

Constructions for Clumps Statistics

P. Nicodème

CNRS - LIX École polytechnique, Palaiseau

Joint work with F. Bassino, J. Clément and J. Fayolle

00100

00100
11100

00100
111100
10111111100

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	000101010011110001
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	000101010011110001
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	000101010011110001
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	000101010011110001
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	000101010011110001
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	000101010011110001
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	000101010011110001
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	000101010011110001
20		111101011
21	14	10101001001000
22	15	101001111100
23		0111010111
24	16	0001001010110
25		111110111
26		0110111000011
27	17	1000000000
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

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...

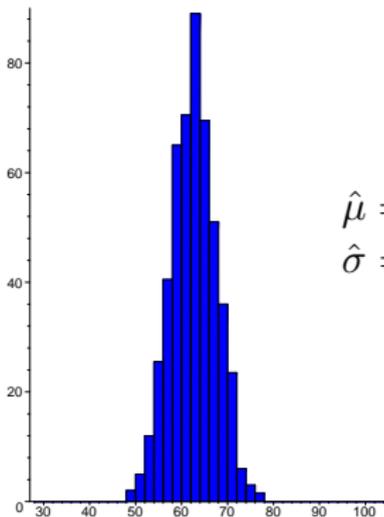
...011110

111...

- ▶ while the 100 NEVER OVERLAP

Expected waiting time: 111 – 14 100 – 7

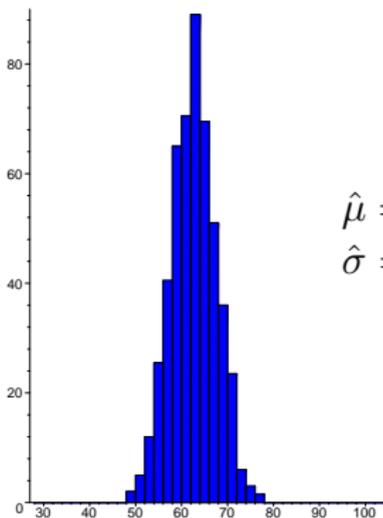
Number of occurrences 1000 texts of size 500



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

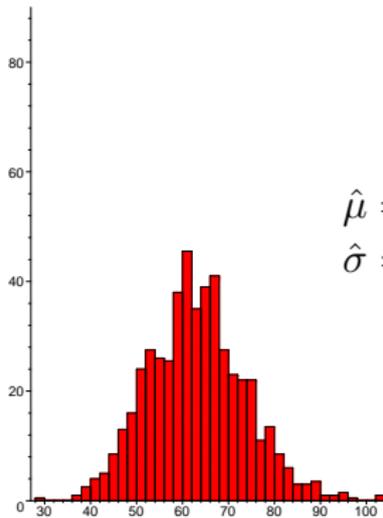
100

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

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
- ▶ future extensions to **tandem repeats**?

... GGG**GATCGA**|**GATCGA**|**GATCGA**|**GATCGA**GGG ...

Counting clumps

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

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

Definition of a clump of a word w

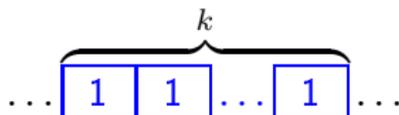
Set of **maximal contiguous** positions such that:

- ▶ each position is inside an occurrence of w
- ▶ if a position is at the beginning or at the end of an occurrence of w , it is covered by two occurrences of w , with exception of the first and last positions of the clump

Probability of start at a position i

(A) word $w = 1^k$

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

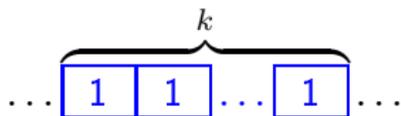


$$\mathbf{P}(\text{start}) = p^k$$

Probability of start at a position i

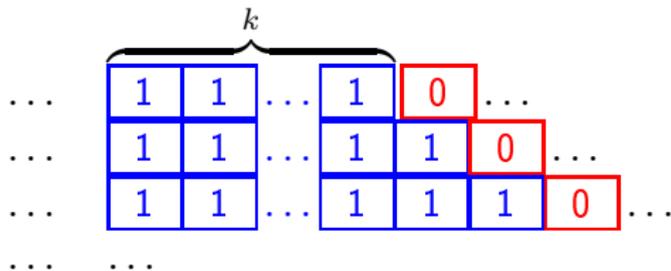
(A) word $w = 1^k$

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



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

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

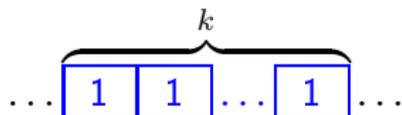


$$\mathbf{P}(\text{start}) = p^k \times q \times (1 + p + p^2 + \dots) = p^k \times \frac{q}{1 - p} = p^k$$

