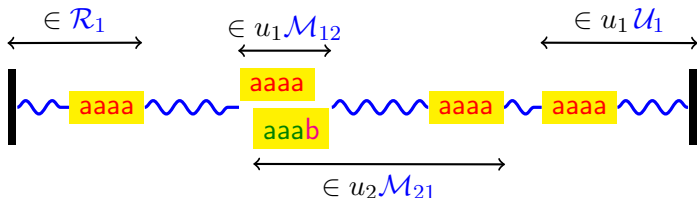
Guibas-Odlyzko decomposition for a pattern (u_1, u_2)

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



Guibas-Odlyzko decomposition for a pattern (u_1, u_2)

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



- ▶ 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 , $\mathcal{N} = \{n \mid \text{there is no } v \in U \text{ such that } n = xvy\}$.

Computing the languages

► Régnier-Szpankowski Equations

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

$$L(z, \mathbf{x}_1, \mathbf{x}_2)$$

$$= N(z) + (R_1(z)\mathbf{x}_1, R_2(z)\mathbf{x}_2) \begin{pmatrix} M_{11}(z)\mathbf{x}_1 & M_{12}(z)\mathbf{x}_2 \\ M_{21}(z)\mathbf{x}_1 & M_{22}(z)\mathbf{x}_2 \end{pmatrix} \begin{pmatrix} U_1(z) \\ U_2(z) \end{pmatrix}$$

► The generating functions are also computable by **automata**

Markov case - One word w

- ▶ **Conditional probability matrix:**

$$\mathbb{P} = (p_{ij}), \quad i, j \in \{1, \dots, |\mathcal{A}|\}$$

- ▶ $w = aaba$

- ▶ $\mathcal{C} = \{baa, abaa\}$

- ▶ $C_w(z) = p_{ab}p_{ba}p_{aa}^2z^3 + p_{aa}p_{ab}p_{aa}^2z^4$

- ▶ $C_w(z)$ is the **conditioned** generating function of the autocorrelation set

Markov case - One word $w = w_1.w_2 \dots .w_m$

Conditional probability matrix: $\mathbb{P} = (p_{ij})$, $i, j \in \{1, \dots, |\mathcal{A}|\}$

Stationary vector: $\mu = (\mu_1, \dots, \mu_{|\mathcal{A}|})$

Stationary matrix: $\mathbb{\Pi} = \begin{pmatrix} \mu \\ \vdots \\ \mu \end{pmatrix} = \lim_{n \rightarrow \infty} \mathbb{P}^n$

$$R(z) = \frac{P(w)z^m}{D}$$

$$M(z) = 1 + \frac{z-1}{D}$$

$$U(z) = \frac{1}{D}$$

$$N(z) = \frac{D - P(w)z^m}{(1-z)D}$$

$$D = (1-z)C_w(z) + P(w)z^m(1 - (1-z)F(z))$$

$$F(z) = \frac{1}{\mu_{w_1}} \left[(\mathbb{P} - \mathbb{\Pi}) \left(\mathbb{I} - (\mathbb{P} - \mathbb{\Pi})z \right)^{-1} \right]_{w_m, w_1}$$

Limit Laws

- ▶ **Poisson** when number of occurrences is $O(1)$
- ▶ **Gaussian** when number of occurrences is $\Theta(n)$

Also **large deviation result**

Bibliography

- ▶ **Lothaire** book (2005), Chapter 7
- ▶ **Régnier** (2000) "A unified approach to word occurrences probabilities"

Part II

Clump Analysis

Counting clumps

$w = aa$ $T = bbbb$ *aaaa* $bbbb$ *aaa*

- ▶ word occurrences counting with overlaps: 5 matches
(*aa|a|a|*, *aa|a|*)
- ▶ clumps counting: 2 matches (*aaaa|*, *aaa|*)

Counting clumps

$w = aa$ $T = bbbb\color{red}aaa\color{red}bbbb\color{red}aaa$

- ▶ word occurrences counting with overlaps: 5 matches
($aa|a|a|$, $aa|a|$)
- ▶ clumps counting: 2 matches ($\color{red}aaaa|$, $\color{red}aaa|$)

A clump of occurrences of a word w is

- ▶ either composed of a **single occurrence** of w
- ▶ or of a **maximal set** of occurrences of w such that **each occurrence of w overlaps at least another occurrence** of w

Generalization to clumps of occurrences of a set of words

$\{u_1, \dots, u_r\}$

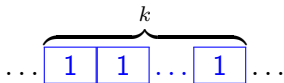
Aim of our work

- ▶ a **combinatorial** analysis for counting clumps of reduced sets of words
- ▶ an **algorithmic** construction by **automata** that solves the counting problem in the general case and implies a normal limit law

Probability of start at a position i

(A) word $w = 1^k$

$$p = \mathbf{P}(1) = 1 - \mathbf{P}(0) = 1 - q$$

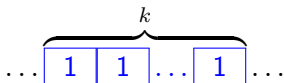


$$\mathbf{P}(\text{start}) = p^k$$

Probability of start at a position i

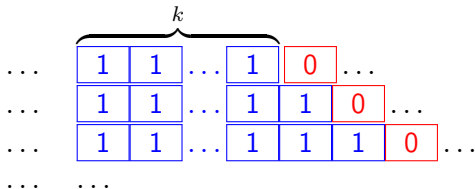
(A) word $w = 1^k$

$$p = \mathbf{P}(1) = 1 - \mathbf{P}(0) = 1 - q$$



$$\mathbf{P}(\text{start}) = p^k$$

(B) clump $\Gamma = 1^k.1^*$

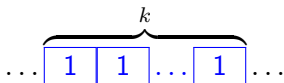


$$\mathbf{P}(\text{start}) = p^k \times q \times (1 + p + p^2 + \dots) = p^k \times \frac{q}{1 - p} = p^k$$

FALSE

Probability of start at a position i

(A) word $w = 1^k$



$$\mathbf{P}(\text{start}) = p^k$$

(B) clump $\Gamma = 1^k.1^*$

$$\dots 1 \mid \overbrace{1 \dots 1}^k \dots$$

no clump **beginning** at position i

$$\dots 0 \mid \overbrace{1 \dots 1}^k \dots$$

a clump **begins** at position i

$$\mathbf{P}(\text{start}) = p^k \times q = p^k \times (1 - p)$$

Probabilistic approach

Prum, Reinert, Schbath, Pape, ...

