

Алгоритмические аспекты анализа символьных последовательностей

Мирошниченко Любовь Александровна
Институт математики СО РАН
luba@math.nsc.ru

Объект исследования: символьные последовательности.

Σ – непустое конечное множество символов (алфавит);

$T = t_1 t_2 \dots t_i \dots t_N$ ($t_i \in \Sigma, 1 \leq i \leq N$) – последовательность символов, цепочка символов, текст, строка, слово.

Примеры:

- двоичные последовательности ($|\Sigma| = 2$);
- последовательность действий, траектория движения;
- тексты программ;
- тексты естественного языка;
- музыкальные тексты (песенные мелодии);
- древнерусские церковные песнопения;
- ДНК, РНК ($|\Sigma| = 4$);
- аминокислотные послед. ($|\Sigma| = 20$);
- порядки генов;
- порядки дисков политенных хромосом ...

Основные задачи анализа символьных последовательностей:

- поиск образцов;
- восстановление структуры текста: выявление повторов (периодичностей, симметрий ...);
- сложность текста

Сравнение последовательностей

Анализ параллельных текстов

Поиск образцов

$P = p_1 p_2 \dots p_m$ – образец, $T = t_1 t_2 \dots t_N$ – текст, $t_i \in \Sigma, p_j \in \Sigma, 1 \leq i \leq N, 1 \leq j \leq m, m < N$

Задача: обнаружить все вхождения образца P в текст T .

Обобщения:

- Образцы с джокерами: $p_j \in \Sigma \cup \{x\}$, x – любой символ
- Образцы, позиции которых заданы множествами символов:
 $P = a-[ag]-c-[cg]-[acg]-[ct]-a, P = a-[ag]-c-[cg]-\neg t-x-a$
- Поиск образца с «ошибками»
- Поиск множества образцов (попевки)
- Поиск множества образцов, позиции которых заданы множествами символов. Кол-во образцов $\sim 10^5$. Алгоритмы используют «ядерную» стратегию в сочетании с построением ДКА.

Бачинский А.Г., Гусев В.Д., Немытикова и др. Новая версия банка образов PROF-PAT 1.0: технология формирования и программа быстрого поиска образов в аминокислотных последовательностях // Молекулярная биология, 1996, Т. 30, вып. 6, С. 1409–1417

- Образцы с переменными
 $P = abXXcX : abttct; ab\underline{abbabb}c\underline{abb}$
- Параметризованные образцы: 2 алфавита: Σ и Π :
Образцы $abcXbbYUccZ$ и $abcZbbXXccY$ π -согласованы

Повтор — пара совпадающих фрагментов текста

Классификация повторов

По характеру расположения в тексте:

- разнесенные ... agttc ... agttc...
- тандемные ... agttcagttc...
- с наложением : ... agttcagttcagttc ...

По наличию искажений:

- совершенные ... agttc ... agttc ...
- несовершенные ... agttc ... aattc ... , ... agttc ... agtttc... ,
- с точностью до агрегирования: ... agttc ... gatct ... ($\{a,g\} \rightarrow Pu, \{c,t\} \rightarrow Py$).
- с точностью до подстановки на элементах алфавита:
... agttc... tcaag... ($a \leftrightarrow t, c \leftrightarrow g$); секвентные переносы в музыке;

По направлению считывания

- прямые : ... agttc ... agttc ...
- симметричные: ... agttc ... cttga...
- инвертированные: ... $\overrightarrow{\text{agttc}}$... $\overleftarrow{\text{gaact}}$...

Представление текста в терминах повторов

Σ – конечный алфавит;

$S = s_1 s_2 s_3 s_4 s_5 \dots s_N$ – текст, составленный из элементов Σ ;

$S [i : j] = s_i \dots s_j$ – фрагмент текста ($1 \leq i < j \leq N$).

l -грамма – фрагмент текста длины l ($S [i : i + l - 1]$).

$$S = \underline{s_1 s_2 s_3 s_4 s_5} \dots s_N$$

Полное число l -грамм: $N - l + 1$.

Число различных l -грамм: $M_l \leq N - l + 1$.

$l_{\max}(S)$ – наибольшее l , при котором в S есть **повторяющиеся** l -граммы.

