



# Context-Free Path Querying with Structural Representation of Result

**Автор:** Semyon Grigorev

Saint-Petersburg University  
Programming Languages and Tools Lab JetBrains

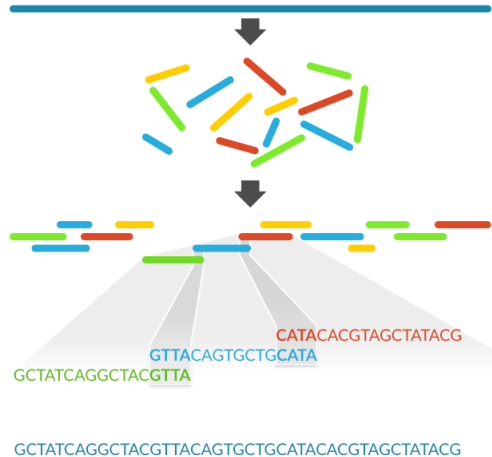
July 7, 2017

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

- Геном — длинная последовательность нуклеотидов
- На деле строка над алфавитом  $\{A, C, G, U\}$

# Получение данных

- Из биологического материала читаются короткие строчки
- Эти кусочки склеиваются в более длинные строки
- Множество строчек — сборка
- Данных очень много, поэтому строится граф, содержащий множество полученных строк



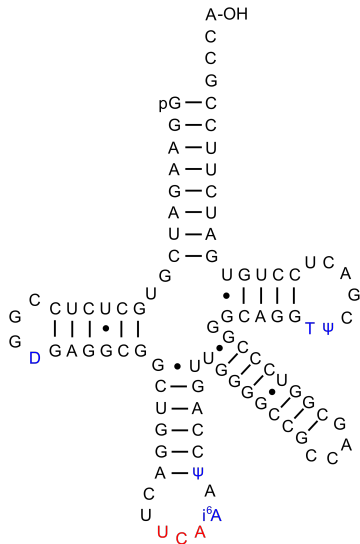
# Метагеномная сборка

- Данные из окружающей среды
- Изучаем набор генов всех микроорганизмов в образце

# Что ищем

- Хочется понять что у нас в сборке
- Такие последовательности как тРНК, рРНК позволяют провести классификацию организма
- 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*

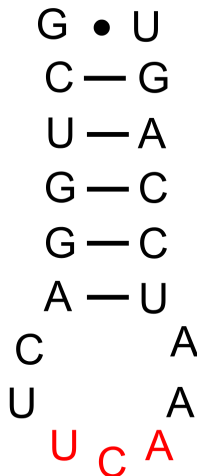
| *g STEM c*

| *g STEM u*

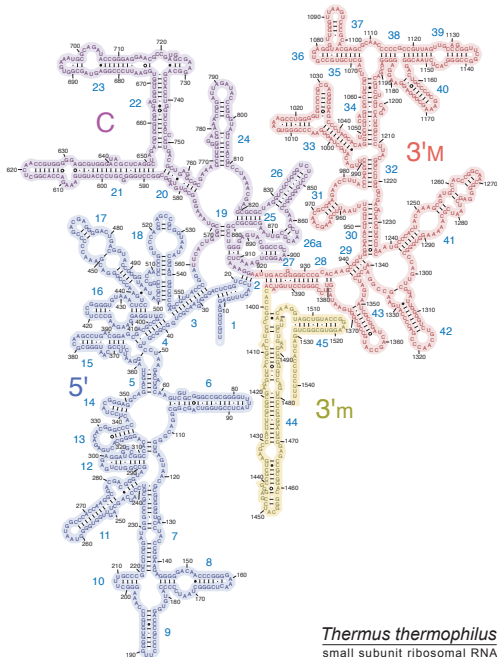
| *u STEM g*

| *ANY\*[4..7]*

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



# Вторичная структура 16s рРНК



*Thermus thermophilus*  
small subunit ribosomal RNA



- Задача поиска линейных цепочек, удовлетворяющих КС-грамматике, в графе

- В лаборатории созданы алгоритмы
- Реализован инструмент, основанный на алгоритме GLL
- Умеет решать задачу поиска линейных цепочек в графе, удовлетворяющих КС-грамматике

**Цель** работы — научиться классифицировать организмы в метагеномной сборке

**Задачи:**

- Адаптировать существующий алгоритм под специфику задачи
- Провести экспериментальные исследования работы алгоритма

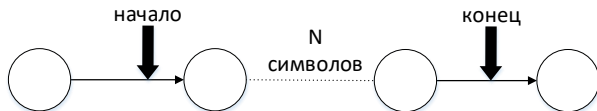
- Разбор осуществляется при помощи дескрипторов
- Дескриптор — четвёрка (слот, позиция во входе, дерево, вершина стека)
- На каждом шаге достаём дескриптор из очереди и разобрав очередной символ создаём новые дескрипторы