$w = aaa$

$\mathbf{P}(a)$ small \rightsquigarrow $\mathbf{P}(\text{start-of-a-clump-at-a-position})$ small

\rightsquigarrow

Poisson law for the number of clumps

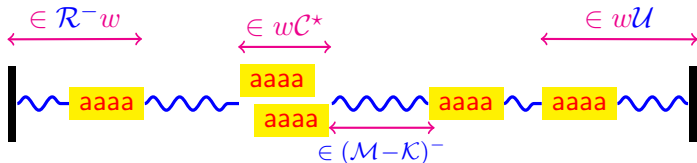
Geometric law for the number of words in a clump

This extends to the general case.

Chen-Stein Poisson approximation provides bounds for the total variation distance to the composed law.

A - Combinatorial approach

Language equations for the clumps



Notation: if $\mathcal{L} = \mathcal{W} \cdot w$ we write $\mathcal{L}^- = \mathcal{W}$

Combinatorial decomposition

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

Clump: $\mathbf{\Gamma} = w \mathcal{C}^* = w \mathcal{K}^*$

Some combinatorial properties

$$w = aaaaaa$$

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

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

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

Properties

- ▶ $\mathcal{K} \subset \mathcal{M}$
- ▶ $\mathcal{M} - \mathcal{K} = \mathcal{L}w$
- ▶ $\mathcal{K}^* = \mathcal{C}^*$ and \mathcal{K}^* is **unambiguous**

A Prefix Code generating the clumps

Lemma.

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

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

Generating functions

$$F(z, \bullet, \bullet, \bullet) = \mathcal{N}(z) + \frac{\mathcal{R}(z)}{\omega_w z^{|w|}} \Gamma(\bullet z, \bullet, \bullet) \frac{1}{1 - \frac{\mathcal{M}(z) - \mathcal{K}(z)}{\omega_w z^{|w|}} \times \Gamma(\bullet z, \bullet, \bullet)} \mathcal{U}(z)$$

- ▶ number of w and number of **clumps**:

$$\Gamma(z, u, x) = ux \omega_w z^{|w|} \frac{1}{1 - x \mathcal{K}(z)}$$

- ▶ number of **clumps** and **total** number of positions inside clumps:

$$\Gamma(tz, u) = u \omega_w (tz)^{|w|} \frac{1}{1 - \mathcal{K}(tz)}$$

- ▶ number of w and **total** number of positions inside clumps:

$$\Gamma(tz, x) = x \omega_w (tz)^{|w|} \frac{1}{1 - x \mathcal{K}(tz)}$$

- ▶ number of “stuttering” w and **total** number of positions inside clumps:

$$\Gamma(tz, x) = \frac{x}{1-x} \omega_w (tz)^{|w|} \frac{1}{1 - \frac{x}{1-x} \mathcal{K}(tz)}$$

One word - Expectation - Variance

clumps

$$\mathbf{E}(O_n^{\hat{\mathbf{K}}}) = (\mathbf{n} - |w| + 1)\omega_w(1 - \mathcal{K}(1)) - \omega_w\mathcal{K}'(1)$$

$$\mathbf{Var}(O_n^{\hat{\mathbf{K}}}) = \mathbf{n} \times (1 - \mathcal{K}(1))^2 \mathbf{V}_w - \mathbf{n} \times \omega_w(1 - \mathcal{K}(1))(\mathcal{K}(1) - 2\omega_w\mathcal{K}'(1))$$

one word

$$\mathbf{E}(O_n^w) = (\mathbf{n} - |w| + 1)\omega_w, \quad \mathbf{Var}(O_n^w) = \mathbf{n} \times \mathbf{V}_w + O(1).$$

$$\mathbf{V}_w = \omega_w(2\mathcal{C}(1) - 1 - (2|w| - 1)\omega_w)$$

Putting up equations for clumps of two words

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

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

$$\mathbb{K} = \begin{pmatrix} \mathcal{K}_{11} & \mathcal{K}_{12} \\ \mathcal{K}_{21} & \mathcal{K}_{22} \end{pmatrix}$$

$$\mathbb{E} = \mathbb{K}^* \quad \mathbb{G} = \begin{pmatrix} w_1 \mathbb{E}_{11} & w_1 \mathbb{E}_{12} \\ w_2 \mathbb{E}_{21} & w_2 \mathbb{E}_{22} \end{pmatrix}$$

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

B - Automaton approach

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

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

$$\mathcal{E}_{aabaa, aabaa} = \{baa, abaa\} \qquad \mathcal{E}_{baab, baab} = \{aab\}$$

$$\mathcal{E}_{aabaa, baab} = \{b\} \qquad \mathcal{E}_{baab, aabaa} = \{aa\}$$

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

\mathcal{E}_{w_1, w_2} extension 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

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

\mathcal{E}_{w_1, w_2} extension 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} extension 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} extension 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

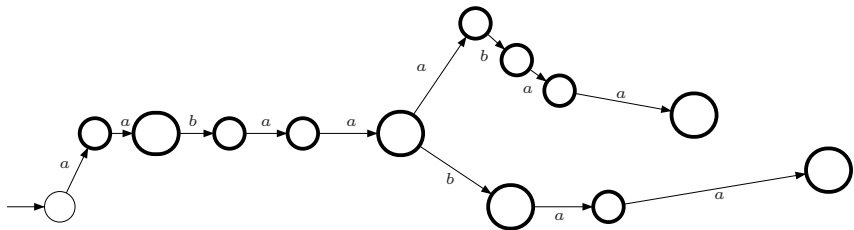
1. Build the set of strings

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

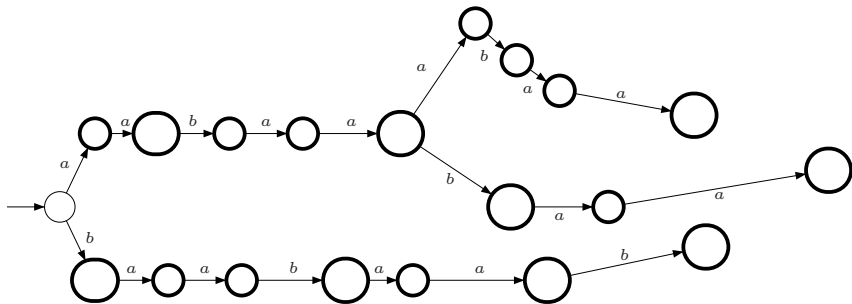
2. Build a trie \mathcal{T} on X
3. Build a Aho-Corasick like automaton upon \mathcal{T} . For each node ν of \mathcal{T} with “access word” v , use the transition function δ