Частотная характеристика l -го порядка текста S –

множество элементов $\Phi_l(S) = \{ \phi_{l1}, \phi_{l2}, \dots, \phi_{lM_l} \}$

где ϕ_{li} ($1 \leq i \leq M_l$) : « i -я l -грамма – \mathbf{x}_i , ее частота в тексте – $\mathbf{F}_l(\mathbf{x}_i)$ ».

Совокупность частотных характеристик

$$\Phi(S) = \{ \Phi_1(S), \Phi_2(S), \dots, \Phi_{l_{\max} + 2}(S) \}$$

называется **полным частотным спектром** текста S .

Наиболее важные параметры частотного спектра:

l_{\max} — длина максимального повтора в тексте.

M_l — размер «словаря» l -грамм.

$F_l^{\max} = \max_{1 \leq i \leq M_l} F_l(x_i)$ — максимальная частота встречаемости l -грамм в тексте ($1 \leq l \leq l_{\max}$);

E_l^k — Число различных l -грамм, встречающихся в тексте ровно k раз

E_l^0 — число l -грамм, ни разу не встретившихся в тексте.

$$E_l^0 = |\Sigma|^l - M_l. \quad M_l = \sum_{k=1}^{F_l^{\max}} E_l^k \quad N - l + 1 = \sum_{k=1}^{F_l^{\max}} k \cdot E_l^k$$

Позиционные аномалии:

- кластеры l -грамм
- сверхравномерное распределение
- «изолированные» точки
- «ГЭПЫ»

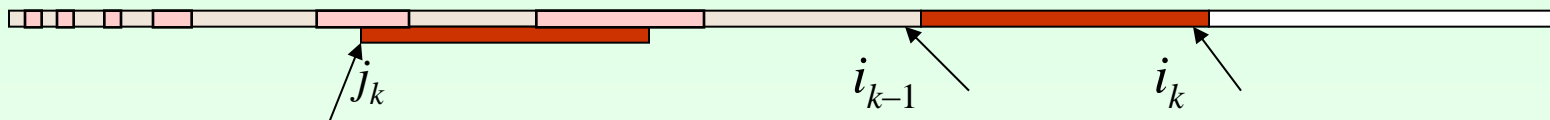
Сложность по Лемпелю и Зиву

Сложность $c_{LZ}(S) = m$ последовательности S : **число шагов процесса ее «порождения»**. Каждый шаг: **копирование** максимально длинного фрагмента из синтезированного префикса и/или **генерация** символа.

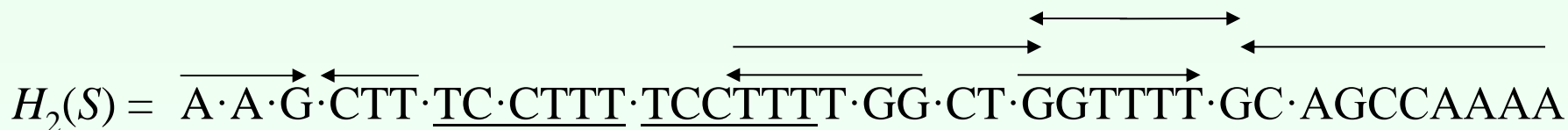
Сложностное разложение S есть конкатенация

$$H(S) = S[1:i_1] S[i_1 + 1:i_2] \dots S[i_{k-1} + 1:i_k] \dots S[i_{m-1} + 1:N],$$

