

# Применение сжатых индексов для полнотекстового поиска

Студент: Вадим Соколов  
Научный руководитель:  
Алексей Неганов

# 1

## Цели и Задачи

---



**Сравнение текущих подходов**



**Suffix array & Radix tree**



**Реализация Compressed suffix array**



**Поиск подстроки**

## Inverted Index

- Для текста с разделением на слова
- 5-10% от оригинального текста
- ~50% с учетом позиции
- Невозможно восстановить оригинал

"it is what it is"

"what is it"

"it is a banana"

"a": {2}

"banana": {2}

"is": {0, 1, 2}

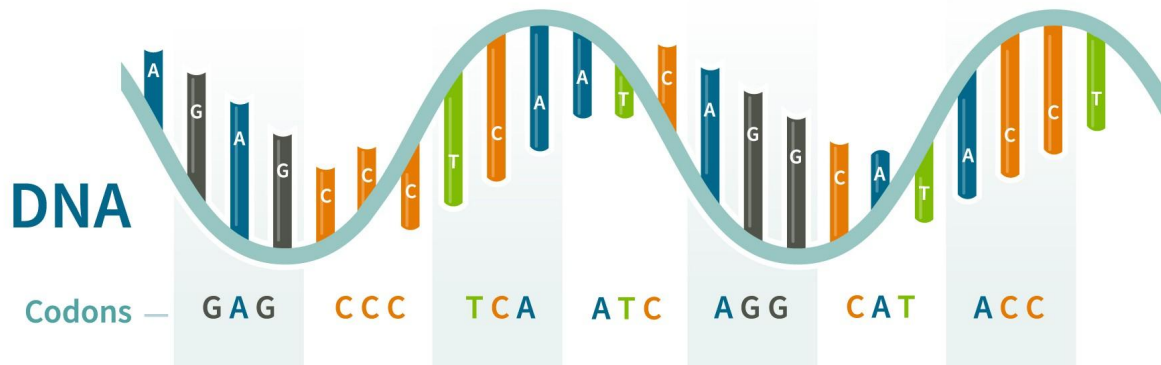
"it": {0, 1, 2}

"what": {0, 1}

# 3 Проблемы Inverted Index

---

## ● ДНК и белковые структуры



# 3 Проблемы Inverted Index

---



## Восточные языки

俄罗斯人前鋒力于发展  
科学和艺术是美妙的仅

- Китайский

これが未来なので私たちは科学と芸  
術を緒に開発する必要があります

- Японский



# Проблемы Inverted Index

---



Fuzzy search

## Search results

---

This wiki is using a new search engine. ([Learn more](#))

Search

[Content pages](#) [Multimedia](#) [Translations](#) [Everything](#) [Advanced](#)

Did you mean: *andré emotions*

## 4

# Suffix Array



Поиск подстроки



Занимает большой размер



$I = n \log n + n \log \sigma$

Idx	Suffixes	SA-Idx	Idx	Sorted Suffix
0	BANANA\$	0	6	\$
1	ANANA\$	1	5	A\$
2	NANA\$	2	3	ANA\$
3	ANA\$	3	1	ANANA\$
4	NA\$	4	0	BANANA\$
5	A\$	5	4	NA\$
6	\$	6	2	NANA\$

Suffix Array [6, 5, 3, 1, 0, 4, 2]

# 5

## Succinct Data Structures

---



$Z + o(Z)$  бит



Self-index



Rank and Select



# 6

## Elias-Fano Encoding

---



$Z + o(Z)$  бит



Self-index



Rank and Select