Probability of start at a position i

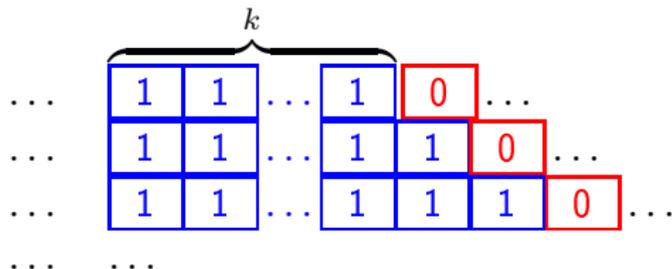
(A) word $w = 1^k$

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



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

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

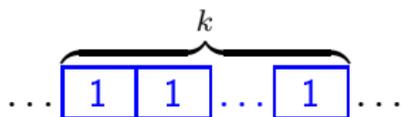


$$\mathbf{P}(\text{start}) = p^k \times q \times (1 + p + p^2 + \dots) = p^k \times \frac{q}{1 - p} = p^k$$

FALSE

Probability of start at a position i

(A) word $w = 1^k$



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

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

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

no clump **beginning** at position i

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

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.

Combinatorial approach

Jacquet-Szpankowski

I- Combinatorial approach

1 word - Bernoulli model

\mathcal{A} alphabet, w considered word

Polynomial of autocorrelation

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

$w = ababa$

$ababa$

$ababa|$

$ababa$

$ababa$

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

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

Régnier and Szpankowski decomposition

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

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

Régnier and Szpankowski decomposition

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

aaaaa**ababa** $\in \mathcal{R}$, bbbbabab**aba** $\notin \mathcal{R}$

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

ababa aaaa**ababa** $\in \mathcal{M}$ ababa bbbbbbbbab**aba** $\notin \mathcal{M}$ ababa**ba** $\in \mathcal{M}$

Régnier and Szpankowski decomposition

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

$aaaaa$ *ababa* $\in \mathcal{R}$, $bbbb$ *abababa* $\notin \mathcal{R}$

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

ababa $aaaa$ *ababa* $\in \mathcal{M}$ *ababa* $babbbbbbbb$ *ababa* $\notin \mathcal{M}$ *ababa* ba $\in \mathcal{M}$

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

ababa $aabbbabbbbbbb$ $\in \mathcal{U}$ *ababa* $babbbbbbbb$ $\notin \mathcal{U}$

Languages \Rightarrow generating functions

$$\begin{aligned} \text{(I)} \quad \Sigma^* &= \mathcal{U} + \mathcal{M}\Sigma^* &\Rightarrow \quad \frac{1}{1-z} &= U(z) + \frac{M(z)}{1-z} \\ \text{(II)} \quad \Sigma^*w &= \mathcal{R}\mathcal{C} + \mathcal{R}\Sigma^*.w &\Rightarrow \quad \frac{\pi_w z^m}{1-z} &= R(z) \left(C(z) + \frac{\pi_w z^m}{1-z} \right) \\ \text{(III)} \quad \mathcal{M}^+ &= \Sigma^*.w + \mathcal{C} - \epsilon &\Rightarrow \quad \frac{M(z)}{1-M(z)} &= \frac{\pi_w z^m}{1-z} + C(z) - 1 \\ \text{(IV)} \quad \mathcal{N}.\Sigma &= \mathcal{R} + \mathcal{N} - \epsilon &\Rightarrow \quad zN(z) &= R(z) + N(z) - 1 \end{aligned}$$

Add the mark x after each match \Rightarrow marked language \mathcal{L}

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

Translation into generating function

$$L(z, x) = N(z) + \frac{xR(z)U(z)}{1-xM(z)}$$

Solving the system

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

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

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

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

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

Language equations for the clumps

Notation:

