



# Зачем биологам синтаксический анализ

**Автор:** Артём Горохов

Санкт-Петербургский государственный университет  
Лаборатория языковых инструментов JetBrains

15 октября 2016г.

- Множество задач, связанных с обработкой и пониманием биологических данных
- Одна из задач - поиск организмов в метагеномных сборках

- Геном (последовательность ДНК) - длинная последовательность нуклеотидов
- На деле строка над алфавитом  $\{A, C, G, T\}$

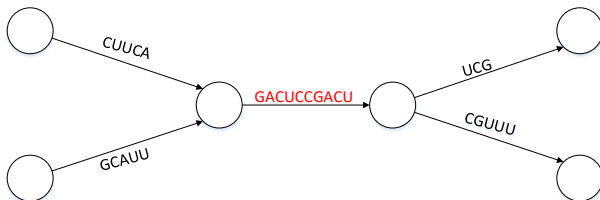
- Из биологического материала получают риды - последовательности строчек длиной около 100 символов
- Риды склеиваются в более длинные строки
- Множество полученных строк - сборка

В зависимости от биологического материала, из которого получают данные, сборки бывают:

- single-cell. Сборки, для получения которой была взята одна или несколько (до четырёх) клеток колонии
- multi-cell. В качестве биологического материала были взяты тысячи или даже десятки тысяч клеток одного штамма
- метагеномные. Данные взяты из среды обитания целевой колонии, в которой были как её представители, так и соседствующих

# Метагеномная сборка

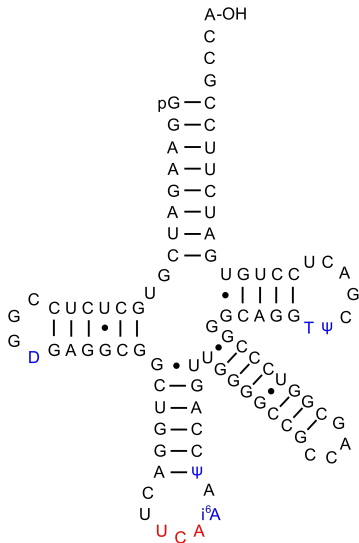
CUUCA GAC UCCGACU UCG  
UCCGACU CGUUU  
GCAUU GAC UC



# Что ищем

- Такие последовательности как тРНК, рРНК могут определить организм которому они принадлежат
- Есть у этих последовательностей есть вторичная структура, которая может быть описана КС-грамматикой

GGAAGAUCG...GCA... =>



# Грамматика для кусочка тРНК

*START* = *STEM*

*STEM* = *a STEM u*

| *u STEM a*

| *c STEM g*

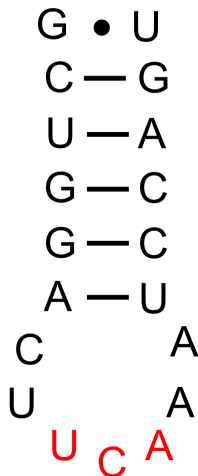
| *g STEM c*

| *g STEM u*

| *u STEM g*

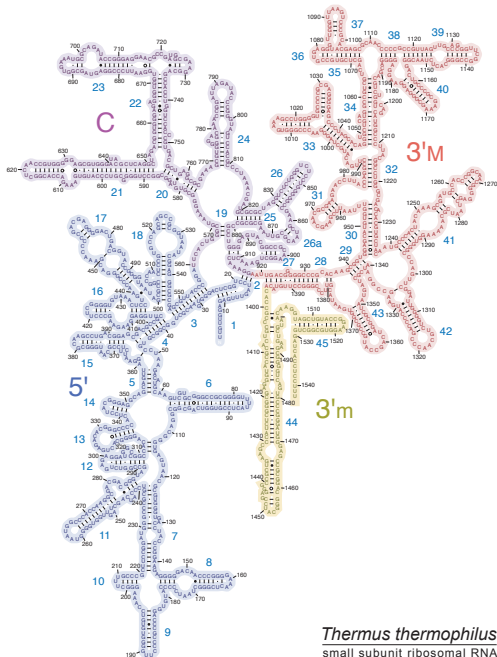
| *ANY\*[3..6]*

*ANY* = *a | u | g | c*





# Вторичная структура 16s рРНК



*Thermus thermophilus*  
small subunit ribosomal RNA

- Generalized LL
- Нисходящий синтаксический анализатор
- В лучшем случае работает за линейное время, в худшем - за  $O(n^3)$
- Строит все возможные выводы цепочки