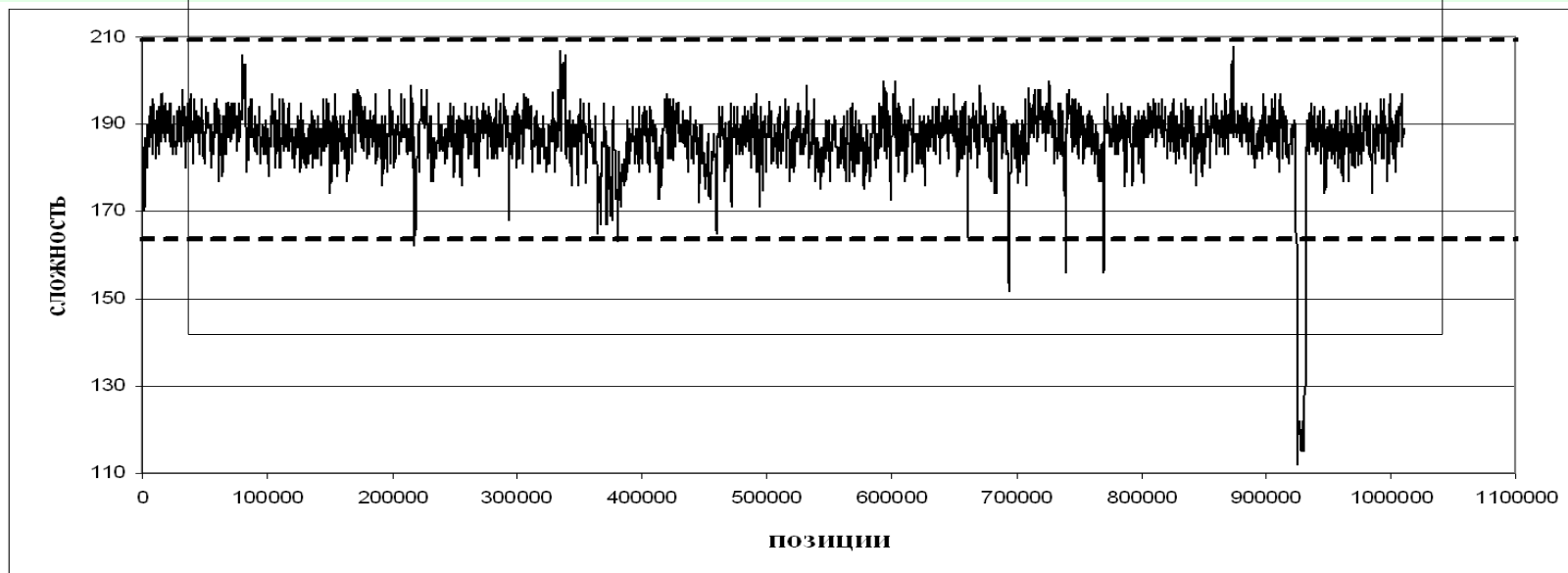
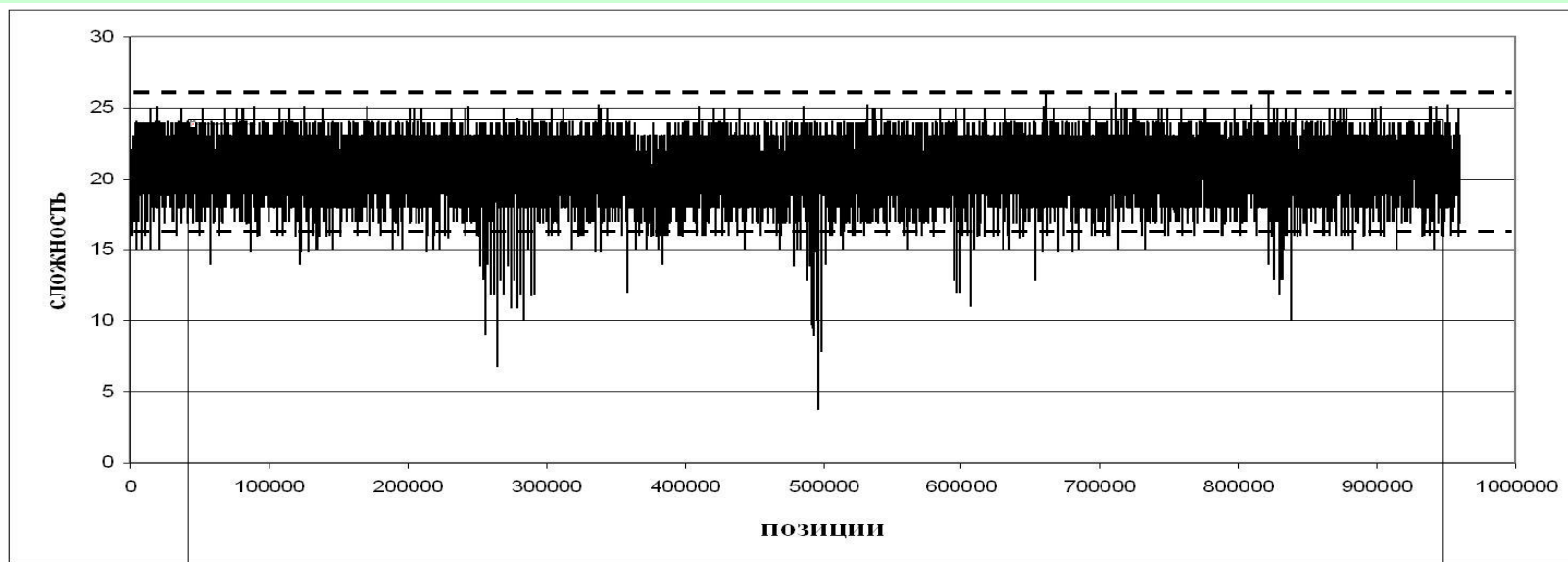
где $S[i_{k-1} + 1:i_k] = S[j_k:j_k + l_k - 1]$ — фрагмент, синтезируемый на k -м шаге, $j_k < i_{k-1} + 1$



