

**СТРУКТУРНЫЕ АНАЛОГИИ В
СИМВОЛЬНЫХ
ПОСЛЕДОВАТЕЛЬНОСТЯХ
РАЗЛИЧНОЙ ЯЗЫКОВОЙ ПРИРОДЫ**

*Гусев В.Д., Мирошниченко Л.А.,
Саломатина Н.В.*

Новосибирск, Институт математики СО РАН

Объект исследования

Σ – непустое конечное множество символов (алфавит);

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

Примеры:

- ДНК, РНК ($|\Sigma| = 4$); аминокислотные послед. ($|\Sigma| = 20$);
- порядки генов ($|\Sigma| \sim 10^4$); порядки дисков политенных хромосом;
- музыкальные тексты (песенные мелодии);
- знаменные песнопения;
- тексты программ;
- слова, предложения, ..., тексты естественного языка;
- последовательности действий животных;
- формальные последовательности.

Повтор – структурообразующий элемент

Повтор (в широком смысле) — пара фрагментов текста, совпадающих с точностью до переименования элементов алфавита и (или) изменения направления считывания.

Типы повторов:

- разнесенные ... **AGTTC** ... **AGTTC** ...
- тандемные ... **AGTTCAGTTC**...
- симметричные: ... **AGTTC** ... **CTTGA** ...
- прямые комплементарные: ... **AGTTC** ... **TCAAG** ...
- инвертированные: ... **AGTTC** ... **GAACT** ...
- с произвольным переименованием элементов:

if $j > 1$ *then* $k := k + 1$ *else* $h := h * 2$;

if $pr2 > 1$ *then* $tam := tam + 1$ *else* $tut := tut * 2$

Сложность последовательности

- I. определения учитывающие частоту встречаемости символов или относительно коротких l -грамм (энтропийные меры, SIMPLE...)
- II. разнообразие l -граммного состава (комбинаторная сложность, лингвистическая ...)
- III. определения, основанные на учете длинных повторов (грамматическая сложность, Лемпеле-Зив-подобные меры...).

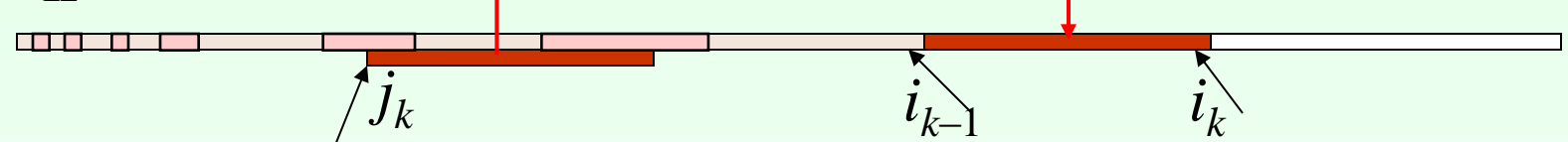
Представление текста в виде конкатенации повторяющихся фрагментов

Сложность последовательности S по Лемпелю и Зиву измеряется **числом шагов порождающего ее процесса**. Каждый шаг: **копирование** максимально длинной "заготовки" из уже синтезированной части текста или **генерация** символа.

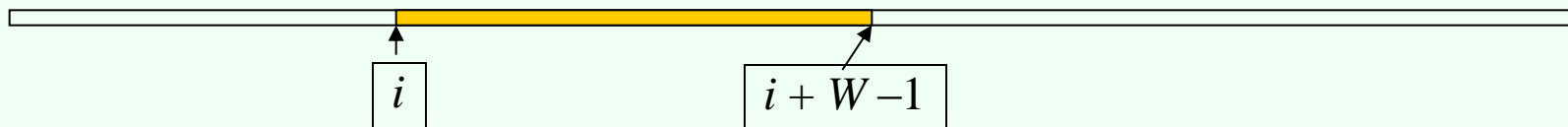
Сложностное разложение S есть конкатенация фрагментов