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

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

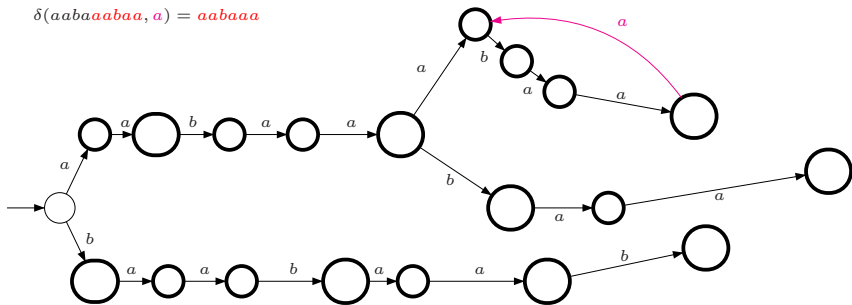


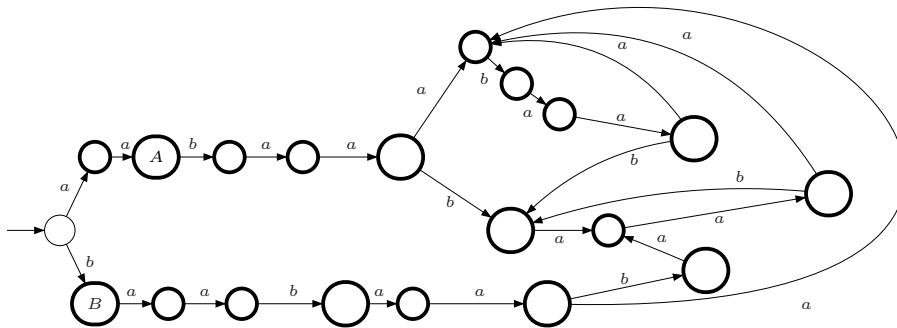
$X = \{a b a a, a b a a b a a, a a b a a b a a, a b a a b, b a a b, b a a b a a b\}$



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

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

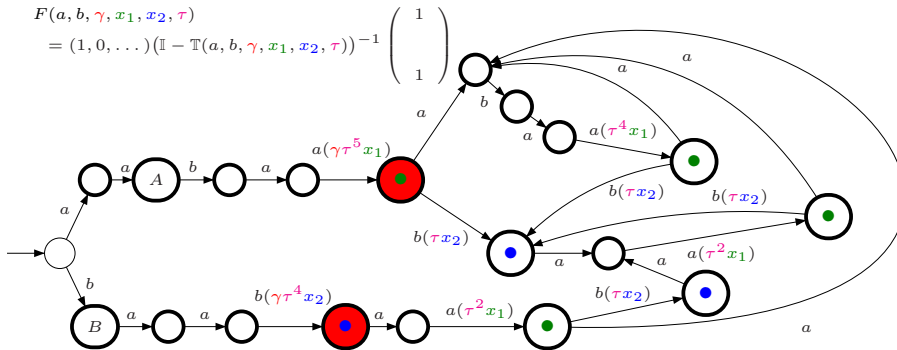




An automaton for $\mathcal{V} = \{v_1 = aabaa, v_2 = baab\}$. All transitions labeled by a and b ending respectively on state A and B are omitted.

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

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



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

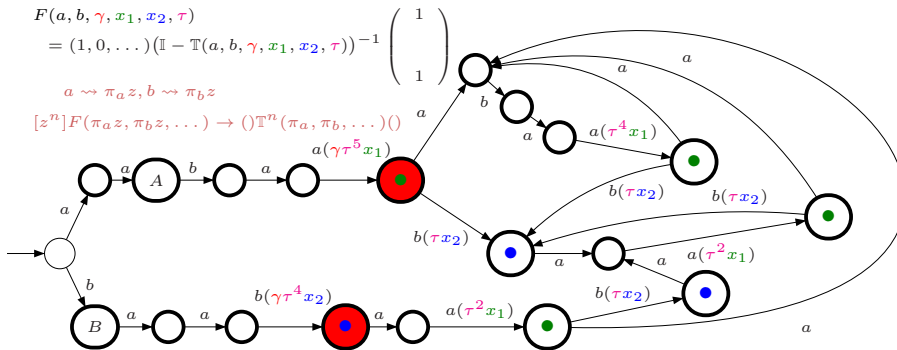
- ▶ $\gamma \rightarrow$ the **number of clumps**
- ▶ $\tau \rightarrow$ total **length of clumps**
- ▶ $x_1, x_2 \rightarrow$ occurrences of $aabaa, baab$

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

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

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

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



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

Formal weights on transitions

- ▶ γ → the **number of clumps**
- ▶ τ → total **length of clumps**
- ▶ x_1, x_2 → occurrences of $aabaa, baab$

Asymptotic Limit Laws

- ▶ **One word**

- $O_n^{\mathfrak{K}} = O(1)$ Poisson law for the number of clumps

- ▶ **General non-reduced sets**

- $O_n^{\mathfrak{K}} = \Theta(n)$ Normal limit law (number of clumps, size covered)

Proofs

- ▶ Poisson law: Rouché theorem, singularity analysis

- ▶ Normal law: automaton, Perron-Frobenius for $\mathbb{T}(\dots)$, singularity analysis, large powers theorem

Complexity of computing the prefix code(s)

- ▶ **one word**

$$|\mathcal{K}| \log(|\mathcal{K}|)$$

- ▶ **several words** (reduced case)

$$\left(\sum_{i,j} |\mathcal{K}_{i,j}| \right) \log \left(\sum_{i,j} |\mathcal{K}_{i,j}| \right)$$

Complexity of **insertion** of random keys in **tries**

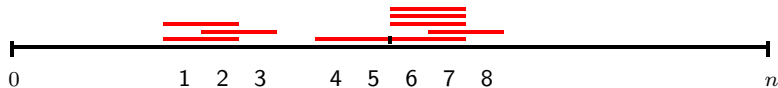
Throwing flat dimers on an “integral” segment

1 2



Throwing flat dimers on an “integral” segment

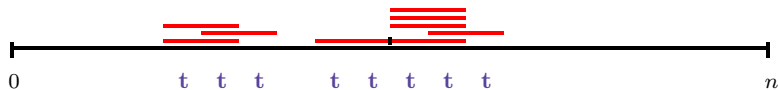
1 2



Size covered by the **dimers**?