ДНК-ориентированный вариант меры сложности наряду с прямыми повторами учитывает также симметричные и комплементарные (прямые и симметричные).



Профили сложности генома микоплазмы «R» при $W = 60$ и $W = 1000$



Регулярные повторы

gttttgggggttgtagaattatthttgttagtaaaac aatgataaataacgcttaacttgcttact
gttttgggggttgtagaattatthttgttagtaaaac cctataaacaatcaggattatatgtacta
gttttgggggttgtagaattatthttgttagtaaaac ttagtcaagattthttaataccaggggtgca
gttttgggggttgtagaattatthttgttagtaaaac tccatathtttcccttactattactatgct
gttttgggggttgtagaattatthttgttagtaaaac acgattthtaaaaattatgatataataataac
gttttgggggttgtagaattatthttgttagtaaaac tagaatctctthtaattcccaccaagct
gttttgggggttgtagaattatthttgttagtaaaac cthttthtaattthtattataactthgt
gttttgggggttgtagaattatthttgttagtaaaac aattgcatcattaacgthtaagacgthtact
gttttgggggttgtagaattatthttgttagtaaaac aatcaaacaactcgctthtctaaatcatcaa
gttttgggggttgtagaattatthttgttagtaaaac thtgagcataatggcgctthtgagctthtag
gttttgggggttgtagaattatthttgttagtaaaac aatgcagattthaaaagattcaggaacgatt
gttttgggggttgtagaattatthttgttagtaaaac attagccccacaattatattaacctccct
геном "*Mycoplasma synoviae* 53" (ID AE017245), поз. 690229

CRISPRs (*Clustered Regularly Interspaced Short Palindromic Repeats*)

Тексты песен: куплет – припев...

Фрактальные структуры:

$\overleftrightarrow{\text{cat}} \overleftrightarrow{\text{tac}} \overleftrightarrow{\text{cat}} \overleftrightarrow{\text{tac}}$

$\overleftrightarrow{\text{cat}} \overleftrightarrow{\text{atg}} \overleftrightarrow{\text{cat}} \overleftrightarrow{\text{atg}}$

Класс 1 : (ag-ga)(ag-ga)(ag-ga), поз. 5870, 25/80 (низковир.);
 (gtg)(gtg)(gtg)(gtg) = V^4 , поз. 805, 10/80 (высоковир.);

Класс 2 : ggaagg ggaagg gg = EGEG, поз. 9411, $27/46 + 1/28$

Фракталоподобные структуры:



Все ВКЭ: gagctcaaaactggagagctc, 146/161

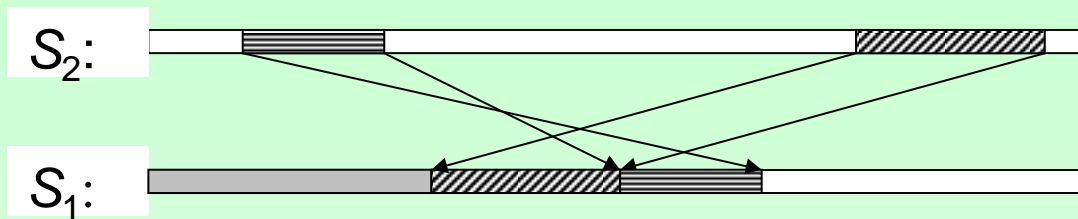
Класс 1 : gttggttgctggttg, поз. 308, 70/80

Класс 2 : cggcaccaaccggctcggcggc, поз. 8471, 41/46

Arabidopsis th (4-я хромосома), 194 вхождения:

$\overleftrightarrow{\text{tgtcga}} \overleftrightarrow{\text{tcgaca}} \text{tcaccatgag} \overleftrightarrow{\text{tgtcga}} \overleftrightarrow{\text{tcgaca}}$
 $\overleftrightarrow{\text{tgtcga}} \overleftrightarrow{\text{tcgaca}} \text{gaggtagtaa} \overleftrightarrow{\text{tgtcga}} \overleftrightarrow{\text{tcgaca}}$

Сравнение символьных последовательностей



- **Сложностное разложение** S_1 по S_2 : покрытие S_1 фрагментами S_2 .
- На каждом шаге копируется максимальный фрагмент S_2 , совпадающий с префиксом непокрытого участка S_1
- **Сложность** $c(S_1 / S_2)$: число компонентов в разложении S_1 по S_2

Инверсионное расстояние $d_I(\pi, \sigma)$ между последовательностями π и σ определяется минимальным числом инверсий, переводящих одну из них в другую

$$\left(\begin{array}{cccccccc} 1 & 2 & \dots & i-1 & i & i+1 & \dots & j-1 & j & j+1 & \dots & N \\ & & & & \longleftarrow & & & & \longrightarrow & & & \\ 1 & 2 & \dots & i-1 & j & j-1 & \dots & i+1 & i & j+1 & \dots & N \end{array} \right)$$

Точки разрыва для перестановок π и σ :