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

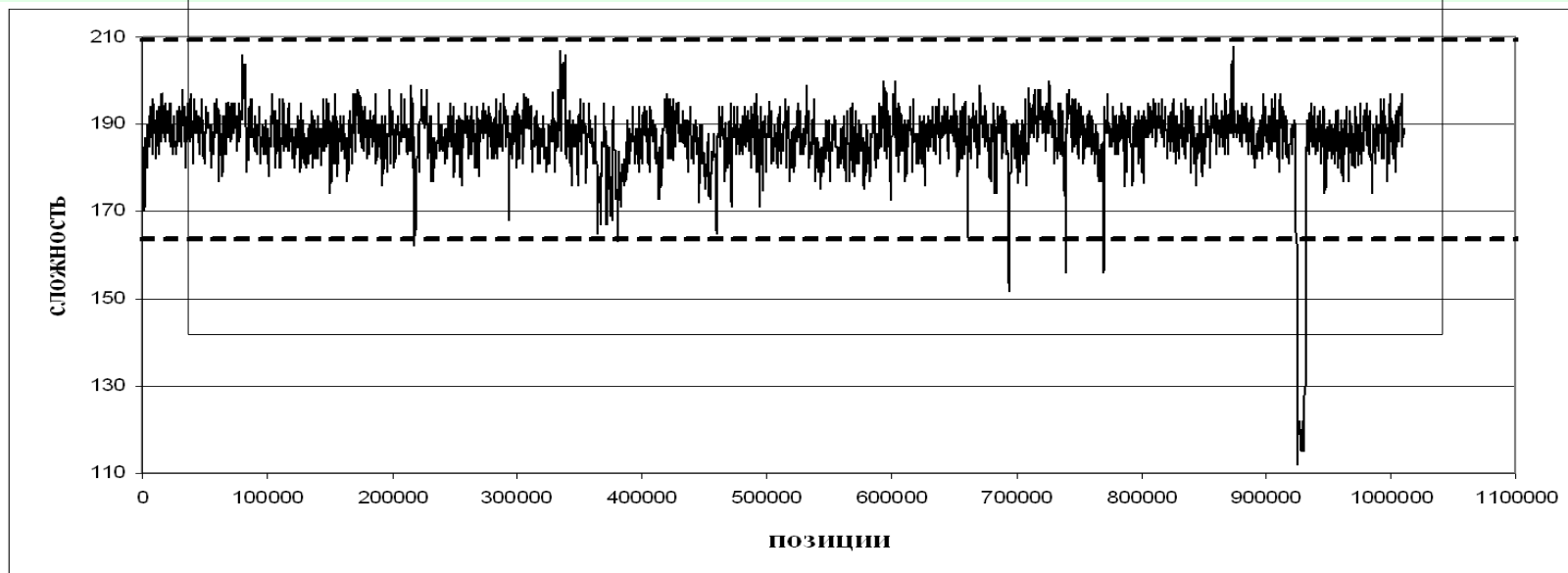
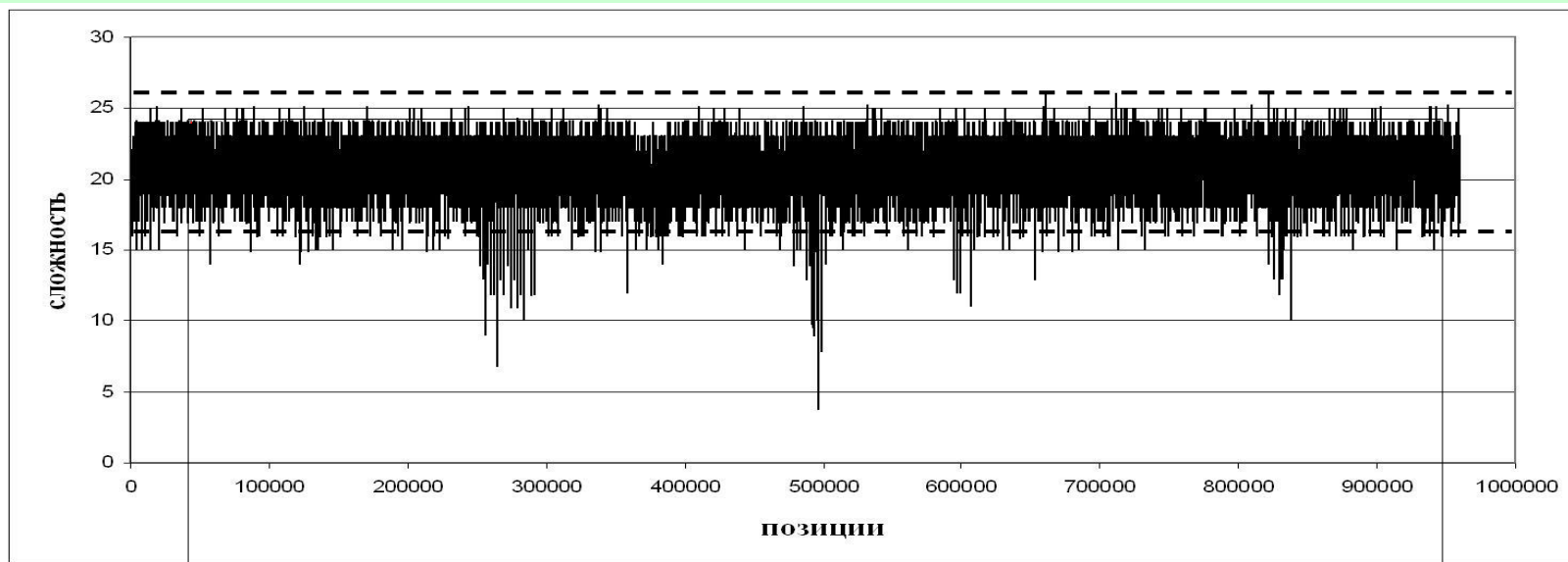
$C_{LZ}(S) = m$ — сложность (число шагов процесса).