Throwing flat dimers on an “integral” segment

aa stuttering occurrences of **aa** on $(a + b)^n$

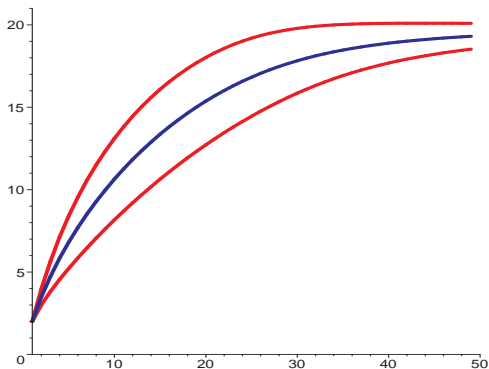


$$P(a) = P(b) = \frac{1}{2} \quad P(aa) = \frac{1}{4}$$

stuttering by substituting $x \rightsquigarrow \frac{x}{1 - \frac{1}{4}x}$ in $\Gamma(zt, x)$

Throwing flat dimers on an “integral” segment

Size covered



Number of dimers

$$n = |\text{segment}| = 20$$

Part III

Application of Clump Analysis to Genomics

Revisiting waiting times in DNA evolution

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)**, $\alpha, \beta \in \{A, C, G, T\}$

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)

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

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

Numerical remarks

- ▶ **length** of promoters $n \in [500 - 2000]$
- ▶ $k \in \{5, 6, 7, 8, 9, 10\}$ for the k -mers
- ▶ **Mutation probability** $\pi_{\alpha \rightarrow \beta} \approx 10^{-9}$

Numerical remarks

- ▶ **length** of promoters $n \in [500 - 2000]$
- ▶ $k \in \{5, 6, 7, 8, 9, 10\}$ for the **k -mers**
- ▶ **Mutation probability** $\pi_{\alpha \rightarrow \beta} \approx 10^{-9}$

We have

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

Numerical remarks

- ▶ **length** of promoters $n \in [500 - 2000]$
- ▶ $k \in \{5, 6, 7, 8, 9, 10\}$ for the k -mers
- ▶ **Mutation probability** $\pi_{\alpha \rightarrow \beta} \approx 10^{-9}$

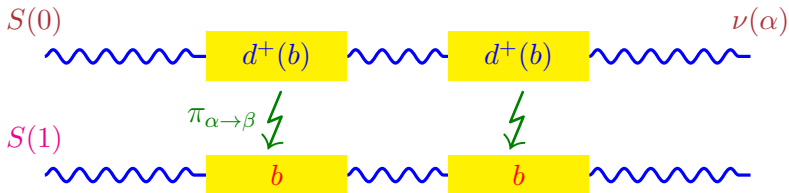
We have

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

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

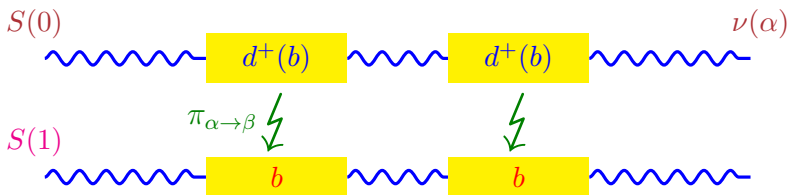
Behrens-Vingron 2010

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



Behrens-Vingron 2010

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



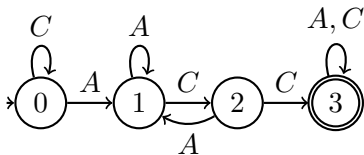
$$\left\{ \begin{array}{l} \mathfrak{P}_n \approx \sum_{i=1}^{\lfloor n/k \rfloor} (-1)^{i+1} \binom{n - i(k-1)}{i} \Phi^i \\ \Phi = \sum_{(a_1, \dots, a_k) \in \mathcal{A}^k \setminus \{b_1, \dots, b_k\}} \nu(a_1) \times \dots \times \nu(a_k) \cdot \prod_{j=1}^k \pi_{a_j \rightarrow b_j}(1) \end{array}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

Behrens-Nicaud-P.N. 2012

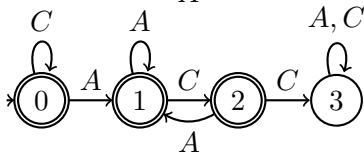
Construct an **automaton**

- ▶ 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



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

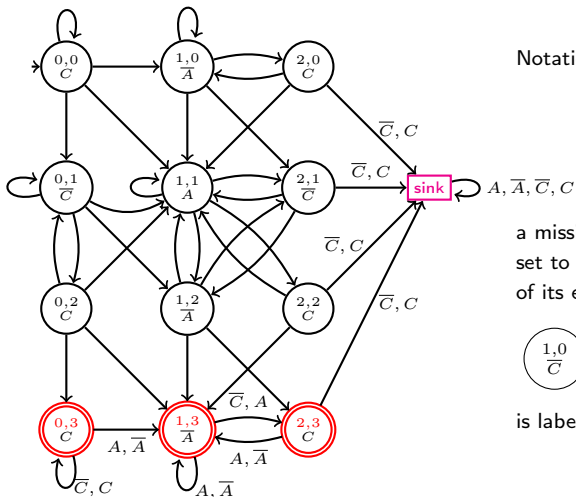


$$\overline{\mathcal{M}}_{\text{ACC}} = \{Q, \delta, s = 0, Q \setminus F\}$$

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

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

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



Notations for the transitions:

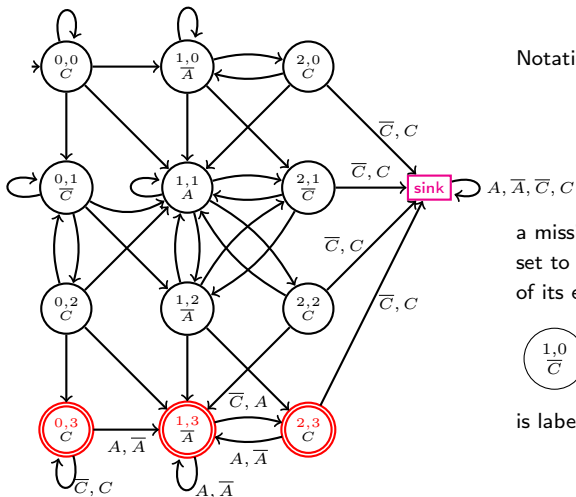
$$\begin{cases} A = \begin{pmatrix} A \\ A \end{pmatrix}, & C = \begin{pmatrix} C \\ C \end{pmatrix} \\ \overline{A} = \begin{pmatrix} A \\ \overline{A} \end{pmatrix}, & \overline{C} = \begin{pmatrix} C \\ \overline{C} \end{pmatrix} \end{cases}$$

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:

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

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		$\frac{\mathbf{E}_{\text{BNN}}(T_{1000})}{\mathbf{E}_{\text{BV}}(T_{1000})}$
	$\mathbf{E}_{\text{BNN}}(T_{1000})/10^6$	Rank	$\mathbf{E}_{\text{BV}}(T_{1000})/10^6$	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
...

