

Реализация возможности сжатия строки в КС-грамматику в YaccConstructor

Автор: Зиновьева А.Г.

Руководитель: к.ф.-м.н. Григорьев С.В.

Введение

- Информация может содержать повторяющиеся фрагменты, которые также могут содержать повторяющиеся фрагменты и так далее
- Можно описать повторяющийся фрагмент только один раз, а не для каждого случая отдельно
- Можно создать грамматику с правилами, которые будут описывать эти повторяющиеся фрагменты
- В YaccConstructor нет возможности построить грамматику для строки

Задачи

Для реализации возможности сжатия текстовой строки в грамматику в YaccConstructor были поставлены следующие задачи

- Реализовать алгоритм сжатия строки Sequitur
- Реализовать возможность построения общей грамматики для нескольких строк
- Реализовать конечное представление грамматики в формате YARD.IL
- Оформить фронтенд для данного алгоритма к YaccConstructor
- Проверить эффективность данного решения

Алгоритм Sequitur

- Автор - К. Невилл-Манин, 1997 год
- Принимает на вход последовательность дискретных символов: например, текстовую строку
- Посимвольно обрабатывает строку
- Ищет повторяющиеся диграмы (пары рядом стоящих символов) и заменяет их на нетерминальные
- Результатом алгоритма является контекстно-свободная грамматика

Алгоритм Sequitur

Свойства грамматики

- Уникальность: не может быть двух правил с одинаковой правой частью
 - Невозможна такая ситуация: $A \rightarrow ab$ и $B \rightarrow ab$
- Полезность: каждое правило используется более одного раза
 - $S \rightarrow AA$ $A \rightarrow Bc$ $B \rightarrow ab$: правило B не является полезным, поэтому правило A должно быть преобразовано в $A \rightarrow abc$

Алгоритм Sequitur

Входные данные: “abcdabc”

1. $S \rightarrow a$ 2. $S \rightarrow ab$ 3. $S \rightarrow abc$ 4. $S \rightarrow abcd$ 5. $S \rightarrow abcda$

Новых правил не образовалось, т.к. нет повторяющихся диграмов

6. $S \rightarrow abcdab$: $S \rightarrow AcdA$ $A \rightarrow ab$

Встретился повторяющийся диграм, добавляем новое правило

7. $S \rightarrow AcdAc$ $A \rightarrow ab$: $S \rightarrow BdB$ $B \rightarrow Ac$ $A \rightarrow ab$: $S \rightarrow BdB$ $B \rightarrow abc$

Создали новое правило, но правило A стало бесполезным, поэтому избавляемся от него

Результат алгоритма: $S \rightarrow BdB$ $B \rightarrow abc$

YARD

- Язык спецификаций грамматик, являющийся частью YaccConstructor
- Поддерживает EBNF (Расширенная Форма Бэкуса-Наура)
- Позволяет описать результат работы алгоритма с помощью двух выражений
 - выбор ($A|B$)
 - конкатенация ($A \cdot B$)

Реализация

- Выбор структуры данных для описания грамматики – двусвязный список
 - добавление нового символа в правило
 - замена двух символов на нетерминал
- Основная функция реализации – обработка стека диграм

Реализация

- Модификация алгоритма для сжатия нескольких строк: обрабатываются только диграмы, не содержащие специальный символ
 - Для строки $ab\&ac\&ab\&ac$ результат будет $A\&B\&A\&B$ $A \rightarrow ab$ $B \rightarrow ac$
- Представление в YARD.IL – построение грамматики с помощью конструкторов
 - PTok – для описания терминалов PTok(a)
 - PRef – для описания нетерминалов PRef(A)
 - PSeq – для описания последовательностей PSeq([PTok(a); PTok(b)])
 - PAlt – для разделения по специальным символам PAlt(PRef(A), PAlt(..))

