



Зачем биологам синтаксический анализ

Автор: Артём Горохов

Санкт-Петербургский государственный университет
Лаборатория языковых инструментов JetBrains

15 октября 2016г.

- 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:

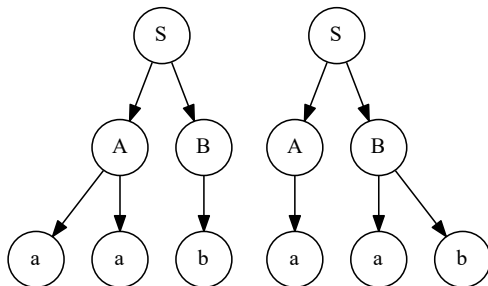
$S = A B$

$A = a a$

| a

$B = b$

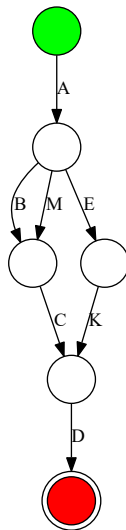
| a b



GLL для графов

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

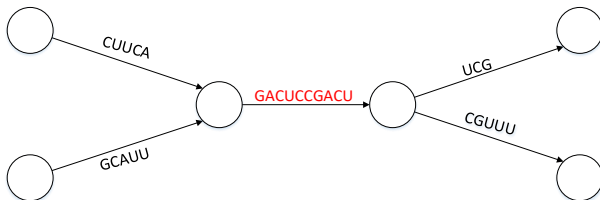
$\{ABCD; AMCD; AEKD\} \Rightarrow$



Метагеномная сборка

- Есть множество цепочек, подлежащих анализу
- Все объединяются в граф

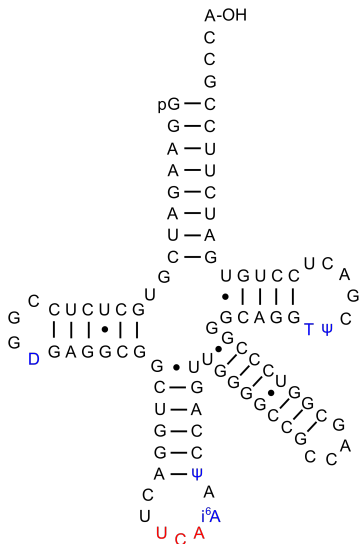
CUUCA**GAC**UCC**GACU**UCG
 UCCGACUCGUUU
GCAUU**GAC**UC

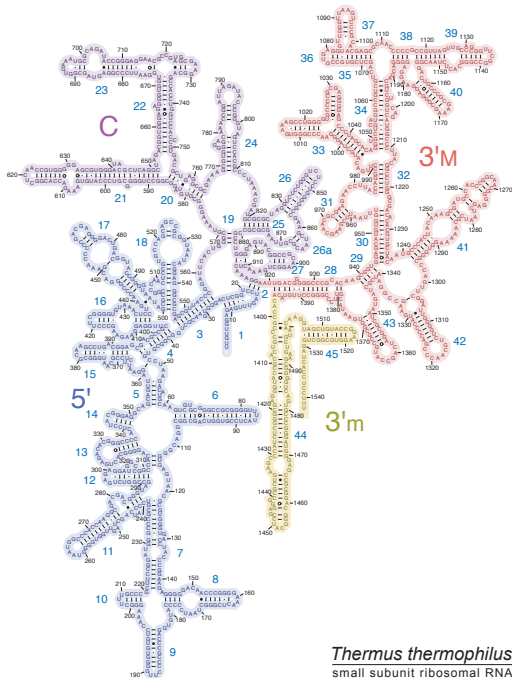


Структура цепочек

- КС грамматика может описать вторичную структуру

GGAAGAUCG...GCA... =>





Thermus thermophilus

small subunit ribosomal RNA

Увеличение производительности

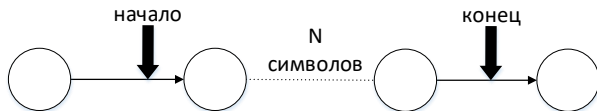
- Полученные метагеномные сборки не поддаются анализу без предварительных преобразований
- Сам алгоритм нуждается в модернизации

- Infernal позволяет распознавать структуры в линейном входе
- Рёбра, длинее 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 infernal



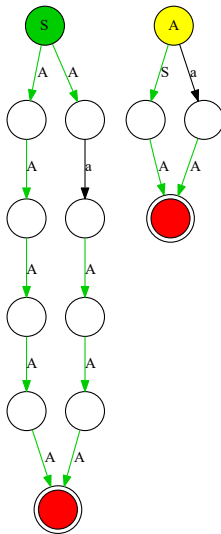
Преобразование грамматики к автомату

Грамматика

Автомат

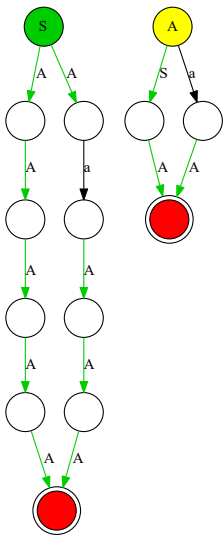
$$S = A A A A A$$

$A \ a \ A \ A \ A$

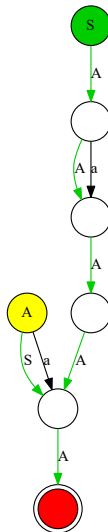
$$A = S A$$
 $| \ a \ A$ $| \quad a$ 

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 `infernai`
- Поиск полноразмерных 16s
- Поиск других структур