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

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.

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

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.

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

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

In a random sequence of length n , 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

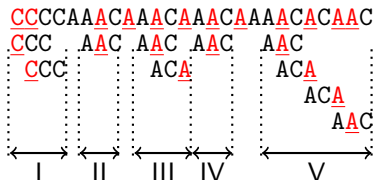
- ▶ i putative-hit positions $A \rightarrow C$
- ▶ and j putative-hit positions $C \rightarrow A$

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

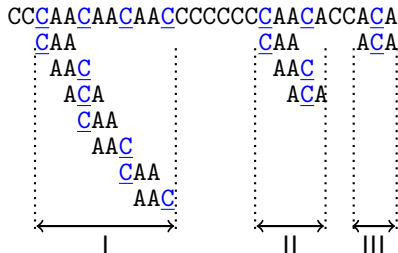
$$\mathcal{A} = \{A, C\} \quad b = ACC \longrightarrow d(ACC, 1) = \{\underline{C}CC, A\underline{A}C, AC\underline{A}\}$$



- ▶ (left) $b = ACC$ - 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

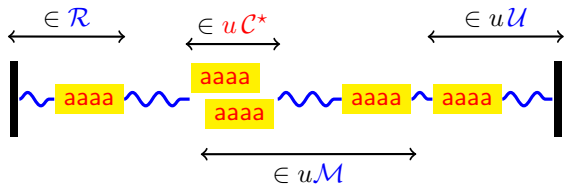
$$\mathcal{A} = \{A, C\} \quad b' = AAA \longrightarrow d(AAA, 1) = \{\underline{C}AA, A\underline{C}A, AA\underline{C}\}$$



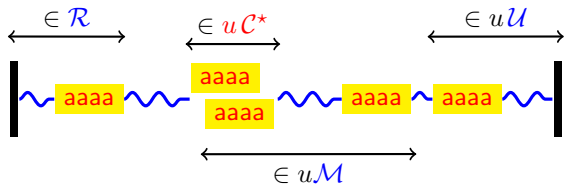
- ▶ (right) $b' = AAA$ - clump I contains 7 occurrences of $d(AAA)$, but only 4 putative-hit positions for $b' = AAA$. The number of word occurrences is not the relevant statistics for counting putative-hit positions

A - Language approach

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



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

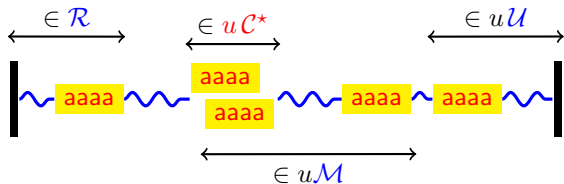


- ▶ residual language $\mathcal{D} = \mathcal{L}.u^-$: $\mathcal{D} = \{h, h \cdot u \in \mathcal{L}\}$
- ▶ $\mathcal{L}_2 - \mathcal{L}_1 = \mathcal{L}_2 \setminus \mathcal{L}_1 = \{h; h \in \mathcal{L}_2, h \notin \mathcal{L}_1\}$

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

Constrained Guibas-Odlyzko languages

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

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

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

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

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

Constrained clumps

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

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

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

Counting putative hit-positions

$$b = \text{ACAC} \quad \mathbf{P(A)} = \mathbf{P(C)} = 1/2$$

$$d_\ell(b) = (\text{AAAC}, \text{ACAA}, \text{ACCC}, \text{CCAC}) \quad d_\ell(b)(z, t) = \left(\frac{z^4 t}{16}, \frac{z^4 t}{16}, \frac{z^4 t}{16}, \frac{z^4 t}{16} \right)$$

$$\widehat{\mathbb{K}}_b(z, t) = \begin{pmatrix} 0 & \frac{z^2 t}{4} & \frac{z^2 t}{4} & \frac{z^3 t}{8} \\ \frac{z^3 t}{8} + \frac{z^2}{4} & \frac{z^3 t}{8} & \frac{z^3 t}{8} & 0 \\ 0 & 0 & 0 & \frac{z^3 t}{8} + \frac{z^2}{4} \\ 0 & \frac{z^2 t}{4} & \frac{z^2 t}{4} & \frac{z^3 t}{8} \end{pmatrix}$$

The extension $\text{AC} \in \widehat{\mathcal{K}}_{\text{ACAA}, \text{AAAC}}$ as in $\text{ACAA}|\text{AC}$ leads to *no new* putative-hit position, in contrast to $\text{ACAA}|\text{AAC}$

$$\widehat{\mathcal{K}}_{\text{ACAA}, \text{AAAC}}(z, t) = \frac{z^3 t}{8} + \frac{z^2}{4}$$

Generating function of the number of putative-hit positions

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

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

Generating function of the number of putative-hit positions

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

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

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

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

Generating function of the number of putative-hit positions

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

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

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

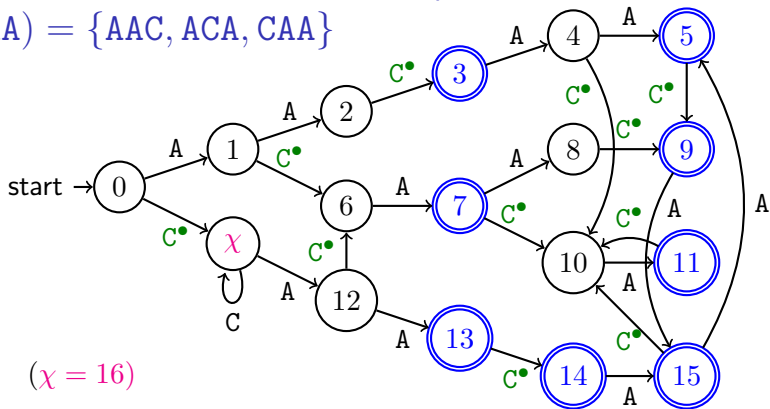
$$\widehat{\mathbb{G}}(z, t) = \left(v_i(z, t) \widehat{\mathbb{S}}_{ij}(z, t) \right).$$

