

**Программа ChromoGGL**  
**для вычисления расстояний между геномными структурами**  
**и построения эволюционных сценариев**

**Руководство пользователя**

**(<http://lab6.iitp.ru/ru/chromoggl>)**

**(авторы: В.А. Любецкий, К.Ю. Горбунов, Р.А. Гершгорин)**

**2019 год**

## Оглавление

|  |    |
|--|----|
| Общие сведения.....  | 3  |
| Алгоритм вычисления расстояния между двумя хромосомными структурами .....  | 3  |
| Алгоритм приведения общего графа к финальному виду в случае различных цен операций и дополнительных операций удаления и вставки..... | 6  |
| Алгоритм построения дерева эволюции хромосомных структур .....   | 15 |
| Установка и запуск в среде Windows программы ChromoGGL.....  | 16 |
| Входные данные программы ChromoGGL.....  | 16 |
| Файл структур.....   | 16 |
| Файл с ценами операций .....   | 17 |
| Файл с весами расстояний разных типов.....   | 17 |
| Параметры программы ChromoGGL .....  | 18 |
| Командная строка.....  | 18 |
| Выходные данные программы ChromoGGL в задаче 1 .....   | 18 |
| Примеры решения задачи 1 .....   | 20 |
| Выходные данные программы ChromoGGL в задаче 2 .....   | 22 |
| Примеры решения задачи 2.....  | 23 |
| Вспомогательная утилита gbk_converter .....  | 24 |
| Рекомендуемые стандартные программы .....  | 26 |
| Литература.....  | 26 |

## Общие сведения

Программа ChromoGGL предназначена для решения следующих двух задач:

**1) Вычисление расстояния между двумя данными хромосомными (синоним: геномными) структурами и построение кратчайшей последовательности операций, преобразующих одну структуру в другую.** Рассматривается общее определение структуры как произвольного множества путей и циклов, представляющих линейные и кольцевые хромосомы, вместе с операциями, которые преобразуют одну структуру в другую. Структуры могут включать паралогичные гены, если указано их правильное соответствие. Последовательность операций допускает переменный генный состав. Допускаются произвольные цены отдельных операций. Задача состоит в минимизации суммарной цены операций в последовательности, которая преобразует одну структуру в другую. Последовательность, на которой достигается минимум, называется *кратчайшей*. Среднее арифметическое суммарных цен при переходе от одной структуры к другой и наоборот называется *расстоянием* между ними. Несмотря на столь общую постановку, предложенный нами оригинальный алгоритм является *линейным* по времени работы и по используемой памяти; и также весьма быстрым при практическом счёте на суперкомпьютере. При этом для него доказываемая *точность*, если выполняются некоторые условия на цены отдельных операций. Эти условия отличаются от известных условий.

**2) Вычисление матрицы попарных расстояний между хромосомными структурами из данного множества и построение дерева эволюции этих структур, наилучшим образом согласованного с полученной матрицей.** Исходные структуры соответствуют листьям дерева. Применяется простой линейный алгоритм, основанный на алгоритме решения задачи 1. Построение дерева на основании вычисленной матрицы осуществляется известными филогенетическими алгоритмами.

В следующем разделе описан алгоритм решения задачи 1, который является ключевым и для решения задачи 2.

### Алгоритм вычисления расстояния между двумя хромосомными структурами

**Модель хромосомной структуры** описывается как конечное множество ориентированных цепей и циклов, включая петли. Такое множество можно рассматривать как ориентированный граф, который будем называть *хромосомной структурой*. Ребро графа будем называть *геном*; отдельную цепь или отдельный цикл графа – *хромосомой* или *компонентой*. Каждому гену приписано имя, обычно *номер  $i$*  этого гена, который может повторяться (в случае паралогов) и тогда номер принимает вид  *$i,j$* . Такая модель, как обычно,

означает, что не учитываются длины генов и межгенных участков, как и состав генов и межгенных участков; направление ребра показывает, на какой цепи лежит ген. Вершина графа показывает «место» соединения соседних генов, независимо от их цепи, т.е. в вершине *отождествляются* (мы говорим, *склеиваются*) два края соседних генов. Обычно в структуре много цепей и циклов, что приводит к их своеобразному взаимодействию, поэтому ситуация многих хромосом в структуре решительно отличается от ситуации одной хромосомы.

