

Анализ символьных последовательности различной языковой природы

Мирошниченко Любовь
Александровна

Институт математики СО РАН

luba@math.nsc.ru

Объект исследования: символьные последовательности различной языковой природы.

Σ – непустое конечное множество символов (алфавит);

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

Примеры:

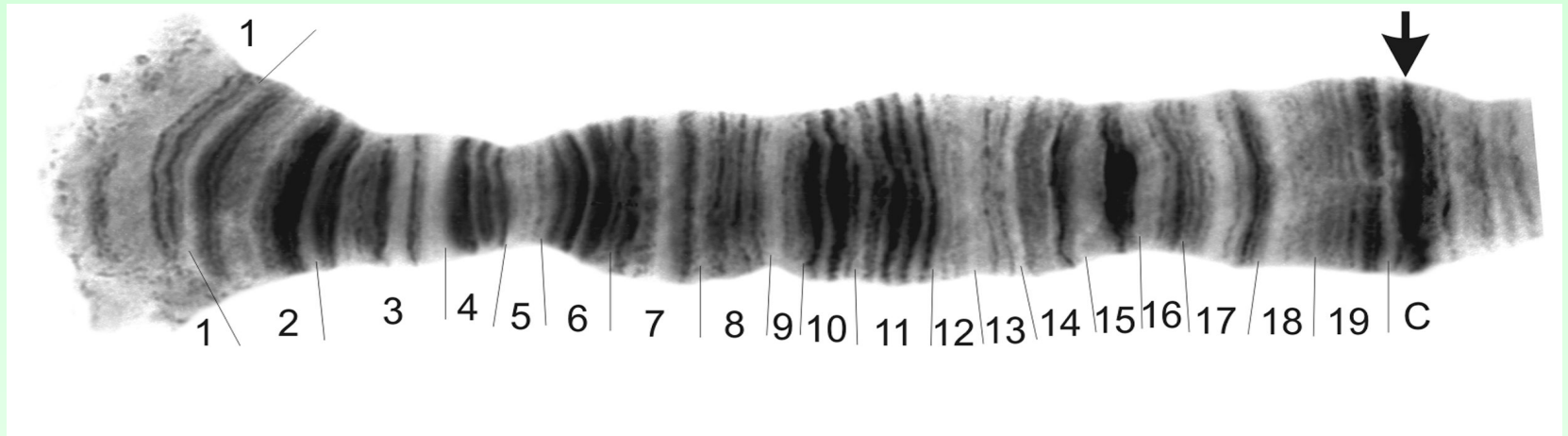
- слова, предложения, ..., тексты естественного языка;
- музыкальные тексты (песенные мелодии);
- древнерусские церковные песнопения;
- тексты программ;
- ДНК, РНК ($|\Sigma| = 4$); аминокислотные послед. ($|\Sigma| = 20$);
- порядки генов; порядки дисков политенных хромосом;
- последовательность действий;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

ДНК и аминокислотные последовательности

- ДНК: $\Sigma = \{A, C, G, T\}$, РНК: $\Sigma = \{A, C, G, U\}$;
- Белки практически всех живых организмов построены из аминокислот всего 20 видов.