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

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

B - Automaton approach

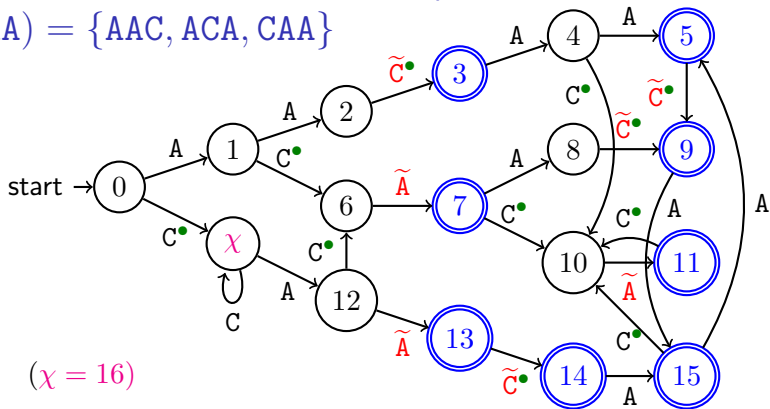
Automaton for constrained clumps of $d(AAA) = \{AAC, ACA, CAA\}$



($\chi = 16$)

- ▶ **Double circles** signal 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**

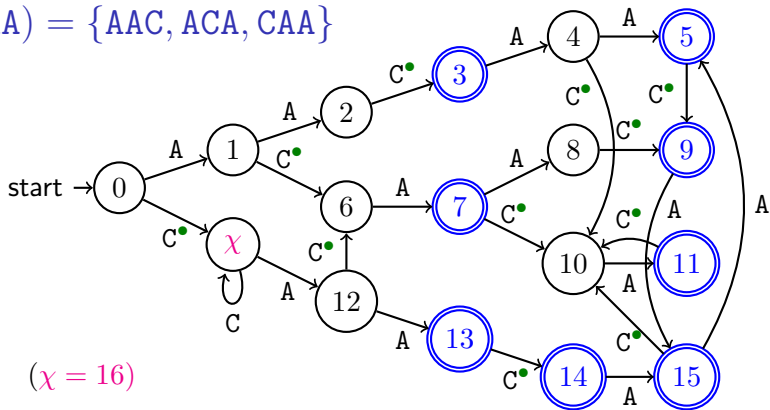
Automaton for constrained clumps of $d(AAA) = \{AAC, ACA, CAA\}$



($\chi = 16$)

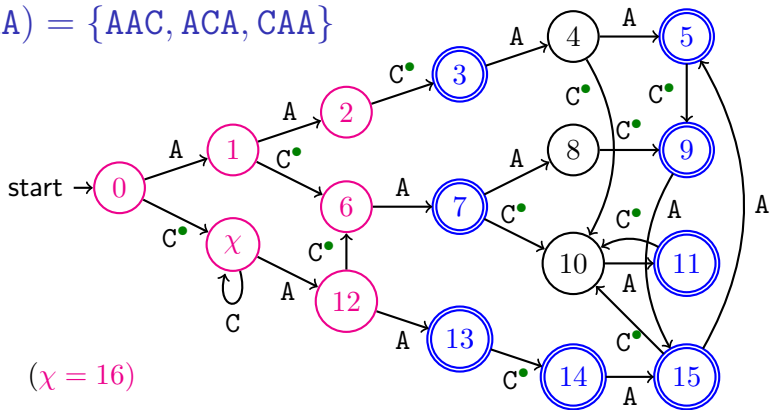
- ▶ **Double circles** signal an occurrence of a word of $d(aaa)$.
- ▶ **Avoiding** AAA leads to **missing transitions** A
- ▶ The **missing transitions** C **point** to the **state** χ .
- ▶ **• characters** mark **putative-hit-positions**
- ▶ Transitions covered by tildes (\tilde{A}, \tilde{C}) emit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

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



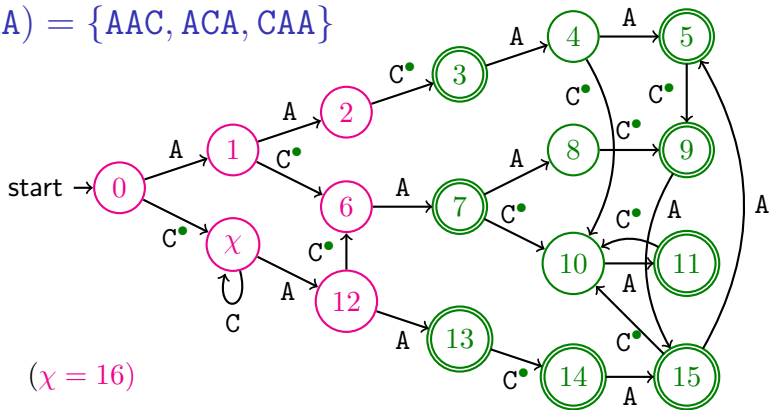
- $O = \{q, \delta(0, w) = q, w \in X\}$, (**occurrence** of a word of $d(aaa)$).

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



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

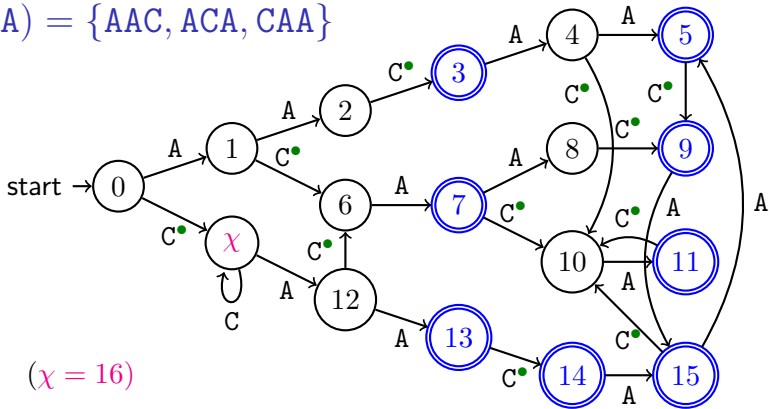
Automaton for constrained clumps of $d(AAA) = \{AAC, ACA, CAA\}$