Модель включает следующие *стандартные* операции над хромосомной структурой. *Двойная переклейка* – расклейка двух склеек краёв генов и новая переклейка четырёх краёв; *полупорная переклейка* – расклейка двух склеенных краёв и склеивание одного края с каким-то несклеенным краем, второй край остаётся *свободным*; *разрез* или *склейка* – соответственно расклейка двух склеенных краёв с образованием двух свободных краёв или склейка двух свободных краёв. Пусть даны хромосомные структуры  $a$  и  $b$ , *общим* (*особым*) называется ген, который принадлежит обеим структурам (только одной из них); ген из структуры  $a$  называется *a-геном*, соответственно определяется *b-ген*. Модель включает две *дополнительные* операции (подразумевается преобразование  $a$  в  $b$ ): *удаление* (связного максимального) участка особых  $a$ -генов и *вставка* участка особых  $b$ -генов. При удалении, если участок находился строго внутри цепи или цикла, два образовавшихся свободных конца общих генов склеиваются между собой; если он находился с краю цепи, край общего гена становится свободным; наконец, если он являлся отдельной хромосомой, она удаляется целиком. Если участок вставляется строго внутри цепи или цикла, место вставки предварительно расклеивается; вставка может выполняться с краю цепи или как новая хромосома. Нетрудно доказать, что использование немаксимальных удалений особых генов не приводит к расширению возможностей, как и разрезание участка особых  $a$ -генов в первых трёх операциях или операции вставки. Поэтому эти возможности не рассматриваются.

Итак, даны шесть операций и каждой из них приписано положительное рациональное число, которое называется её *ценой*. Включение в модель цен операций является её принципиальной особенностью нашей модели.

Напомним, задача состоит в поиске *кратчайшей* последовательности из этих операций, которая переводит структуру  $a$  в структуру  $b$ . Здесь «кратчайшая» означает последовательность, у которой суммарная цена составляющих её операций минимальна. В последовательности каждая операция рассматривается вместе с хромосомной структурой, к которой она применяется.

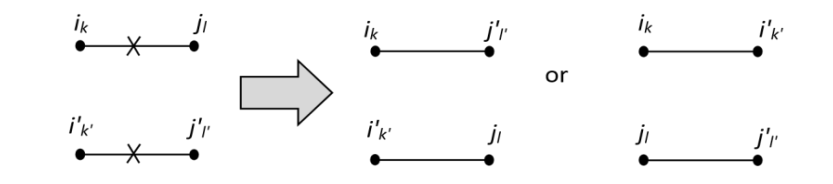
**Определение общего графа и его финального вида.** В общем графе  $a+b$  двух структур  $a$  и  $b$  вершины – края общих генов и, кроме того, все максимальные участки

особых генов. Каждый край берётся один раз; более формально: вместо края пишется имя гена с индексом 1 или 2, указывающим на его начало или конец. Вершины первого типа называются *обычными*, а второго – *особыми*. Ребро соединяет две обычные вершины, если в одной из структур края склеены, т.е. примыкают друг к другу на хромосоме. Ребро соединяет обычную вершину с особой, если край общего гена склеен с крайним геном участка особых генов. Рёбра из первого случая называются *обычными*, а из второго – *особыми*. В цепи крайнее ребро с особым краем назовём *висячим*. Ребро помечается  $a$  или  $b$  в зависимости от того, в какой из них имеется склейка; вершины могут соединяться двумя рёбрами. Особые вершины делятся на  $a$ - и  $b$ -вершины. В графе могут быть изолированные вершины – участки из особых генов: если такой участок – цикл, то проводим в нём петлю, которую назовём *особой*. Получается неориентированный граф.