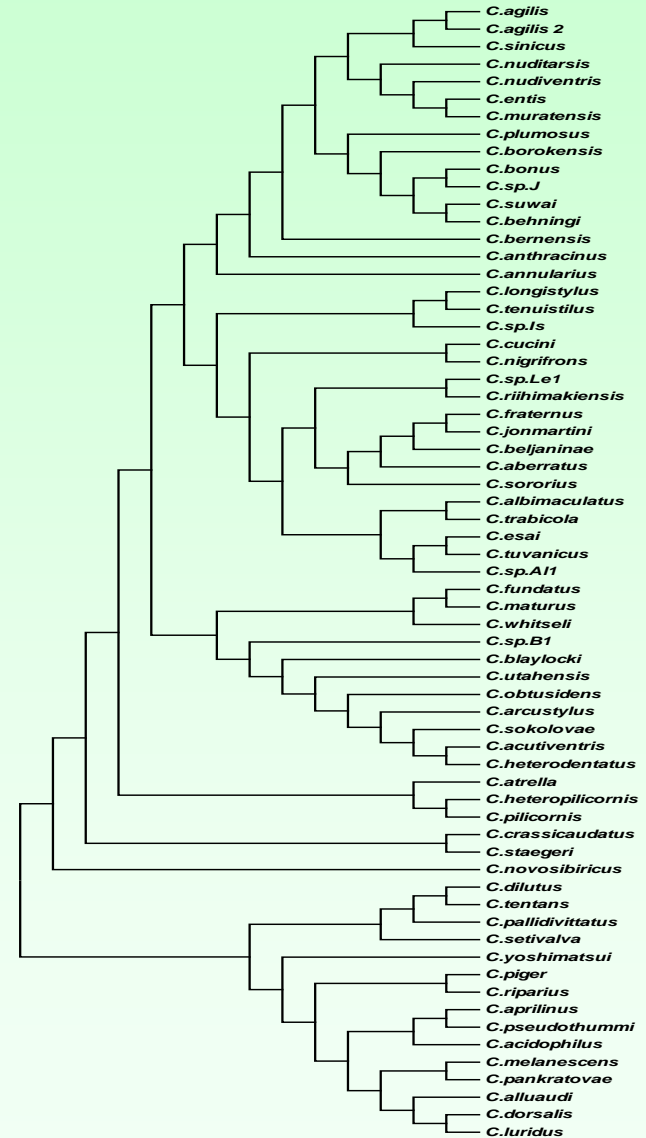
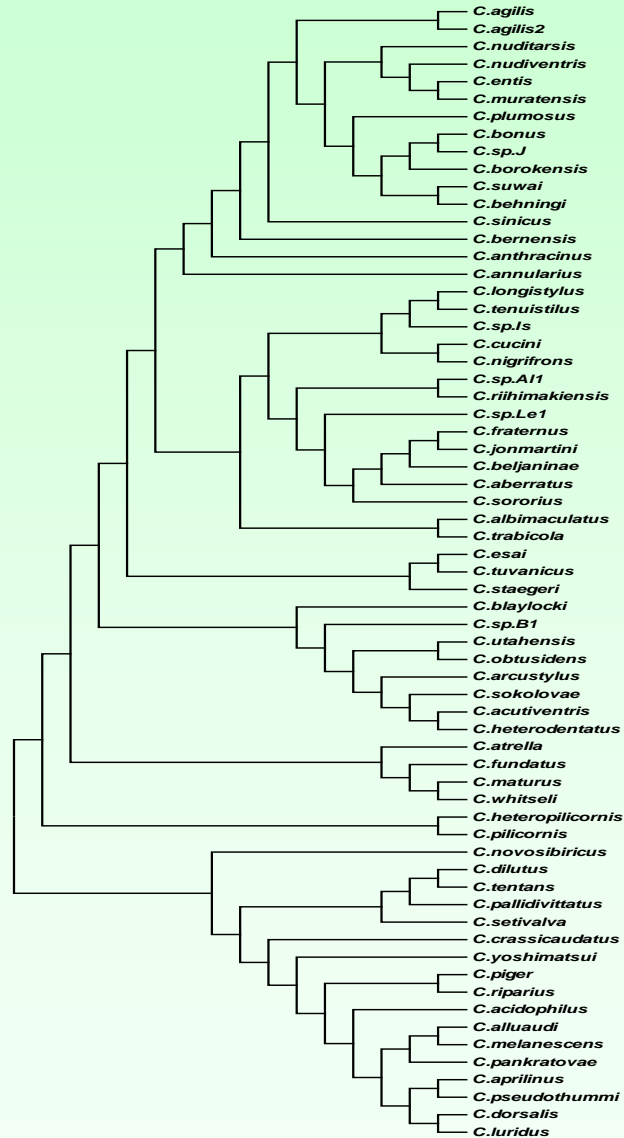
$\pi_i = a, \pi_{i+1} = b$, в σ нет биграмм ab и ba

$\pi = 0 | 6 4 | 1 8 5 | 3 | 2 9 7 | 10$

$$r(\pi, \sigma) = 5$$

$\sigma = 0 | 5 8 1 | 2 9 7 | 6 4 | 3 | 10$

Phylogenetic trees for 65 species of the genus *Chironomus*



Сравнение «близких» последовательностей

$K = \{P_1, P_2, \dots, P_k\}$ – множество текстов, разбитое на k классов

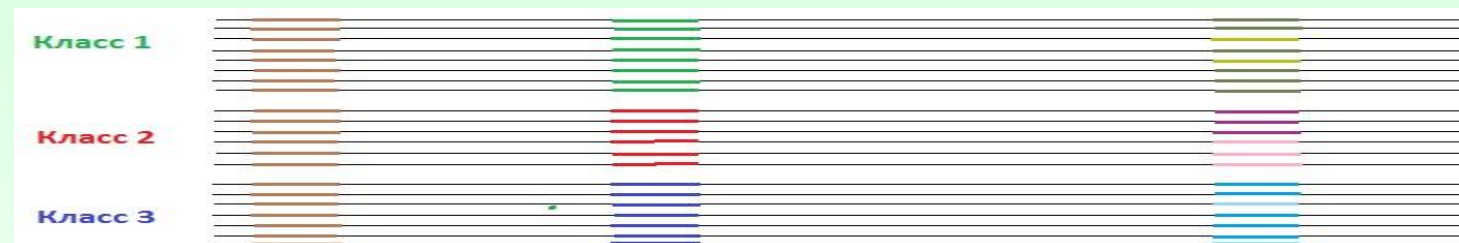
$P_i = \{T_{i1}, T_{i2}, \dots, T_{imi}\}$ ($1 \leq i \leq k$), m_i – количество текстов в i -м классе.