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

Combinatorial decomposition

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

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

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

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

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

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

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

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

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

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

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

One word - Expectation - Variance

clumps

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

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

one word

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

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

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 $W = \{aab, abaa\}$

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

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

The set of language equations

$$\bigcup_{k \geq 1} (\mathbb{M}^k)_{i,i} = \mathcal{A}^* \cdot w_i + \mathcal{C}_{ii} - \epsilon, \quad 1 \leq i \leq 2,$$

$$\bigcup_{k \geq 1} (\mathbb{M}^k)_{i,j} = \mathcal{A}^* \cdot w_j + \mathcal{C}_{ij}, \quad 1 \leq i \neq j \leq 2,$$

$$\mathcal{U}_i \cdot \mathcal{A} = \bigcup_j \mathcal{M}_{ij} + \mathcal{U}_i - \epsilon, \quad 1 \leq i \leq 2,$$

$$\mathcal{A} \cdot \mathcal{R}_j - (\mathcal{R}_j - w_j) = \bigcup_i w_i \mathcal{M}_{ij}, \quad 1 \leq j \leq 2,$$

$$\mathcal{N} \cdot w_j = \mathcal{R}_j + \bigcup_i \mathcal{R}_i (\mathcal{C}_{ij} - \delta_{ij} \epsilon), \quad 1 \leq i, j \leq 2$$

$$F(z, x, y) = \mathbb{R}(z, x, y) (\mathbb{I} - \mathbb{M}(z, x, y))^{-1} \mathbb{U}(z)$$

Putting up equations for clumps of two words

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

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

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

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

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

II- 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

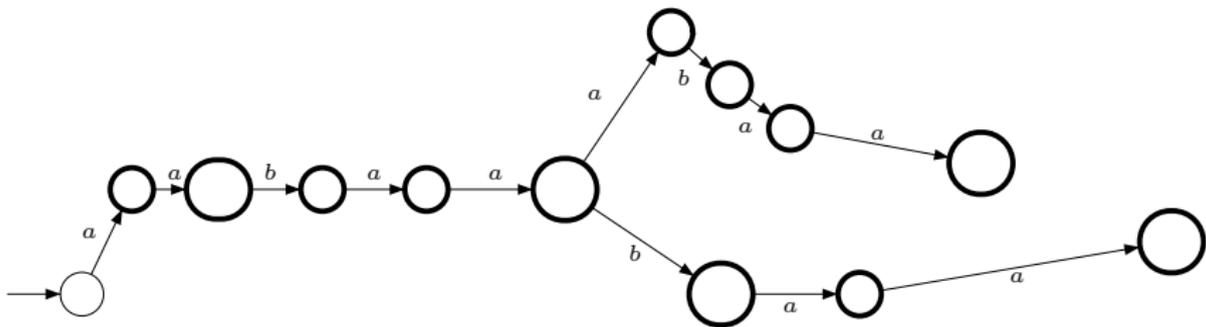
1. Build the set of strings

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

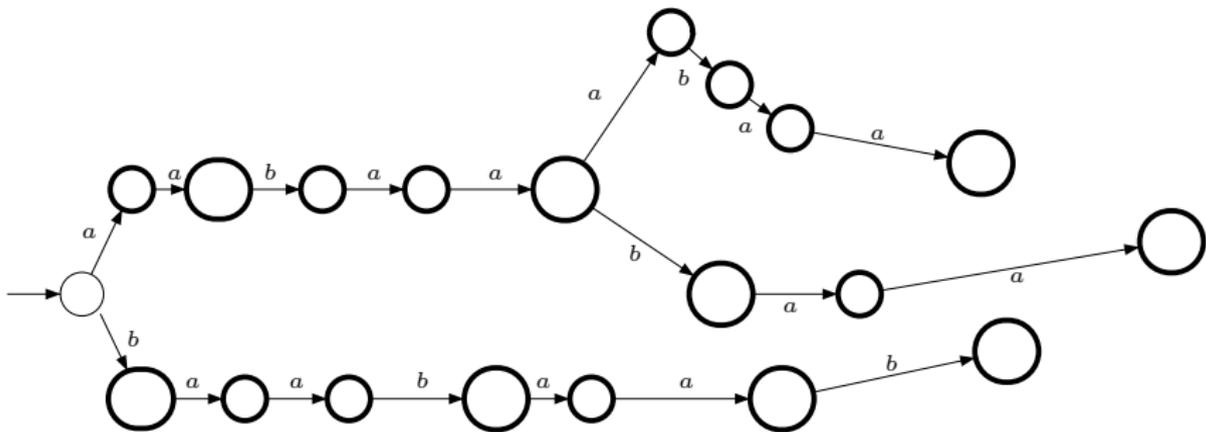
2. Build a trie \mathcal{T} on X
3. Build a Aho-Corasick like automaton upon \mathcal{T} . For each node ν of \mathcal{T} with “access word” v , use the transition function δ