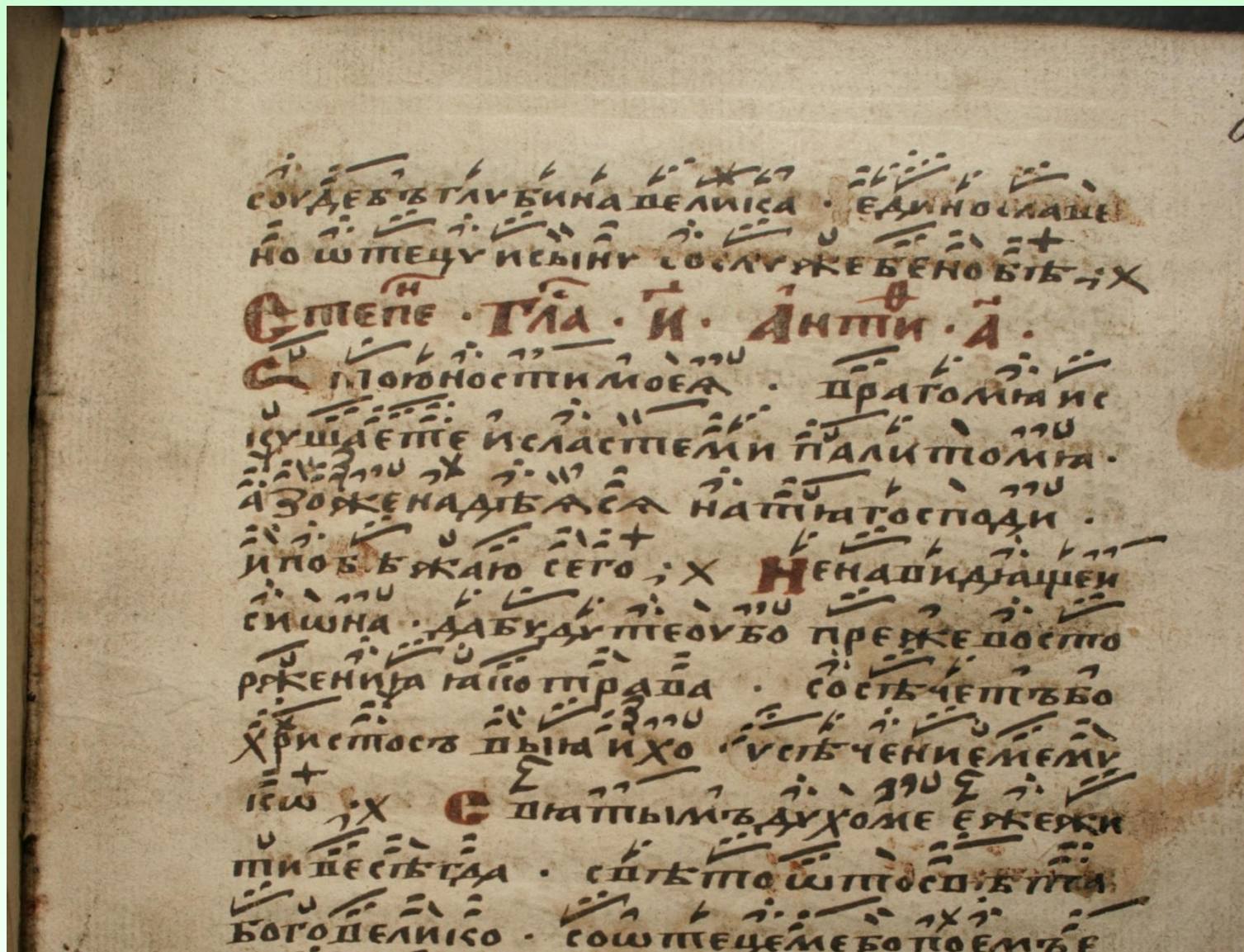
		ВТОРАЯ БУКВА					
		U	C	A	G		
ПЕРВАЯ БУКВА	U	UUU } Фенил-аланин F UUC } UUA } Лейцин L UUG }	UCU } UCC } Серин S UCA } UCG }	UAU } Тирозин Y UAC } UAA } Стоп-кодон UAG } Стоп-кодон	UGU } Цистеин C UGC } UGA } Стоп-кодон UGG } Триптофан W	ТРЕТЬЯ БУКВА	U
	C	CUU } Лейцин L CUC } CUA } CUG }	CCU } CCC } Пролин P CCA } CCG }	CAU } Гистидин H CAC } CAA } Глутамин Q CAG }	CGU } CGC } Аргинин R CGA } CGG }		C
	A	AUU } Изолейцин I AUC } AUA } AUG } Метионин M старт-кодон	ACU } ACC } Треонин T ACA } ACG }	AAU } Аспарагин N AAC } AAA } Лизин K AAG }	AGU } Серин S AGC } AGA } Аргинин R AGG }		A
	G	GUU } Валин V GUC } GUA } GUG }	GCU } GCC } Аланин A GCA } GCG }	GAU } Аспарагиновая кислота D GAC } GAA } Глутаминовая кислота E GAG }	GGU } GGC } Глицин G GGA } GGG }		G