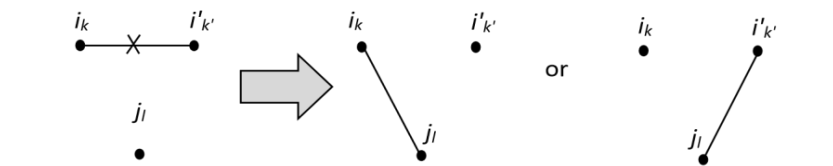
К общему графу  $a+b$  применяются аналоги операций над структурами, которые описываются следующим образом. (1) Удалить два одинаково помеченных ребра и четыре образовавшихся конца соединить двумя новыми неинцидентными рёбрами с той же пометкой. (2) Удалить ребро (скажем, с пометкой  $a$ ) и соединить  $a$ -ребром один из его концов с обычной вершиной, не инцидентной  $a$ -ребру или с особой  $a$ -вершиной, имеющей не более одного инцидентного  $a$ -ребра. (3) Удалить любое ребро. (4) Добавить ребро (скажем, с пометкой  $a$ ) между вершинами, не инцидентными  $a$ -ребру. Если в результате операции получаются две инцидентные особые вершины, они сливаются в одну вершину, что входит в определение операции; получаемой вершине приписывается объединение имён исходных вершин. (5) Удаления особой вершины или особой петли; если эта вершина имела две инцидентные ей обычные вершины, они соединяются ребром. Легко определить аналог операция вставки, но оказывается, что без неё можно обойтись без потери общности, что является нетривиальным утверждением.

## Стандартные операции

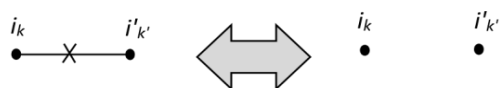
1) Double-cut-and-paste (DP). Исходные рёбра одноимённые.



2) Sesqui-cut-and-paste (SP)



3) Удаление  $a$ -ребра или добавление  $b$ -ребра (C) и удаление  $b$ -ребра или добавление  $a$ -ребра (J)



*Финальный вид* общего графа  $a+b$  определяется как общий граф, состоящий из изолированных обычных вершин и *финальных 2-циклов*. Последнее определяется как граф из двух обычных вершин, соединённых обычными рёбрами, одно из  $a$  и одно из  $b$ . Легко показать, что *исходная задача эквивалентна* приведению графа  $a+b$  к финальному виду при некоторых ограничениях на цены операций, в том числе, если операции, кроме вставки и удаления, имеют одинаковые цены.

### Алгоритм приведения общего графа к финальному виду в случае различных цен операций и дополнительных операций удаления и вставки

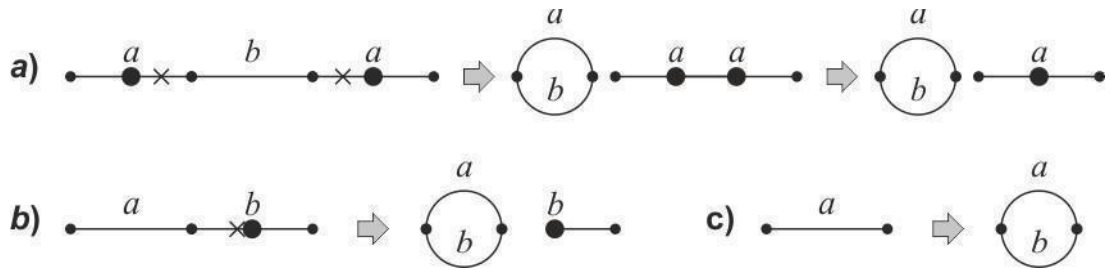
Напомним: в этом разделе соответствие паралогов предполагается известным, но допускаются все операции, неравный генный состав и любые цены операций. По исходным структурам  $a$  или  $b$  тривиально строится общий граф  $a+b$ . Цель алгоритма – привести  $a+b$  к финальному виду.