Метагеномные сборки довольно большие, поэтому их необходимо предварительно обрабатывать

- Infernal позволяет распознавать структуры в линейном входе
- Рёбра, длиннее искомым структур можно делить на части и проверять с помощью Infernal
- После фильтрации рёбер граф распадается на компоненты связности, на которых алгоритм можно запускать анализатор независимо

# Отказ от построения дерева

- Синтаксический анализатор возвращает лишь границы и длину найденных цепочек
- Восстановление цепочки идёт путём извлечения подграфа, состоящего из путей заданной длины
- Ложные фильтруются с помощью Infernal



- Грамматика для 16s рРНК сильно неоднозначная и довольно большая
- Из-за этого количество слотов в грамматике очень много
- В процессе разбора создаётся огромное количество дескрипторов

# Преобразование грамматики к автомату

## Грамматика

$START = STEM$

$STEM = a STEM u$

$| u STEM a$

$| c STEM g$

$| g STEM c$

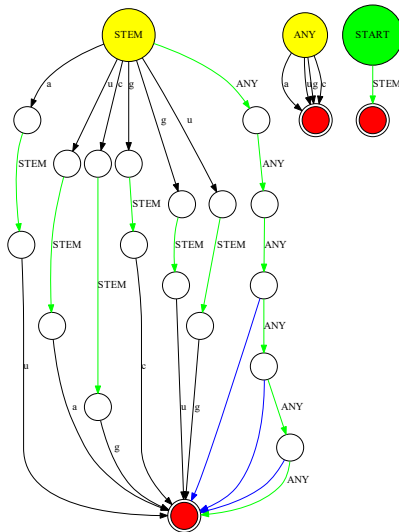
$| g STEM u$

$| u STEM g$

$| ANY^*[3..6]$

$ANY = a | u | g | c$

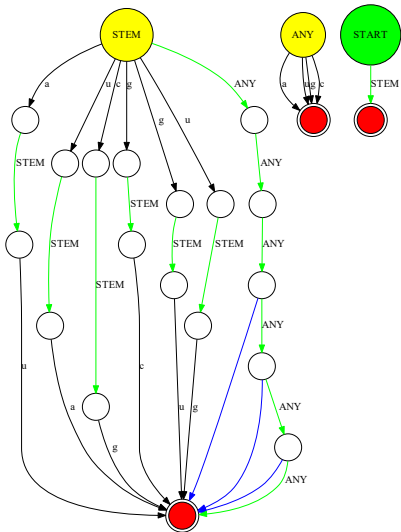
## Автомат



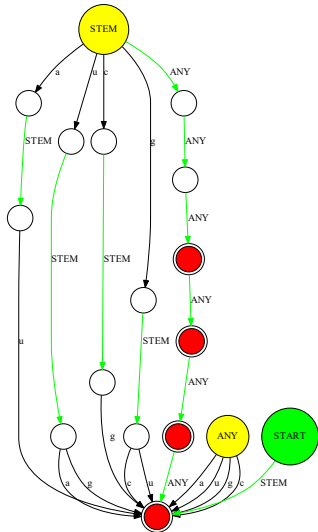


## Минимизация автомата

## Изначальный автомат



## Минимизированный автомат



Результаты работы на сборке, состоящей из 59000 вершин и 87000 рёбер и грамматике кусочка 16s pPНК, длиной около 300 символов

	начальная грамматика	мин. автомат
Время работы	10 ч.	3 ч. 40 мин.
Кол-во слотов /состояний	41	17

Эксперименты проводились на машине с 32 ГБ RAM и CPU core i7-4790

[b]0.4

- 0 :  $S \rightarrow subClassOf^{-1} S subClassOf$
- 1 :  $S \rightarrow type^{-1} S type$
- 2 :  $S \rightarrow subClassOf^{-1} subClassOf$
- 3 :  $S \rightarrow type^{-1} type$

Рис.: Grammar for query 1

[b]0.4

- 0 :  $S \rightarrow B subClassOf$
- 0 :  $S \rightarrow subClassOf$
- 1 :  $B \rightarrow subClassOf^{-1} B subClassOf$
- 2 :  $B \rightarrow subClassOf^{-1} subClassOf$

Рис.: Grammar for query 2

Таблица: Evaluation results for Query 1 and Query 2

Ontology	#triples	Query 1		Query 2	
		time(ms)	#results	time(ms)	#results
skos	252	10	810	1	1
generations	273	19	2164	1	1
travel	277	24	2499	1	1
univ-bench	293	25	2540	11	11
foaf	631	39	4118	2	2
people-pets	640	89	9472	3	3
funding	1086	212	17634	23	23
atom-primitive	425	255	15454	66	66
biomedical-measure-primitive	459	261	15156	45	45
pizza	1980	697	56195	29	29
wine	1839	819	66572	8	8