Polytene chromosomes



Cytophotomicrograph of arm A of the species *C. piger*

Пример древнерусской церковной рукописи



СОУДЕВЪ ТЛУБИНА ДЕЛИКА • ЕДИНОСЛАВ
НО ШПЕЦУ ИСЫНУ СОСЛУЖЕБЕНОВЪ Х
С ПЕПЕ • **Г**ЛА • **И** • **А**НТИ • **А** •
С ПЛОЮНОСТИ МОЕА • ПРАГОМЪ ИС
КУШАЕТЕ И СЛАСТЕМИ ПАЛИТОМА •
А ЗОЖЕНА ДЬ АСА НАТНЪ ГОСПОДИ •
И ПОВЪЖАЮ СЕГО • Х **Н**ЕНАДИДЪЩЕИ
СИШНА • ДА БУДУТЕ ОУБО ПРЕЖЕ ПОСТО
РЖЕНИЯ ТАКО ПРАДА • СОСЪЧЕПЪ БО
ХРИСТОСЪ ПЫНЪ ИХО • УСЪЧЕНИЕМЪ
ИСО • Х **С** ПЯТЫМЪ ДУХОМЕ ЕЖЕЖИ
ТИ ВЕСЪТЪА • СЪТЪ ПО ШПОСПЪТЪА
БОГО ДЕЛИКО • СО ШПЕЦЕМЪ БО ПОЕМЪ Е


Знамена (крюки)

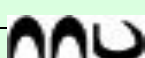
Примеры начертаний:

— юк: e2; — пал а; — чаи а: d4c4

— ла простая: f1 (e1...)

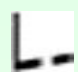
— ла поводная с облачком и оттяжкой: d4e4f2.d4


— олубчик борзый: c4d4 (d4e4, e4f4 ...)

— ла: H4H4A2

— ийца со статьей: d4e4d4.c8H4

Примеры толкований:

 — *столица с очком*: назад отшибнуть гортанью, вскочить и опуститься на голубчик или на скамейцу: e4d4 (d4c4...)

 — *сложитие*: покудрить гортанью: f8e8f4, g4f4...

Кодировка песнопений из двознаменника

Первый и шестой символ кода – степенные и указательные пометы

- **Степенные** – указывают высоту распева знамен.
- **Указательные пометы** (**Т** – тихая, **Б** – борзая...) определяют характер исполнения распева знамен.

Знамена кодируются четырехсимвольным кодом.



The image shows a musical staff with a treble clef and a key signature of one flat. The melody consists of 12 notes: G4, A4, B4, G4, A4, B4, C5, B4, A4, G4, F4, E4. Below the staff is a four-symbol code: Г, Н, Ц, Г. Below the code is a table of consonant categories and their corresponding letters.

ПРОСТОЕ СОГЛАСИЕ	МРАЧНОЕ СОГЛАСИЕ	СВѢТЛОЕ СОГЛАСИЕ	ТРЕСВѢТЛОЕ СОГЛАСИЕ
Г А Н	с d e	f g a	b С D



Длительности звуков: – 1 (целая), – 2, – 4, – 8

Н4 – четвертная нота «си» малой октавы

Пример кодировки песнопения из двознаменника

(m0401-c2Bo) (v0121-e2нми) (r0121-e2зе) (r0111-e2мле)
(r0211-e4d4и) (r1941-c4d4e2не) (p1011-d1бо) (v0901-c4e4и)
(p0302-d4c4вну) /
(-0501Td2e2ши)
(*1021-f1) (-0511-d4e4гла) (#0141-f2го) (-1601Ld4e4лы)
(-0901-d4c4мо) (-1002-d1я) (-1001-c1) (m0211-c4H4воз)
(-0511-c4d4гла) /
(v0121-e2го) (r0121-e2лю) (r0211-e4d4бо) (-0511-c4d4на)
(v0301-e2зе) (p1001-d1мли) (v0905Td2e2бо) (p0111-d2жи)
(p1861-c2d1я) (p0201-d2чю) (m0301-c2де) (-2801-H1ca.) /@

Основные задачи анализа текста

- поиск образцов;
- восстановление структуры текста: выявление повторов (периодичностей, симметрий ...);
- сравнение последовательностей: разные определения расстояний и мер близости;
- сложность текста
- сегментация, фрагментация, выделение структурных единиц...

Формальные языки и грамматики

Σ – алфавит;

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

$N = |T|$ – длина строки T ;

$T[1 : p] = t_1 t_2 \dots t_p$ – префикс слова ($1 \leq p \leq N$),

$T[k : N] = t_k t_{k+1} \dots t_N$ – суффикс ($1 \leq k \leq N$),

$T[k : p] = t_k t_{k+1} \dots t_p$ – подслово ($1 \leq k \leq p \leq N$);

e – пустая строка ($|e| = 0$);

Σ^* – множество всех слов (строк) в алфавите Σ , включая e .