Опишем алгоритм приведения общего графа к финальному виду. При единой общей схеме алгоритм зависит от соотношения цен операций, которые заранее фиксируются.

**Шаг 1.** Удалить все особые  $a$ -петли.

**Шаг 2.** Удалить обычное ребро, не входящее в 2-цикл, и замкнуть его в финальный 2-цикл двойной (если внутреннее) или полуторной (если крайнее) переклейкой или склейкой (если единственное). Повторять действие, пока возможно. Если цена двойной переклейки не

больше цены полуторной, то сначала выполняются все двойные переклейки, иначе – все полуторные.



Перед формальным описанием Шагов 3-4 поясним их смысл. На шаге 3 определённые комбинации от двух до четырёх цепей несколькими операциями заменяются на одну цепь. Каждая применяемая операция уменьшает число особых вершин на 1. Принципиальный момент состоит в определении этих комбинаций. Это делается через понятие типа, которое мы введём чуть ниже. Шаг 4 применяется, если удаление  $b$ -вершины дороже по цене всех других операций. Полученное к этому шагу множество хромосом (т.е. путей и циклов) разбивается на пары, совместная обработка которых заменяет операцию удаления  $b$ -вершины на более дешёвую операцию переклейки со слиянием двух  $b$ -вершин; при этом общее число операций не меняется. Нами доказано, что порядок действий, указанный ниже в шагах 3-4, обеспечивает оптимальность результата.

Итак, перейдём к определению типа цепи и цикла. *Длиной* цепи или цикла назовём число внутренних особых вершин сложенное с числом обычных рёбер. *2-Цикл* – цикл длины 2. *Нечётной (чётной)* назовём цепь, у которой длина нечётная (чётная). *a-Цепью* называется изолированная  $b$ -вершина или нечётная цепь, у которой крайние ребра с обычным краем помечены  $a$ ; симметрично определяется *b-цепь*. Оставшимся после шагов 1-2 цепям и циклам (исключая финальные 2-циклы и изолированные обычные вершины), *приписываем тип*: 2-циклу, содержащему  $a$ -вершину, но не  $b$ -вершину, – « $a$ -цикл»; симметрично – « $b$ -цикл». Циклу, в котором имеются как  $a$ -, так и  $b$ -вершины, приписываем тип «цикл». Особой  $b$ -петле приписываем тип «петля». *a-Цепи приписываем тип*:  $1a$ , если в ней одно висячее ребро;  $2a$ , если в ней два таких ребра;  $2a'$ , если она – изолированная  $b$ -вершина;  $3a$ , если в ней нет висячих рёбер, но имеются как  $a$ -, так и  $b$ -вершины (тогда длина цепи строго больше 1);  $3a'$ , если нет ни висячих рёбер, ни  $b$ -вершин (тогда длина цепи равна 1). *b-Цепям тип приписывается аналогично*. Чётной цепи *приписываем тип*: 1, если в ней одно висячее ребро и имеются  $b$ - и  $a$ -вершины;  $1'$ , если она состоит из одной обычной вершины и инцидентной ей  $a$ -вершины;  $1''$ , если состоит из одной обычной вершины и инцидентной ей  $b$ -вершины; 2, если в ней два висячих ребра и имеется невисячее ребро;  $2'$ , если в ней только два висячих ребра; 3, если ребра имеются, но нет висячих рёбер. Цепи типа 1 разделим на цепи типов  $1a$ , если крайняя особая вершина –  $a$ -вершина, и  $1b$ , если она –  $b$ -вершина.

