

# Анализ групп метагеномных образцов с использованием графовых признаков для решения задач сравнительной метагеномики

Попов Владимир Витальевич, группа 21.M03-мм

Санкт-Петербургский государственный университет  
Математическое моделирование, программирование и  
искусственный интеллект

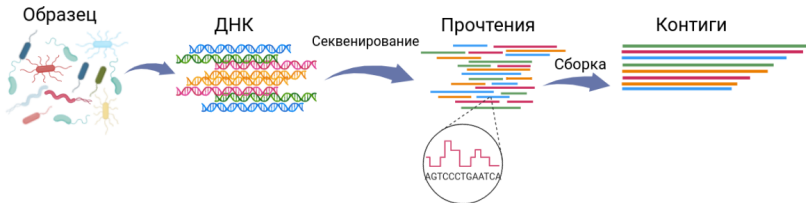
Научный руководитель: к.ф.-м.н., доцент Шпилев П.В.

Консультант: к.ф.-м.н., доцент Коробейников А.И.



Санкт-Петербург  
2023 г.

- **Геном** – хранилище наследственной информации, строка над алфавитом  $\{A, C, G, T\}$ .
- **Прочтение** – предполагаемая подстрока генома полученная в результате секвенирования.
- **Контиг** – непрерывный участок ДНК собранный из нескольких прочтений.
- **Метагеном** – набор геномов организмов, населяющих некоторую среду.

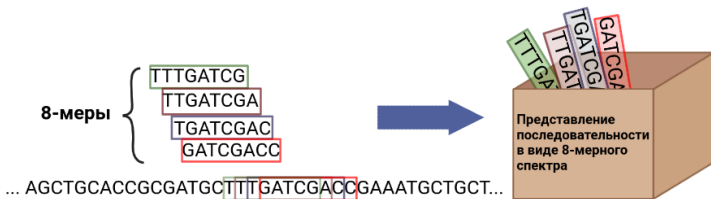


**Мотивация:** Клинические исследования, анализ изменений в составе микроорганизмов в почве и водоемах.

**Методы** сравнения групп метагеномных образцов – референсные/безреференсные.

**$k$ -мерный подход:**

- $k$ -мер – подстрока прочтения длины  $k$ .
- Подсчёт  $k$ -меров – получение словаря  $\{k\text{-мер} : \text{его встречаемость в образце}\}$ .



## Цель:

Улучшение качества методов сравнительного анализа групп метагеномных образцов.

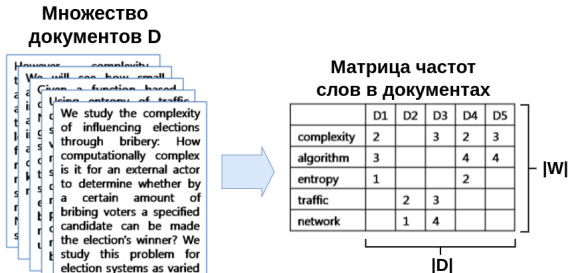
## Метод:

Построение признаков на основе подмножества  $k$ -меров, описывающих скрытую структуру наборов образцов.

## Задачи:

- Разработка и реализация эффективного алгоритма для извлечения наиболее подходящего подмножества  $k$ -меров.
- Проведение вычислительных экспериментов по классификации метагеномных образцов методами машинного обучения.

Дано:



Предположения:

- Тема — некоторое распределение на множестве слов  $W$ .
- Каждый документ описывается некоторым распределением на множестве тем.
- Порядок слов в документе не важен.
- Вероятность появления слова в документе зависит только от темы (не зависит от документа).

## Термины:

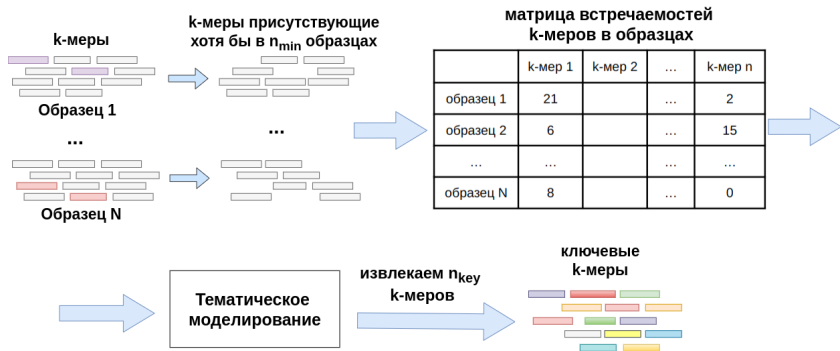
- Документы – метагеномные образцы.
- Слова –  $k$ -меры.
- Темы – скрытая структура групп образцов.

## Предположения:

- Порядок  $k$ -меров в образцах не важен.
- Вероятность появления  $k$ -мера в образце  $d$  по теме  $t$  зависит только от темы.
- Каждый образец имеет несколько смешанных в некоторой пропорции тематик.