Язык L над Σ – произвольное множество слов в Σ ($L \subseteq \Sigma^*$).

Конкатенация языков L_1 и L_2 есть $L_1 L_2 = \{ \alpha \beta : \alpha \in L_1, \beta \in L_2 \}$.

L^* : итерация языка L :

$$L^0 = \{\varepsilon\},$$

$$L^n = LL^{n-1} \text{ для } n \geq 1,$$

$$L^* = \bigcup_{n \geq 0} L^n$$

Порождающей грамматикой называется четверка

$G = (\Sigma, N, P, S)$, где

Σ – алфавит терминальных символов, из которых
составляются «слова» языка ($L(G) \subseteq \Sigma^*$);

N – алфавит нетерминальных символов (или переменных);

$\Sigma \cap N = \emptyset$;

P – конечное множество правил вывода вида $\alpha \rightarrow \beta$, где $\alpha \in (N \cup \Sigma)^* N (N \cup \Sigma)^*$, $\beta \in (N \cup \Sigma)^*$;

S – выделенный символ из N , называемый начальным (или исходным).

Формальная грамматика позволяет получить все цепочки данного языка и только их. Формальные грамматики были введены Хомским (1956г). Им же определена классификация грамматик в зависимости вида применяемых правил вывода (иерархия Хомского).

Иерархия Хомского

Пусть $G = (\Sigma, N, P, S)$ – грамматика. G называется:

- **праволинейной**, если каждое правило из P имеет вид

$$A \rightarrow \alpha B, \text{ где } A, B \in N \text{ и } \alpha \in \Sigma^*;$$

- **праволинейная** грамматика называется **регулярной** (или автоматной), если все ее правила имеют вид

$$A \rightarrow aB \text{ или } A \rightarrow a, \text{ где } A, B \in N \text{ и } a \in \Sigma$$

- **контекстно- свободной**, если каждое правило из P имеет вид

$$A \rightarrow \alpha, \text{ где } A \in N \text{ и } \alpha \in (N \cup \Sigma)^*$$

- **контекстно- зависимой** (или неукорачивающейся), если каждое правило из P имеет вид $\alpha \rightarrow \beta$, где $\alpha, \beta \in (N \cup \Sigma)^*$ и $|\alpha| \leq |\beta|$.

- Грамматика, на которую не накладывается ни одно из указанных ограничений, называется **грамматикой составляющих**.

Языки называются праволинейными, КС или КЗ в зависимости от того, какой грамматикой он порожден.

Пример формальной грамматики

Пусть $G = (\{a,b,c\}, \{A,B,S\}, P, S)$,

где правила вывода P имеют вид:

$S \rightarrow AB, A \rightarrow a, A \rightarrow ac, B \rightarrow b, B \rightarrow cb.$

Данная грамматика позволяет получить всего 4 вывода терминальных строк:

(1) $S \rightarrow AB \rightarrow aB \rightarrow ab$

(2) $S \rightarrow AB \rightarrow aB \rightarrow acb$

(3) $S \rightarrow AB \rightarrow acB \rightarrow acb$

(4) $S \rightarrow AB \rightarrow acB \rightarrow accb$

$L(G) = \{ab, acb, accb\}.$

Для строки acb имеются два разных вывода.

Пример. Арифметические выражения

- $\Sigma = \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9, +, -, *, /, (,)\}$
- $N = \{\text{ФОРМУЛА, ЗНАК, ЧИСЛО, ЦИФРА}\};$
- $S = \text{ФОРМУЛА}$
- Правила:
 1. ФОРМУЛА \rightarrow ФОРМУЛА ЗНАК ФОРМУЛА
 2. ФОРМУЛА \rightarrow ЧИСЛО
 3. ФОРМУЛА \rightarrow (ФОРМУЛА)
 4. ЗНАК $\rightarrow + \mid - \mid * \mid /$
 5. ЧИСЛО \rightarrow ЦИФРА
 6. ЧИСЛО \rightarrow ЧИСЛО ЦИФРА
 7. ЦИФРА $\rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