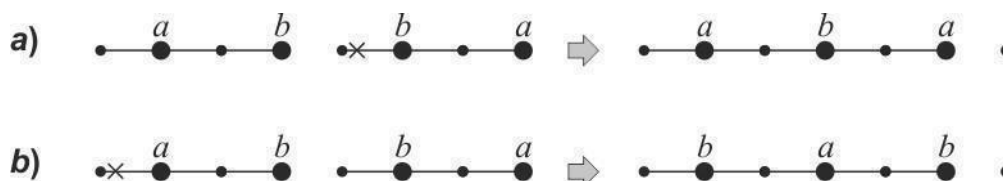
Введём особый тип  $1_c$ : он означает *отложенный выбор между цепями* типов  $1_a$  и  $1_b$  – двумя возможными результатами операции. Алгоритм хранит оба результата до шагов 4.15–4.23, на которых из каждой пары результатов выбирается один и, тем самым, однозначно определяется цепочка операций.

Тип  $2a$  означает объединение типов  $2a$  и  $2a'$ , тип  $3b$  – типов  $3b$  и  $3b'$ , тип  $1_b$  – типов  $1_b$  и  $1''$ ; тип  $2$  – типов  $2$  и  $2'$ .

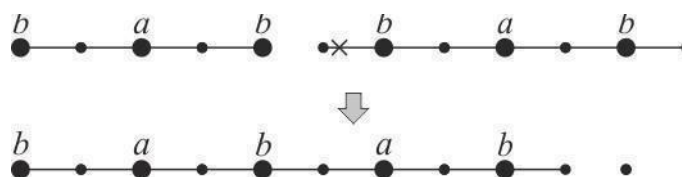
Ниже в описании шага 3 перед знаком « $\Rightarrow$ » перечисляются (разделяемые знаком « $+$ ») типы цепей, их комбинация преобразуется в комбинацию цепей с типами, указанными после « $\Rightarrow$ »; в полученной комбинации изолированные обычные вершины и финальные 2-циклы не указываем, им не приписан тип. Если ниже пункт содержит несколько равенств, то описывается действие, относящиеся к первому равенству; другие действия аналогичны.

**Шаг 3.** В указанном порядке выполняем следующие действия, каждое из них до тех пор, пока оно возможно.

3.1.  $1a+1b=1_c$ . Расклеим крайнее невисячее ребро (назовём его *внешним*) в одной из двух цепей типа  $1a$  или  $1b$  и соответствующую особую вершину склеим с крайней особой вершиной другой цепи (полуторная переклейка). Два варианта показаны на рисунке:

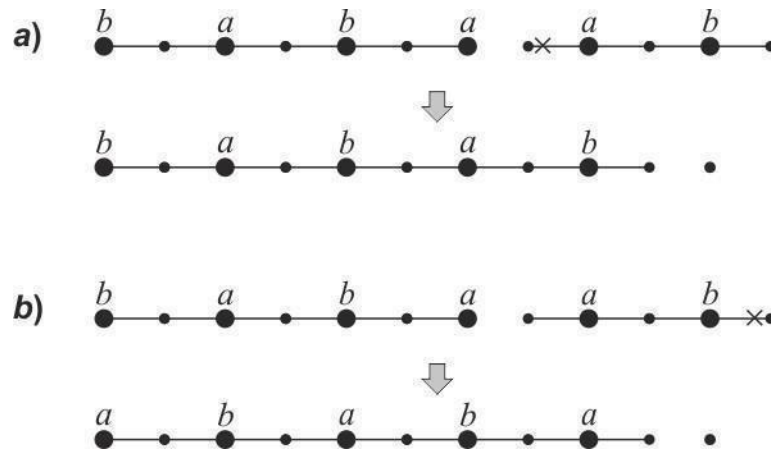


3.2.  $2a+3b=1_b$ ,  $2b+3a=1_a$ ,  $2b'+3a=1_a$ ,  $2b+3a'=1_a$ ,  $2b'+3a'=1'$ . в  $3b$ -цепи расклеим внешнее ребро и особую вершину склеим с крайней особой вершиной  $2a$ -цепи:



3.3.  $2+3=1_c$ . В 3-цепи расклеим внешнее ребро и особую вершину склеим с крайней особой вершиной 2-цепи. В зависимости от того, какое из двух внешних ребро расклеивается, в результате получается цепь типа  $1_a$  или  $1_b$ , рис. *a* и *b*.