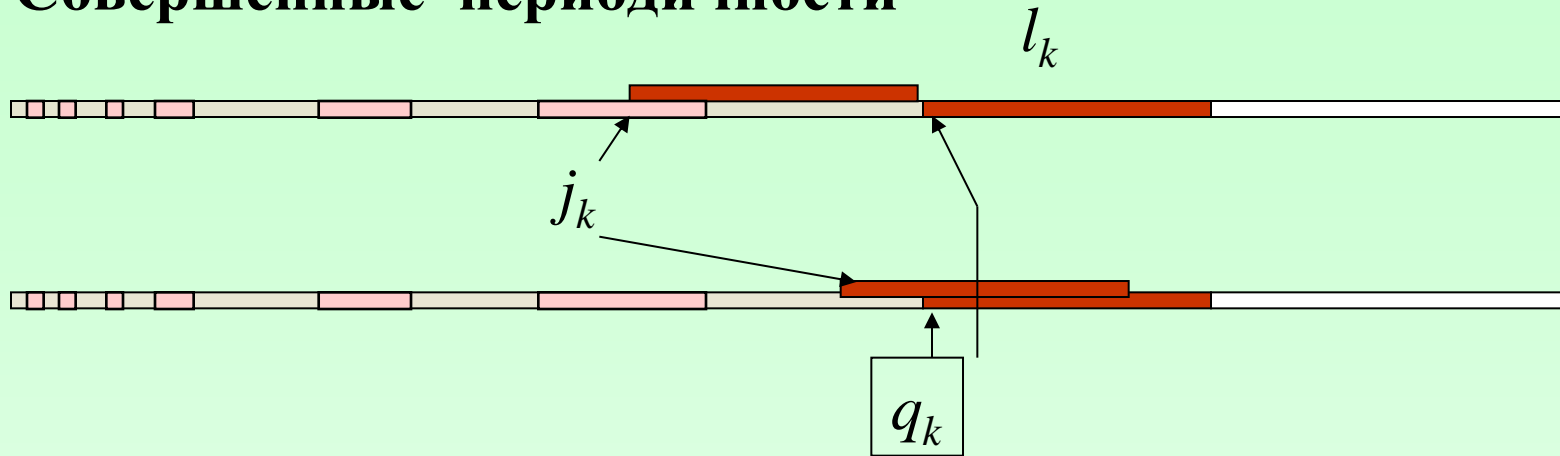
$$c_i = C(S[i : i + W - 1])$$



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



Совершенные периодичности



Если $j_k + l_k \geq q_k$, имеет место периодичность.