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

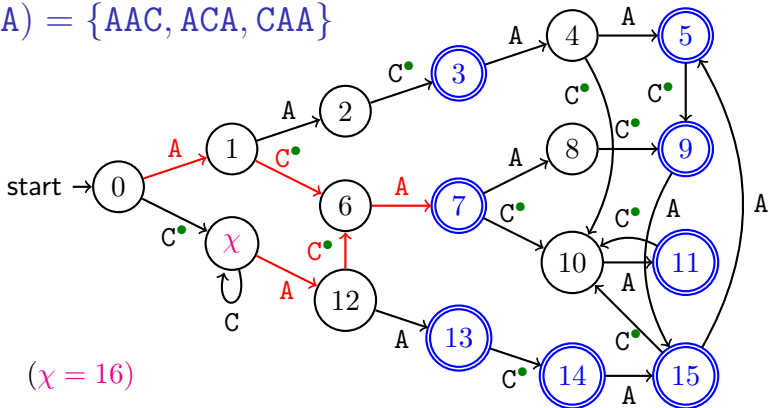
Definition of an auxiliary function θ

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



Definition of an auxiliary function θ

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

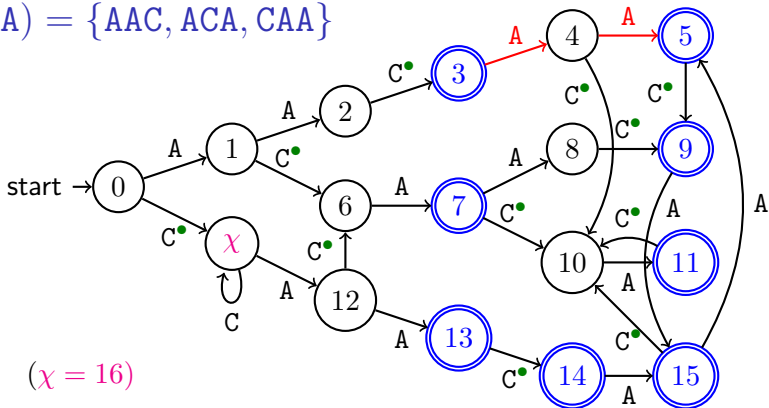


$(\chi = 16)$

$$\theta(7) = ACA$$

Definition of an auxiliary function θ

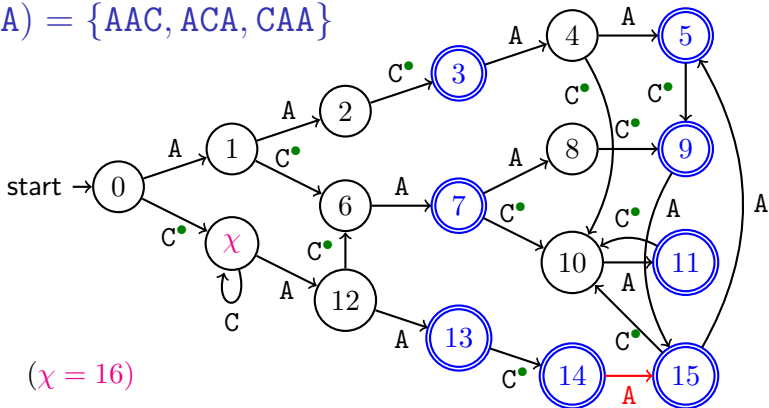
$$d(AAA) = \{AAC, ACA, CAA\}$$



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

Definition of an auxiliary function θ

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



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

Formal definition of θ

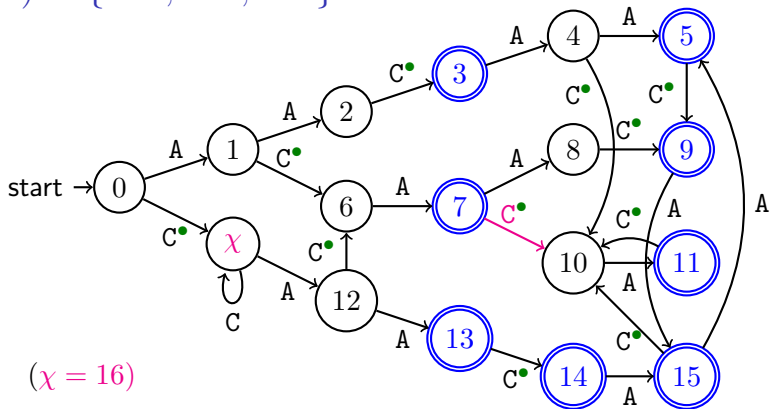
For each state $o \in O$ (recognizing an **occurrence** of $d(b)$),

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

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

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

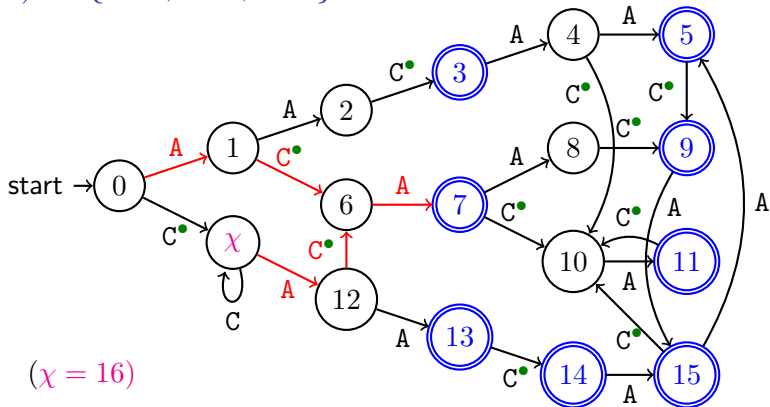
$d(\text{AAA}) = \{\text{AAC}, \text{ACA}, \text{CAA}\}$