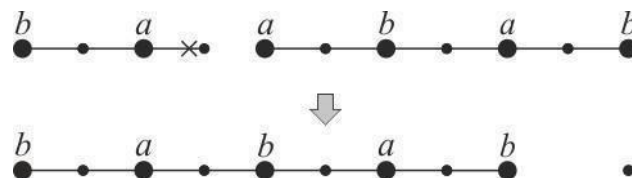




3.4.  $1b+2a+3=2+3=1_c$ ,  $1a+2b+3=2+3=1_c$ ,  $1a+2b'+3=2+3=1_c$ . Сначала выполняем  $1b+2a=2$  (описание ниже); затем  $2+3=1_c$ .

3.5.  $1a+3b+2=3+2=1_c$ ,  $1b+3a+2=3+2=1_c$ ,  $1b+3a'+2=3+2=1_c$ . Сначала  $1a+3b=3$  (описание ниже); затем  $2+3=1_c$ .

3.6.  $1a+2=2a$ ,  $1b+2=2b$ . В  $1a$ -цепи расклеим внешнее ребро и особую вершину склеим с крайней особой вершиной  $2$ -цепи:



3.7.  $1a+3=3a$ ,  $1b+3=3b$ . В  $3$ -цепи расклеим крайнее  $b$ -ребро и особую вершину склеим с крайней особой вершиной  $1a$ -цепи:



3.8.  $1a+1a+2b+3b=2+3=1_c$ ,  $1a+1a+2b'+3b=2+3=1_c$ ,  $1b+1b+2a+3a=2+3=1_c$ ,

$1b+1b+2a+3a'=2+3=1_c$ . Сначала выполняем  $1a+2b=2$  и  $1a+3b=3$ ; затем  $2+3=1_c$ .

3.9.  $1a+1a+2b=3a+2b=1_a$ ,  $1a+1a+2b'=3a+2b'=1_a$ ,  $1b+1b+2a=3b+2a=1_b$ . Сначала  $1a+1a=3a$  (описание ниже); затем  $2b+3a=1_a$ .

3.10.  $1a+1a+3b=1a+3=3a$ ,  $1b+1b+3a=1b+3=3b$ ,  $1b+1b+3a'=1b+3=3b$ . Сначала  $1a+3b=3$ ; затем  $1a+3=3a$ .

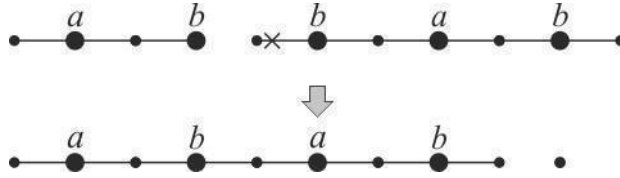
3.11.  $1a+1a=3a$ ,  $1b+1b=3b$ . Склеим крайние особые вершины двух  $1a$ -цепей:



3.12.  $1a+2b=2$ ,  $1a+2b'=2$ ,  $1b+2a=2$ . В  $1a$ -цепи расклеим внешнее ребро и особую вершину склеим с крайней особой вершиной  $2b$ -цепи:



3.13.  $1a+3b=3$ ,  $1b+3a=3$ ,  $1b+3a'=3$ . В  $3b$ -цепи расклеим внешнее ребро и особую вершину склеим с крайней особой вершиной  $1a$ -цепи:



3.14.  $2a+2b+3+3=2+3=1_c$ ,  $2a+2b'+3+3=2+3=1_c$ . Сначала  $2a+2b+3=2$  (описание ниже), затем  $2+3=1_c$ .

3.15.  $3a+3b+2+2=3+2=1_c$ ,  $3a'+3b+2+2=3+2=1_c$ . Сначала  $3a+3b+2=3$  (описание ниже), затем  $2+3=1_c$ .