$\Phi_L(P_i / K)$ – множество «контрастных» L -грамм, присутствующих в каждом тексте i -го класса и отсутствующих во всех остальных текстах.

Класс 1 ($m_1 = 80$): $L_{\max}(\Phi(P_1 / K)) = 17$; $L_{\min}(\Phi(P_1 / K)) = 6$

Класс 2 ($m_2 = 46$): $L_{\max}(\Phi(P_2 / K)) = 35$; $L_{\min}(\Phi(P_2 / K)) = 7$

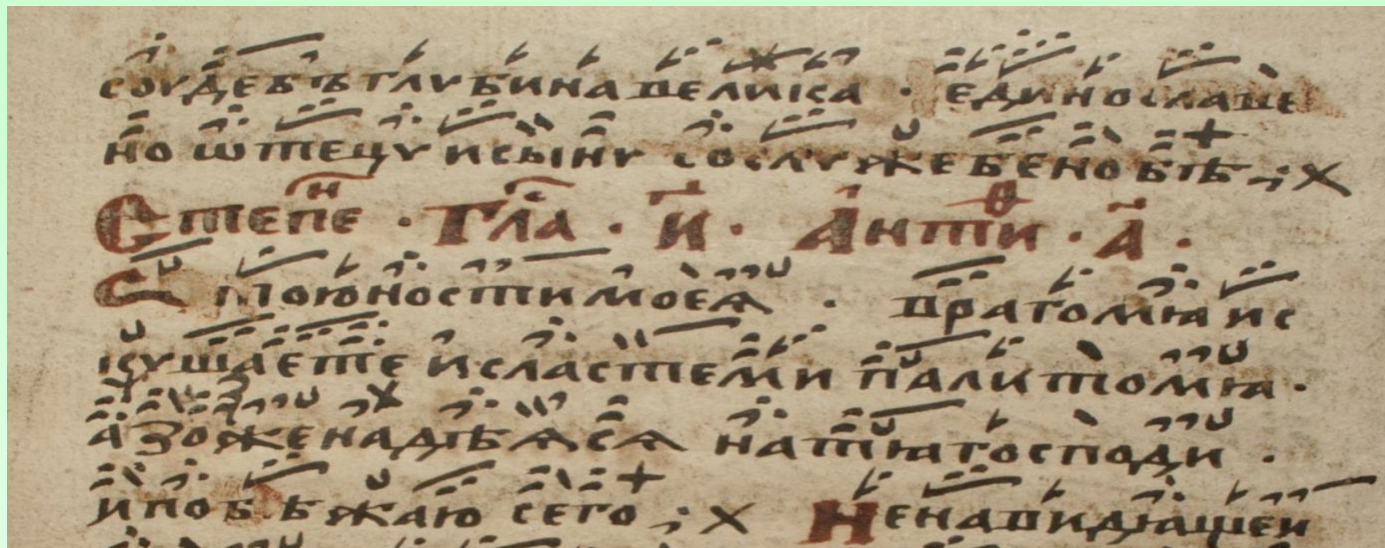
- **L -граммно-позиционный «срез»**



Позиция *pos* «классифицирующая», если множества L -грамм в текстах разных классов не пересекаются