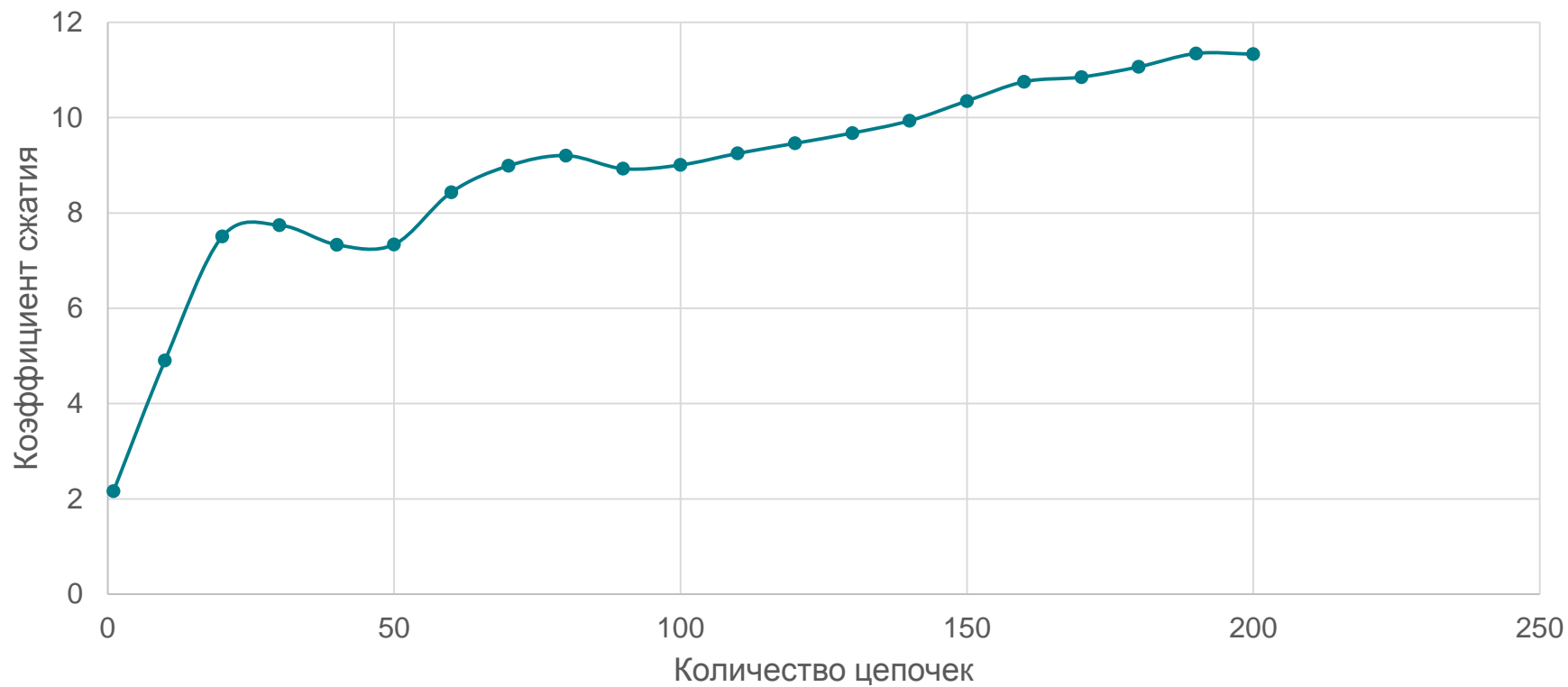
Реализация

- YaccConstructor имеет модульную структуру
- Фронтенды - модули, которые из входных данных пользователя создают представление грамматики в YARD.IL
- Был оформлен фронтенд для данного алгоритма с возможностью
 - Сжать строку без специальных символов
 - Сжать строку со специальными символами
 - Сжать массив строк