Таким образом тематическое моделирование может быть адаптировано для извлечения ключевых  $k$ -меров.

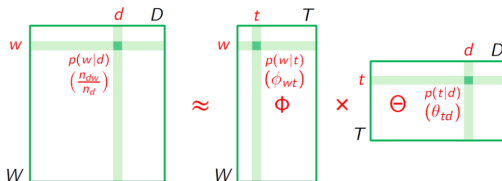
# Алгоритм извлечения ключевых $k$ -меров



Представленный алгоритм был реализован с помощью языков программирования Java и Python.

Обозначим  $T$  – набор тем,  $\varphi_{wt} = p(w|t)$ ,  $\theta_{td} = p(t|d)$ .

**Стохастическое матричное разложение:**  $P \approx \Phi\Theta$



В работе<sup>1</sup> для извлечения ключевых слов предлагается использовать величину «**значимости**» слов:

$$S(w) = p(w) \sum_{t \in T} p(w|t) \ln \frac{p(w|t)}{p(t)},$$

где  $p(t)$  — частное распределение темы,  $p(w)$  — частное распределение слова.

<sup>1</sup>Chuang J., Manning C. D., Heer J. Termite : Proceedings of the International Working Conference on Advanced Visual Interfaces. — ACM, 2012.



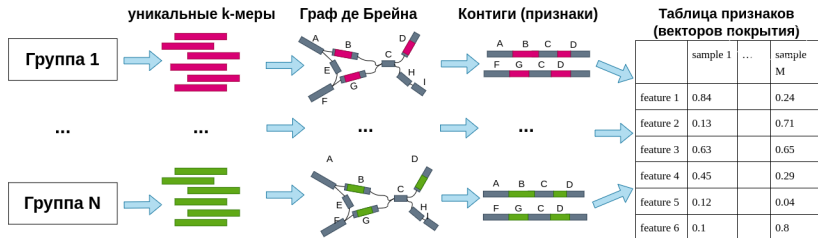
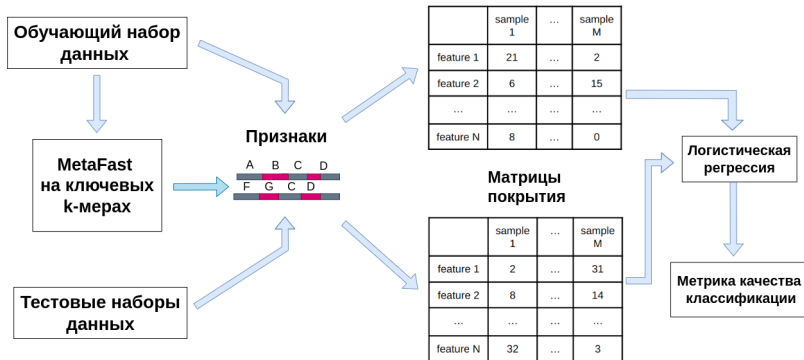


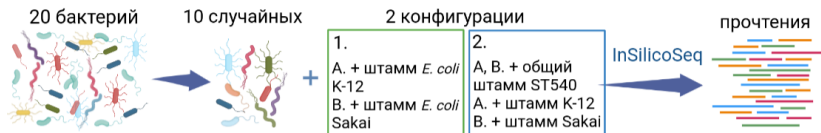
Рис.: Схема работы программы MetaFast<sup>2</sup>

<sup>2</sup>Ulyantsev V.I., Kazakov S.V., Dubinkina V.B., Tyakht A.V., Alexeev D.G. (2016). MetaFast: fast reference-free graph-based comparison of shotgun metagenomic data. Bioinformatics, 32(18), 2760-2767.

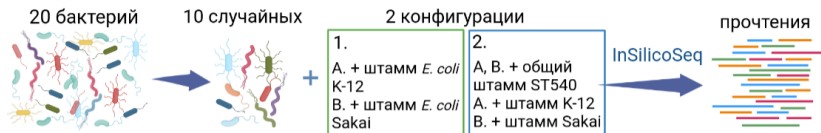
# План вычислительных экспериментов



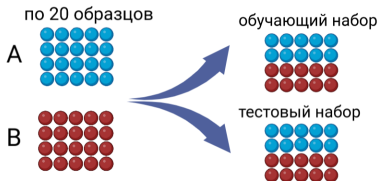
## План моделирования одного образца:



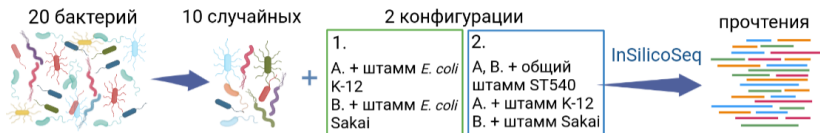
## План моделирования одного образца:



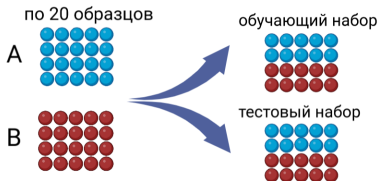
## Для каждой конфигурации:



## План моделирования одного образца:



## Для каждой конфигурации:

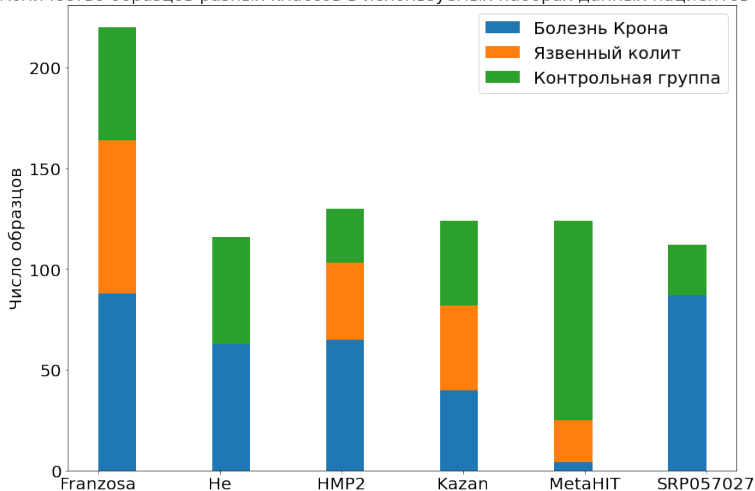


## Результаты классификации:

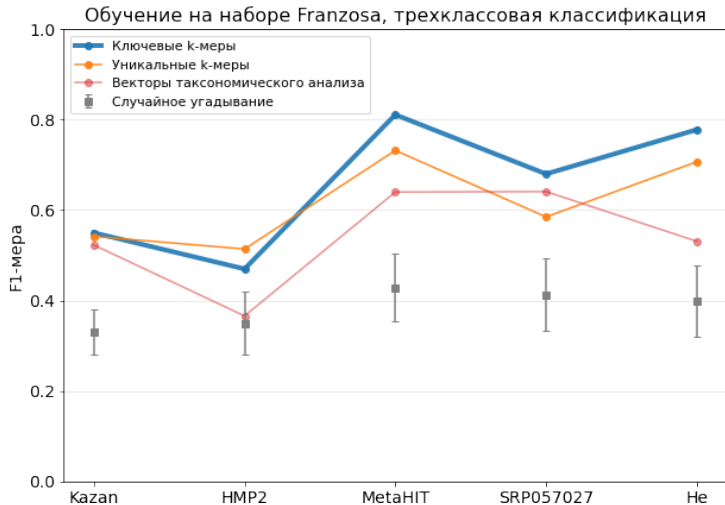
Ассигнатура на тестовом наборе	1	2
Ключевые <i>k</i> -меры	0.95	1.00
Уникальные <i>k</i> -меры	0.95	1.00
Таксономическая аннотация	0.7	0.55

# Наборы реальных метагеномных данных

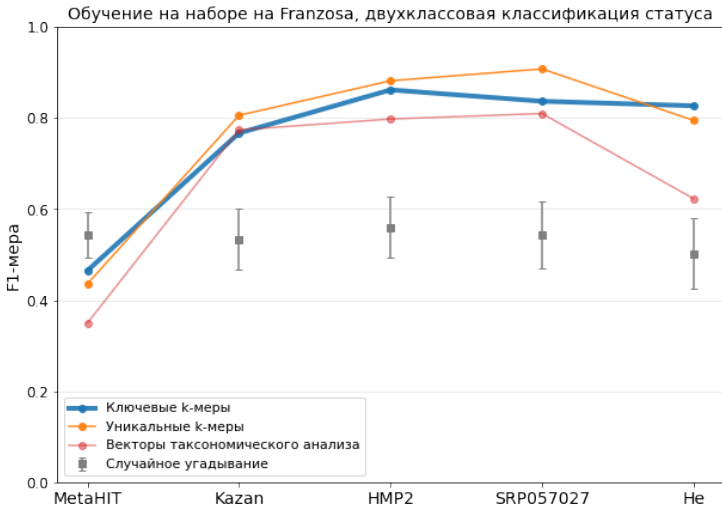
Количество образцов разных классов в используемых наборах данных пациентов с ВЗК



# Вычислительные эксперименты: трехклассовая классификация



# Вычислительные эксперименты: бинарная классификация статуса пациента





## Результаты:

- Разработан и реализован метод извлечения признаков из групп метагеномных образцов на основе ключевых  $k$ -меров.
- Метод внедрен в набор модулей программы MetaFast, работа алгоритма провалидирована на смоделированных метагеномных данных.
- Проведены вычислительные эксперименты по классификации метагеномных образцов пациентов, результаты показали преимущество разработанного метода в некоторых задачах.