- Периодичности
- Фрактальные структуры

Анализ параллельных текстов



Варианты распева крюка светлого в разных гласах

	473	525	442	506	443	334	321	612
R1 =	408	508	422	484	408	319	298	600
	d - 1	e - 3	c - 1	d - 4		c - 2	c - 1	c - 2
	e - 26	e - 3		e - 6		d - 11	d - 2	e - 2
	f - 70	f - 103	f - 33	f - 100	f - 52	e - 28	e - 71	e - 2
	g - 211	g - 191	g - 112	g - 317	g - 233	f - 177	f - 34	f - 189
	a - 100	a - 206	a - 234	a - 57	a - 116	g - 90	g - 171	g - 260
		b - 5	b - 17		b - 3	a - 11	a - 19	a - 142
			C - 25		C - 2			b - 5
					D - 2			
R2 =	53	2			16	1		
I1 (R2) = 2+	g - 53				C - 16			
I2 (R2) = 1-		f - 2				f - 1		
R3 =	11	14	17	11	13	10	17	8
				d - 1			d - 1	d - 2
	f - 4	f - 1		f - 3	f - 3		f - 5	f - 1
	g - 7	g - 2	g - 2	g - 7	g - 5		g - 11	g - 1
		a - 11	a - 15		a - 1	a - 6		a - 4
					b - 1			
					C - 3			
R4 =	1		3	9	5	4	4	3
						d - 2		
	g - 1		f - 1	f - 5	f - 4	f - 1	f - 3	f - 2
			a - 2	g - 4	a - 1	a - 1	g - 1	a - 1
R5 =		1						
I (R5) = 1-		a - 1						
R6 =				g - 2			g - 2	
R7 =					a - 1			
R8 =								1
I (R8) = 2-								g - 1

Вариативность. Сравнение одноименных песнопений в разных певческих книгах

ва из те бе бо я ко

ва из те бе бо я ко

знамя	нота	
+	+	1218
+	-	118
-	+	192
-	-	127

> Sol_619_647; glas1; pesn 12 -> Sol_618_644; glas1; pesn 12

```
(r0201-d2по) (r0211-d4H4по) (-0511-c4d4жде) (s0121-e2стве) (n0201-d2я)
(r0201-d2по) (r0211-d4c4рож) (-0511-c4d4де) (s0111-e2стве) (r0201-e2я)
+++ +-- ++- -++ +++
(r0201-d2ко) (r0201-d2пре) (r0201-d2жде) (r0201-d2рож) (g0402-c2де)
(r0201-e2ко) (r0201-e2пре) (r0201-e2жде) (n0201-d2рож) (g0401-c2де)
++- +-+ +-+ +-+ +-+
(n0111-d2ства) (r0211-d4H4пре) (-0511-c4d4бы) (s0121-e2ла) (g0302-c4H4е)
(n0201-d2ства) (r0211-d4H4пре) (-0511-c4d4бы) (s0111-e2ла) (g0302-c4H4е)
--++ +++ +++ +-+ +++
```

73% совпадений для 21 пары одноименных песнопений (глас 1)

Точность дешифровки коррелирована с оценкой вариативности.

Алгоритмы дешифровки знаменных песнопений.

- Нулевое приближение.
- 3-граммный подход (учет минимального контекста)
- Покрытие длинными фрагментами
- Внутригласовые инварианты (ВИ) и квазиинварианты (КВИ)
- Покрытие попевками в сочетании с ВИ/КВИ
- Связные пути многоуровневого графа

$k = n_+ / N$, где n_+ – число «правильно» интерпретированных знамен,
 N – суммарная длина песнопений.

Подход	Гласы							
	1	2	3	4	5	6	7	8
1	0.59	0.49	0.5	0.5	0.5	0.42	0.54	0.43
2	0.64	0.58	0.62	0.6	0.58	0.52	0.71	0.48
3	0.64	0.63	0.63	0.63	0.63	0.52	0.78	0.5
4	0.68	0.63	0.65	0.69	0.67	0.56	0.79	0.52
5	0.69	0.64	0.66	0.7	0.67	0.57	0.83	0.51
6	0.68	0.63	0.67	0.57	0.61	0.55	0.76	0.53

Исходный текст и дешифровка

У те ню ем у тре ню ю гла бо ку, и вме сто



ми ра песнь при не семь вла ды це, и Хри ста



у зримъ пра вды со лнце, все мь жизнь



во си я ю ща.



Ирмологий, глас 1, песн. 34.

45 знамен

Точность реконструкции $24/45 = 0.53$

3 ошибки в восстановлении ритма



**Алгоритмические аспекты анализа
символьных последовательностей**

Благодарю за внимание!