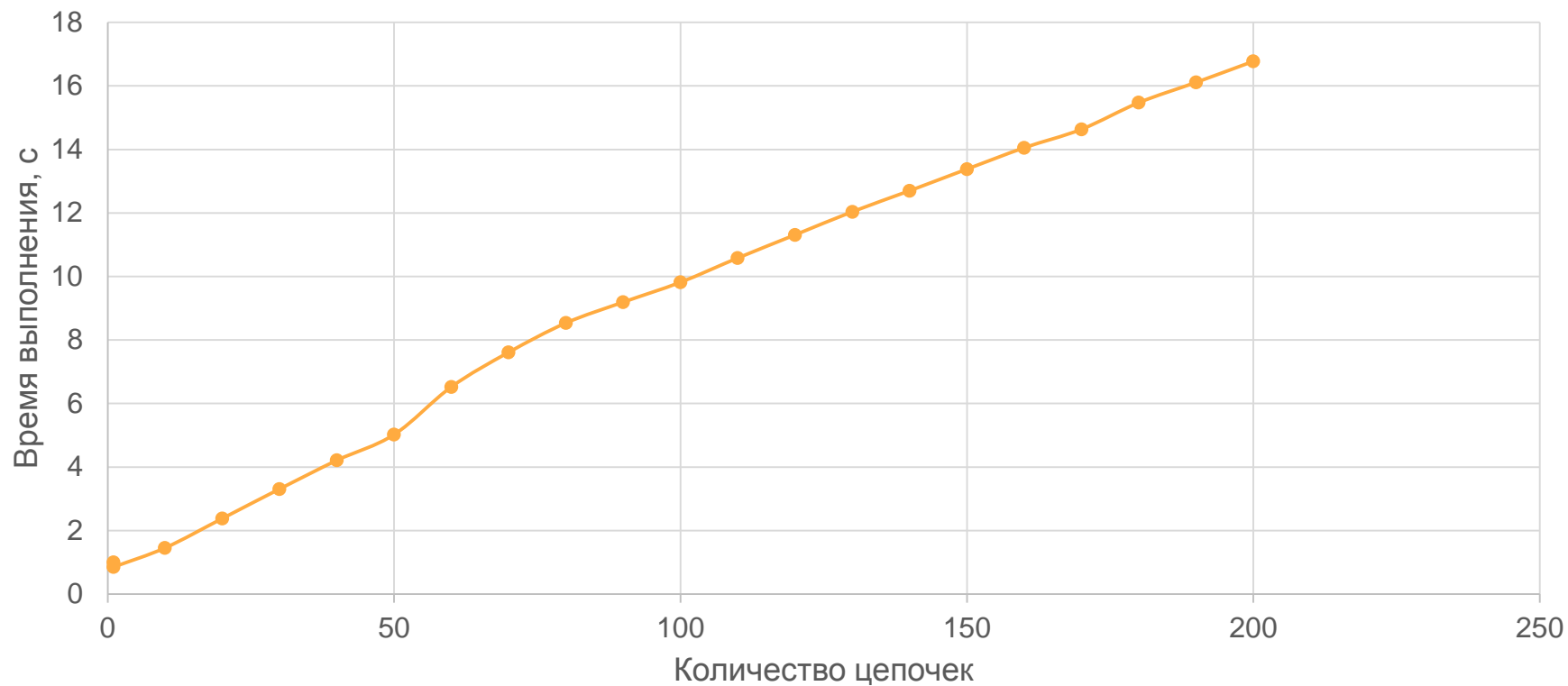
Эксперимент

- Для проверки эффективности решения из базы данных SILVA была взята последовательность 16s РНК бактерий
- Каждая цепочка состоит примерно из 1500 нуклеотидов, то есть последовательность символов из алфавита {A; C; G; T}
- Данные последовательности были склеены через специальный символом & и переданы на вход алгоритму

Эффективность



Производительность



Результаты

- Реализован алгоритм Sequitur
- Реализована возможность построения общей грамматики для нескольких строк
- Оформлен фронтенд к YaccConstructor для построения грамматики из строки в формате YARD.LL
- Проверена эффективность и производительность данного алгоритма на последовательностях РНК