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

$X = \{aabaa, abaabaa, aabaabaa, aabaab, \quad \}$

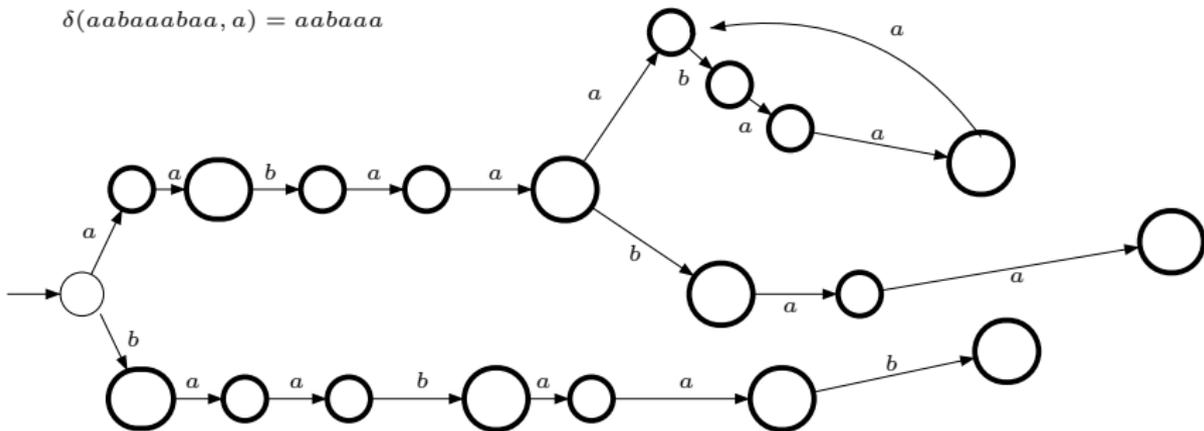


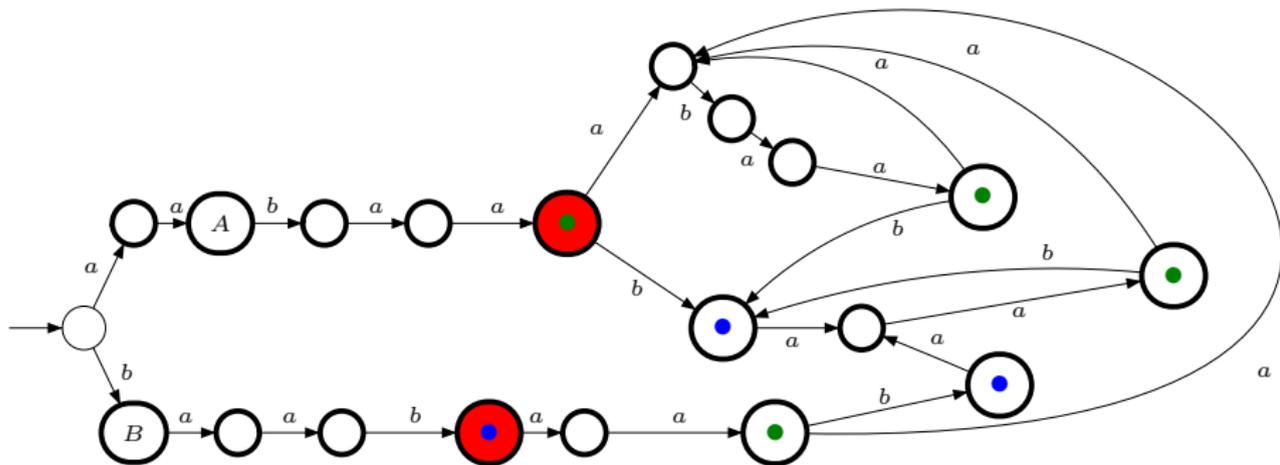
$X = \{aabaa, abaabaa, abaaabaa, abaab, baab, baabaab\}$