Длина периода $p = q_k - j_k$; кратность $\geq \lfloor l_k / p \rfloor + 1$.

$$H(S) = \overbrace{c \cdot a} \cdot \overbrace{a \cdot c \cdot ca} \cdot \overbrace{t \cdot g \cdot at \cdot (ccatgat)^4} \cdot at;$$

l_k	1	1	1	1	2	1	1	2	28	2	$c(S)=10$
j_k	0	0	2	1	1	0	0	6	4	37	

Совершенные периодичности

- ДНК-последовательности (p от 1 до 10^4)
- Сказки, тексты песен.

Дождливым вечером,
Вечером, вечером ...
Пора в путь-дорогу,
Дорогу дальнюю, дальнюю,
Дальнюю идём.,

В чаще чаще меньше пищи, значит
в чаще чаще чище. (Б. Заходер)

- Знаменные песнопения

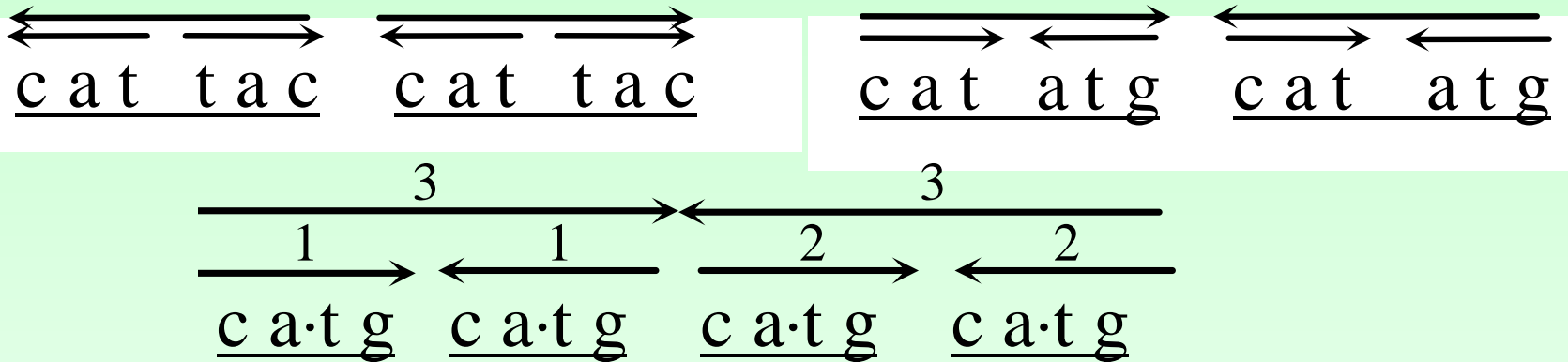
(تأه)²

(ل ه ل)²

(ل تأه ل)²

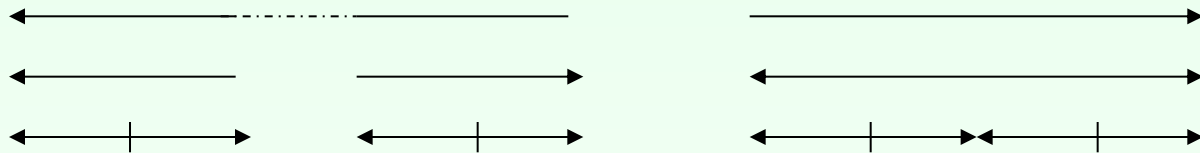
Локальные фракталы в ДНК-последовательностях:

фрагменты текста, обладающие свойством самоподобия.