Вход: a a b

Грамматика:

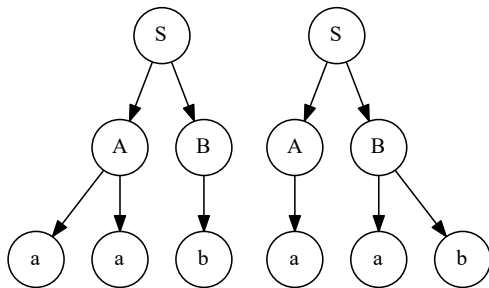
$S = A B$

$A = a a$

| a

$B = b$

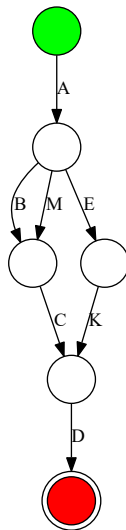
| a b



# GLL для графов

- На вход поступает граф, задающий все входные цепочки
- На рёбрах терминалы

$\{ABCD; AMCD; AEKD\} \Rightarrow$



Полученные метагеномные сборки не поддаются анализу без предварительных преобразований

- Необходимо подготавливать сборки к синтаксическому анализу
- Сам алгоритм нуждается в модернизации

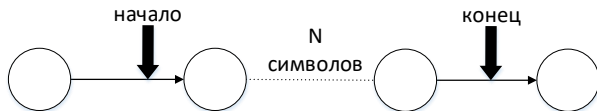
- Infernal позволяет распознавать структуры в линейном входе
- Рёбра, длиннее искомым структур можно делить на части и проверять infernal'ом

- После фильтрации рёбер граф распадается на компоненты связности
- Можно запускать анализатор независимо на разных компонентах



# Отказ от построения дерева

- Парсер возвращает лишь границы и длину найденной цепочки
- Восстановление цепочки идёт путём извлечения подграфа
- Ложные фильтруются *infernai*



# Преобразование грамматики

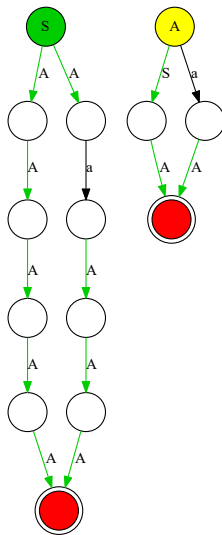
- Грамматика для 16s сильно неоднозначная и довольно большая
- Можно минимизировать её

# Преобразование грамматики к автомату

## Грамматика

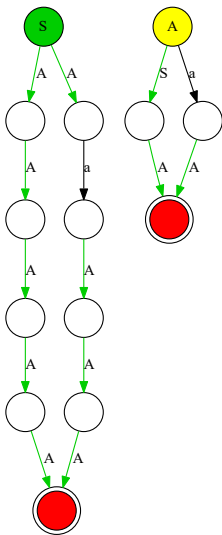
$$\begin{aligned} S &= A A A A A \\ &\quad | A a A A A \\ A &= S A \\ &\quad | a A \\ &\quad | a \end{aligned}$$

## Автомат

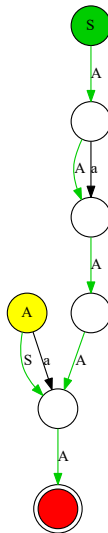


# Минимизация автомата

Изначальный автомат



Минимизированный автомат



	начальная грамматика	мин. автомат
Время работы	10 часов	3ч. 40 мин.

- Детальный анализ качества результата
- Возможно, можно сильнее фильтровать граф, применяя `infernal`
- Поиск полноразмерных 16s
- Поиск других структур