$X = \{aabaa, abaabaa, abaaabaa, abaab, baab, baabaab\}$

$\delta(aabaabaa, a) = aabaaa$





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**

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**

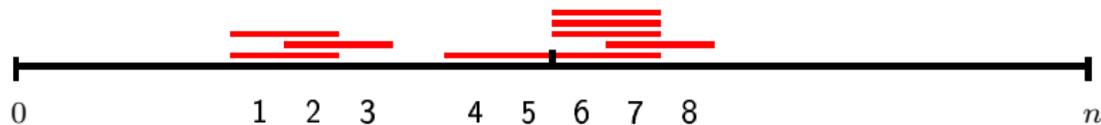
Throwing flat dimers on an “integral” segment

1 2



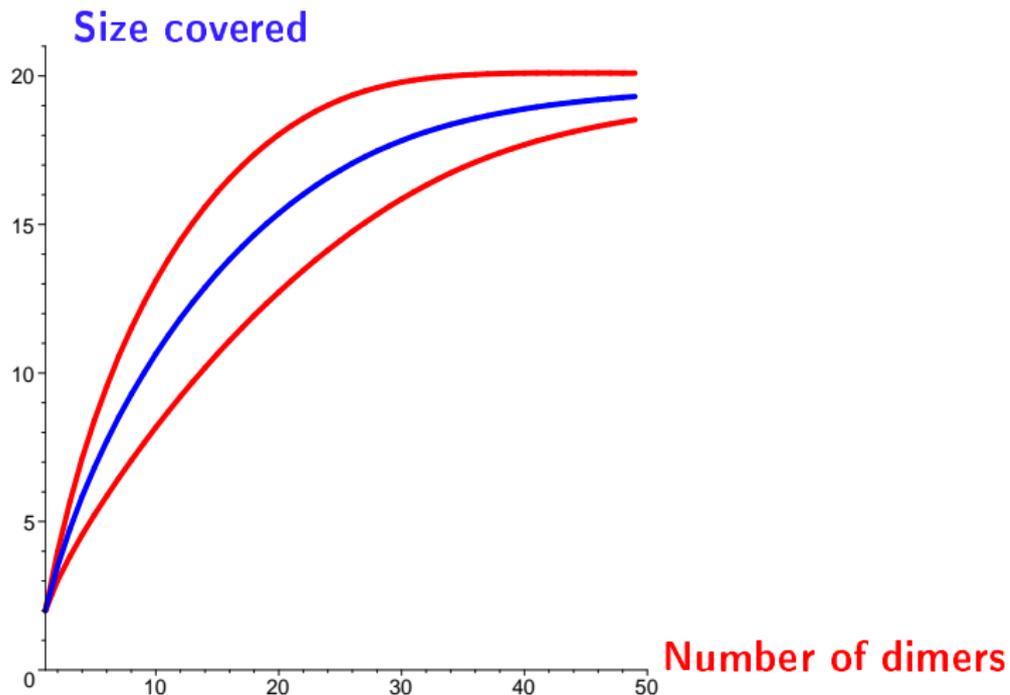
Throwing flat dimers on an “integral” segment

1 2



Size covered by the **dimers**?

Throwing flat dimers on an “integral” segment



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

The Prefix Code $\mathcal{K} = \mathcal{C}_o - \mathcal{C}_o\mathcal{A}^+$ generates \mathcal{C}_o^+

w word, $\mathcal{C}_o = \mathcal{C} - \epsilon$, $v \in \mathcal{C}_o \implies v = \kappa_1 \dots \kappa_r$, $\kappa_i \in \mathcal{K}$

- ▶ $v \notin \mathcal{K} \implies v = \kappa_1 v'$ with $\kappa_1 \in \mathcal{K}$ and $v' \in \mathcal{A}^+$.
- ▶ $v, \kappa_1 \in \mathcal{C}_o \implies \exists p, p', s'$ such that $p = p's'$, $wv = pw$ and $w\kappa_1 = p'w$. Therefore $wv' = s'w$ and $v' \in \mathcal{C}_o$.
- ▶ iterate a **finite** number of time