Фракталоподобные структуры – структуры, в которых позиционно сближены палиндромы

agagaagactagattcaagatcaga



Комбинации периодичностей (компаунды)

- В темном лесе, в темном лесе,
- В темном лесе, в темном лесе,
- За лесью, за лесью,
- Распашуль я распашуль я
- Распашуль я распашуль я
- Пашенку, пашенку.
- Я посею, я посею,
- Я посею, я посею лен,
- Конопель, лен зеленой.

LLLL. ((ä)) ((ä))

серия стоиц, стопица с очком (змийца со статьей, голубчик борзый мальй, палка)², статья мрачная

Периодичности со сложной структурой

MLTR (Multiple Length Tandem Repeat):

$(Xx^n)^m$, где X и x – фрагменты текста, $X \neq x$, $n > 1$, $m > 1$

$((\text{gacctttgg}) - \text{ac-ac-ac-ac})(\text{gacctttgg}) - \text{ac-ac-ac-ac}$

VLTR (Variable Length Tandem Repeat) или Nested Tandem Repeat

$x^{n_1}Xx^{n_2}X\dots Xx^{n_l}$, где $l > 2$, $n_i \geq 1$ для $i = 1, \dots, l - 1$ и по крайней мере одно $n_i \geq 2$

MPTR – Multi-Periodic Tandem Repeat [Hauth]

x_1	x_2	x_1	x_3		
S =	(cagta)	(cagca)	(cagta)	(caaca)	$d(S, (x_1)^{12}) = 9$
	(cagta)	(cagca)	(cagta)	(caaca)	$d(S, (x_1x_2)^6) = 3$
	(cagta)	(cagca)	(cagta)	(caaca)	$d(S, (x_1x_2x_1x_3)^3) = 0$

- Hauth, A.M. and Joseph, D.A. *Beyond Tandem Repeats: Complex Pattern Structures and Distant Regions of Similarity*. Bioinformatics. 2002 Jul;18 Suppl 1:S31-7.
- Matroud A.A, Hendy M. D. , Tuffley C. P. *NTRFinder: a software tool to find nested tandem repeats*, Nucleic Acids Res., Feb 2012; 40(3): e17

Кумулятивные структуры

$A + BA + CBA \dots$ или $A + AB + ABC \dots$ или $\dots ABC + AB + A$

A =	Вот дом, который построил Джек.
BA =	A это пшеница, которая в тёмном чулане хранится в доме, который построил Джек.
CBA =	A это весёлая птица – синица, которая часто ворует пшеницу, которая в тёмном чулане хранится в доме, который построил Джек

Репка, колобок ...

Вариационность типа прорастания (музыкальные тексты)

Иерархическая дупликация (ДНК):

```

21479 g a t t t g - - g g g a c g a a a a c c a c a t c g t c g t t t
      | | | |           | | | | | | | | | | | | | | | | | | | | | | |
21509 g a t t a c - - g g g a c g a a a a c a a c a - g g c a g t t t c
      | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
21539 g a t t a c - - g g c a c c a a a t c g a c g
      | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
21560 a a t a a c a c g g g g
  
```

Регулярные повторы

gttttgggggttgtagaattatthttgttagtaaaac aatgataaataacgcttaacttgcttact
gttttgggggttgtagaattatthttgttagtaaaac cctataaacaatcaggattatatgtacta
gttttgggggttgtagaattatthttgttagtaaaac ttagtcaagattthttaataccagggtgca
gttttgggggttgtagaattatthttgttagtaaaac tccatathtttccttactattactatgct
gttttgggggttgtagaattatthttgttagtaaaac acgaththtaaaaattatgatataataaac
gttttgggggttgtagaattatthttgttagtaaaac tagaatctctthtaattcccaccaagct
gttttgggggttgtagaattatthttgttagtaaaac cthththtaaththtattataactthgt
gttttgggggttgtagaattatthttgttagtaaaac aattgcatcattaacgthtaagacgthtact
gttttgggggttgtagaattatthttgttagtaaaac aatcaaacaactcgctthctaaatcatcaa
gttttgggggttgtagaattatthttgttagtaaaac thtgagcataatggcgctthtgagththtag
gttttgggggttgtagaattatthttgttagtaaaac aatgcagaththaaagathcaggaacgath
gttttgggggttgtagaattatthttgttagtaaaac attagccccacaattatattaacctccct
геном "*Mycoplasma synoviae* 53" (ID AE017245), поз. 690229