Пример вывода $(12 + 5) * 3$

ФОРМУЛА \Rightarrow_1 ФОРМУЛА ЗНАК ФОРМУЛА \Rightarrow_4 ФОРМУЛА *
ФОРМУЛА \Rightarrow_1 ФОРМУЛА * ЧИСЛО \Rightarrow_5 ФОРМУЛА * ЦИФРА \Rightarrow_7
ФОРМУЛА * 3 \Rightarrow_2 (ФОРМУЛА) * 3 \Rightarrow_5 (ФОРМУЛА ЗНАК
ФОРМУЛА) * 3 \Rightarrow_3 (ФОРМУЛА + ФОРМУЛА) * 3 \Rightarrow_1 (ФОРМУЛА
+ ЧИСЛО) * 3 \Rightarrow_4 (ФОРМУЛА + ЦИФРА) * 3 \Rightarrow_2 (ФОРМУЛА + 5) * 3
 \Rightarrow_5 (ЧИСЛО + 5) * 3 \Rightarrow_6 (ЧИСЛО ЦИФРА + 5) * 3
 \Rightarrow_2 (ЦИФРА ЦИФРА + 5) * 3 \Rightarrow_7 (1 ЦИФРА + 5) * 3 \Rightarrow_7 (1 2 + 5) * 3

Конечные автоматы – средство распознавания

Детерминированный конечный автомат – это пятерка

$M = (S, \Sigma, \delta, s_0, F)$, где

S – конечное множество состояний;

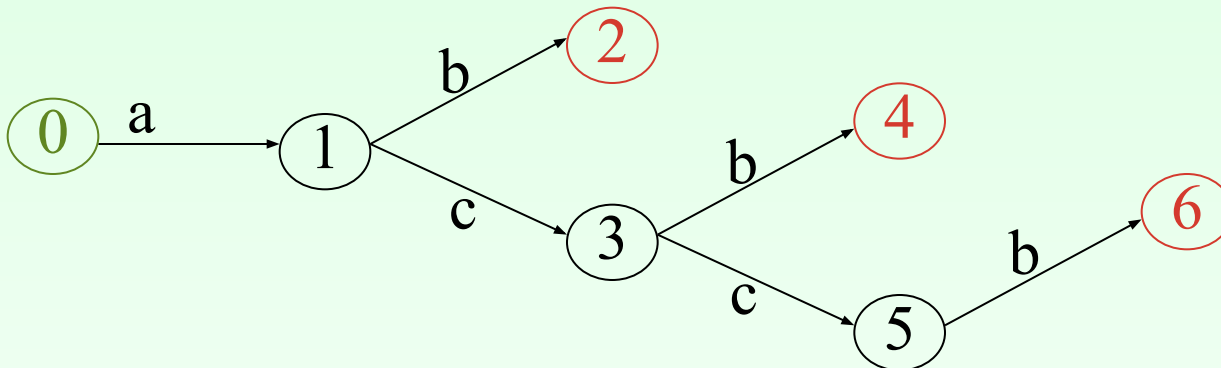
Σ – алфавит;

$\delta : S \times \Sigma \rightarrow S$ – функция переходов;

$s_0 \in S$ – выделенное начальное состояние;

$F \subseteq S$ – множество заключительных состояний;

ДКА, допускающий $\{ab, acb, acsb\}$.



Формальные последовательности

Последовательность Туэ - Морса

Способы задания

1. итерации морфизмов.

$\Sigma = \{a_1 \dots a_q\}$ $\phi : \Sigma^* \rightarrow \Sigma^*$ – морфизм, если $\phi(XY) = \phi(X)\phi(Y) \quad \forall$ слов X и Y .

$\phi = \{0 \rightarrow 01, 1 \rightarrow 10\}$.

$X_0 = 0, X_1 = 01, X_2 = 0110, X_3 = 01101001, X_4 = 0110100110010110 \dots$

2. $X[i] = 0$, если число единиц в двоичной записи числа i чётно,

$X[i] = 1$, в противном случае.

3. Итеративный способ: $X[0] = 0, X[2i] = X[i], X[2i+1] = ((X[i] + 1) \bmod 2)$

Свойства последовательности Туэ-Морса:

1. Отсутствуют подслова вида VVV .

2. $X_{2n} = X_n X_n^R$: слово, полученное на чётном шаге является палиндромом.

Чíсла Фибона́ччи — 0, 1, 1, 2, 3, 5, 8, 13, 21, 34, 55, 89, 144, 233, 377, 610, ...

Последовательность Фибоначчи

$X_0 = 0, X_1 = 01, X_n = X_{n-1} X_{n-2}$

$X_2 = 01.0, X_3 = 010.01, X_4 = 01001.010, X_5 = 01001010.01001$

Морфизм: $\phi = \{0 \rightarrow 01, 1 \rightarrow 0\}$