► $h_{7,10} = \nu_C z$

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

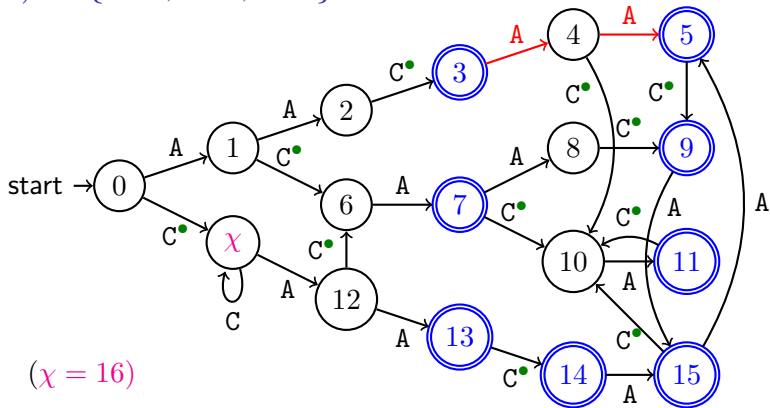
$d(AAA) = \{AAC, ACA, CAA\}$



- ▶ $h_{7,10} = \nu_C z$
- ▶ $h_{6,7}(t) = \nu_{A^t C \rightarrow A} z$

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

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

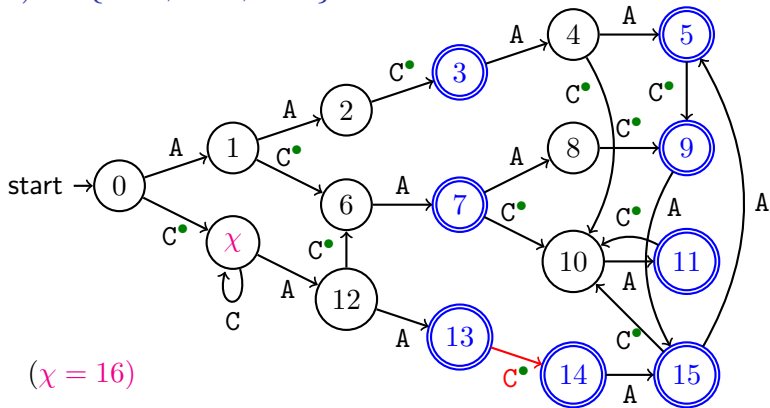


$(\chi = 16)$

- ▶ $h_{7,10} = \nu_C z$
- ▶ $h_{6,7}(t) = \nu_A t_{C \rightarrow A} z$
- ▶ $h_{3,4}(t) = h_{45}(t) = \nu_A z$

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

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



- ▶ $h_{7,10} = \nu_C z$
- ▶ $h_{6,7}(t) = \nu_A t_{C \rightarrow A} z$
- ▶ $h_{3,4}(t) = h_{45}(t) = \nu_A z$
- ▶ $h_{13,14}(t) = \nu_C t_{C \rightarrow A} z$

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

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

From matrix to generating function

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

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)

Rational functions, Taylor coefficient of order n and gfun

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

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

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

Rational functions, Taylor coefficient of order n and gfun

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

recurrence equations \rightarrow **gfun**[rectoproc] \rightarrow **procedure** Proc(n)

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

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

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 ,

Rational functions, Taylor coefficient of order n and gfun

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

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

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

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 ,

► $\widehat{f}_n^{(b)} = \mathbf{P}(S_n(0))$

Rational functions, Taylor coefficient of order n and gfun

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

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

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

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 ,

- ▶ $\widehat{f}_n^{(b)} = \mathbf{P}(S_n(0))$
- ▶ η_n is the **unconditioned 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 **gfun**[diffeqtorec] \rightarrow **recurrence** equations

recurrence equations \rightarrow **gfun**[rectoproc] \rightarrow **procedure** Proc(n)

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

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

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 ,

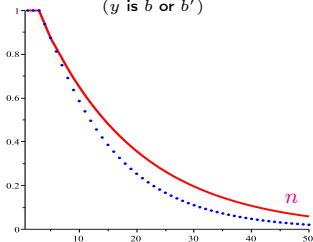
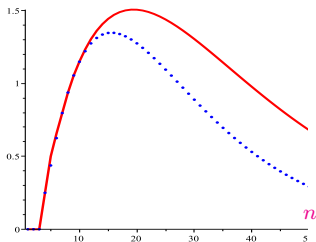
- ▶ $\widehat{f}_n^{(b)} = \mathbf{P}(S_n(0))$
- ▶ η_n is the **unconditionned probability** of the expectation of the count of putative-hit positions
- ▶ **Conditionned expectation:** $\widetilde{\eta}_n = \eta_n / \widehat{f}_n^{(b)}$

An unexpected behaviour

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

$$\hat{f}_n^{(y)} = \mathbf{P}(|S_n(0)|_y = 0)$$

(y is b or b')



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

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

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

$$\mathbf{E}(H) = \left. \frac{\partial F_b(z, t)}{\partial t} \right|_{t=1}$$

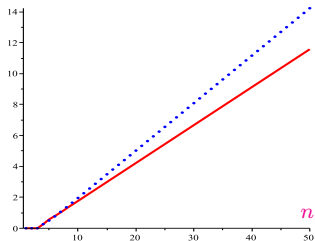
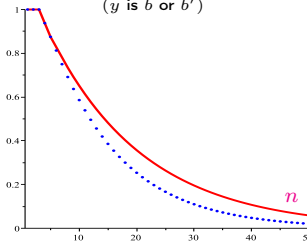
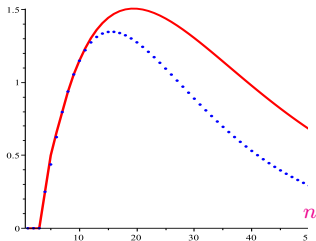
An unexpected behaviour

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

$$\hat{f}_n^{(y)} = \mathbf{P}(|S_n(0)|_y = 0)$$

(y is b or b')

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



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

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

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

$$\mathbf{E}(H) = \left. \frac{\partial F_b(z, t)}{\partial t} \right|_{t=1}$$

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 .

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

A proof by singularity analysis

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

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

$\widehat{f}_n^{(b)}$ 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

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

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

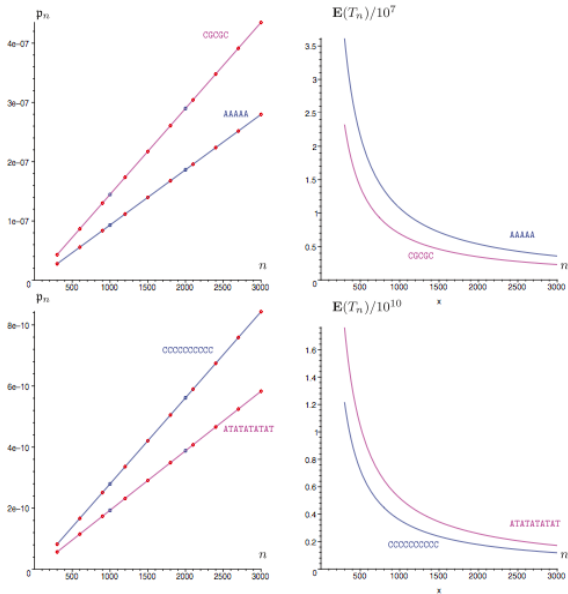
General case

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**



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

FIG. 5. Plots of the probability p_n (left) and of the expected waiting time $E(T_n)$ (right). (Top) $b = AAAAA$ (blue) and $b' = CGCGC$ (magenta). (Bottom) $b = CCCCCC$ (blue) and $b' = ATATATA$ (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.