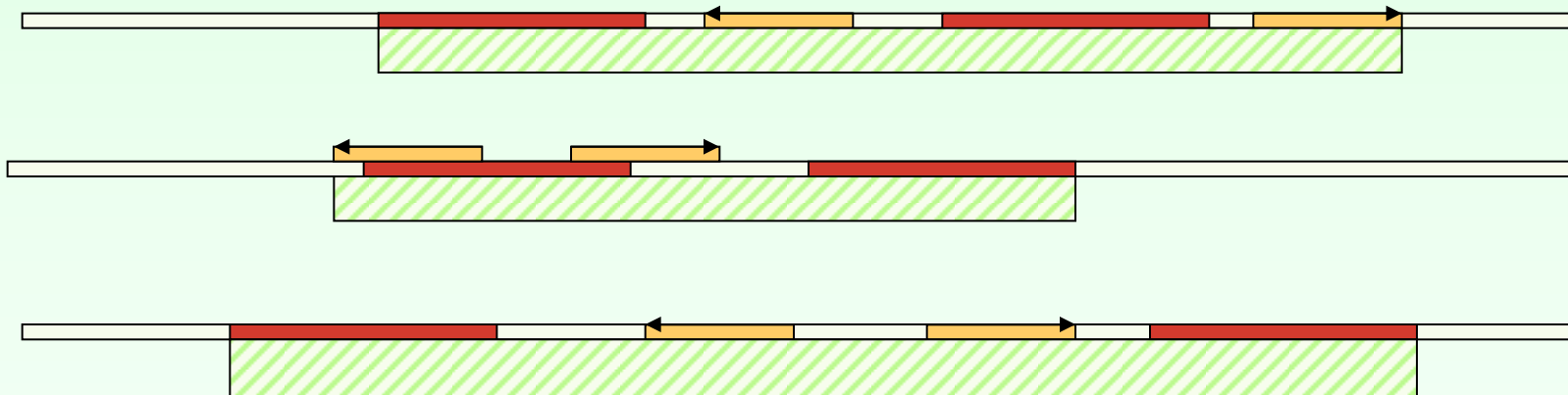
CRISPRs (*Clustered Regularly Interspaced Short Palindromic Repeats*)

Тексты песен: куплет – припев...

Комбинированные структуры

Комбинации разнотипных нерасширяемые повторов
(прямой + инвертированный...) с ограничением на длины
повторов: $l \geq R$ и длины «гэпов» $d \leq r$.

ДНК-последовательности + последовательности
поведения животных



Комбинированные структуры в геноме ВКЭ

Dalnegorsk (и другие вирулентные)

• 769: cgctcagggagatctcacaggaaggggacacaaatgggttagaaggggactcat (49)
 Xs = ctcaggg gggactc = X

Y = gaaggggac

Y = gaaggggac

212: A Q G D L T G R G H K W L E G D S

• 220: tagagtccaaatgccaaatggactcgt (23)

Xs = gagtcca tggactc = X

Y = ccaaatg

Y = ccaaatg

29: R V Q M P N G L

Primorye-183 (и другие инаппарантные)

• 4098: agaagggactgacttggattgtcccttggccgggctacttggagg (42)

Xs = aagggac gtccctt = X

Y = acttggga

Y = acttggga

1322: K G L T W I V P L A G L L G G

Заключение

Каждая «структурная единица» требует разработки «своего» алгоритма

**СТРУКТУРНЫЕ АНАЛОГИИ В
СИМВОЛЬНЫХ
ПОСЛЕДОВАТЕЛЬНОСТЯХ РАЗЛИЧНОЙ
ЯЗЫКОВОЙ ПРИРОДЫ**

Благодарю за внимание!