3.16.  $2a+3+3=1a+3=3a$ ,  $2b+3+3=1b+3=3b$ ,  $2b'+3+3=1b+3=3b$ . Сначала  $2a+3=1a$  (описание ниже в шаге 4), затем  $1a+3=3a$ .

3.17.  $3b+2+2=1b+2=2b$ ,  $3a+2+2=1a+2=2a$ ,  $3a'+2+2=1a+2=2a$ . Сначала  $3b+2=1b$  (описание ниже в шаге 4); затем  $1b+2=2b$ .

3.18.  $2a+2b+3=2a+1b=2$ ,  $2a+2b'+3=2a+1b=2$ . Сначала  $2b+3=1b$ ; затем  $1b+2a=2$ .

3.19.  $3a+3b+2=3a+1b=3$ ,  $3a'+3b+2=3a'+1b=3$ . Сначала  $3b+2=1b$ ; затем  $1b+3a=3$ .

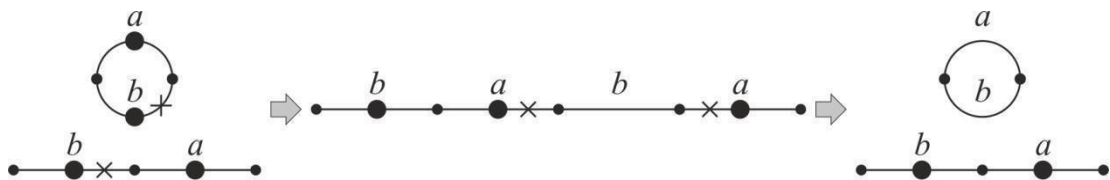
**Шаг 4.** Если цена двойной переклейки больше цены полуторной, то в следующем порядке применяем каждый указанный пункт 4.1–4.24, пока это возможно; иначе таким же образом применяем пункты 4.1'–4.24'.

4.1. «петля»+любой тип  $t$  с  $b$ -вершиной = тип  $t$ . Объединить  $b$ -вершину петли с  $b$ -вершиной компоненты типа  $t$  двойной переклейкой (если эта цепь не изолированная  $b$ -вершина) или полуторной переклейкой (иначе):



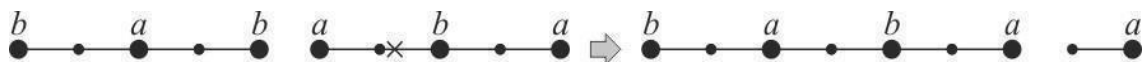
4.1'. То же, что и 4.1.

4.2. «цикл»+любой тип  $t$  с  $b$ -вершиной и  $a$ -вершиной = тип  $t$ . Вставить цикл (двойной переклейкой, сливающей две  $b$ -вершины) рядом с  $b$ -вершиной из компоненты типа  $t$  с той стороны, в которой находится  $a$ -вершина; образовавшееся обычное ребро вырезать:



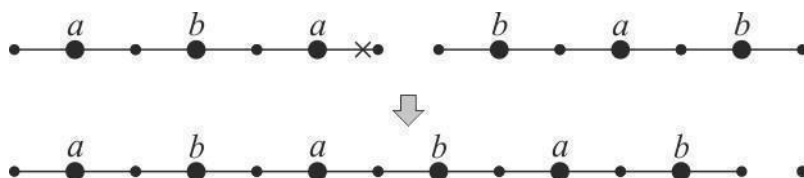
4.2'. То же, что и 4.2.

4.3.  $2a+2b=2+1'$ . Полуторная переклейка с отрезанием двух вершин  $2b$ -цепи (крайней  $a$ -вершины и соседней обычной вершины) и склейкой образовавшегося края с крайней  $b$ -вершиной  $2a$ -цепи:



4.3'.  $2a'+2b=2+1'$ .

4.4.  $3a+3b=3$ . В  $3a$ -цепи расклеить внешнее ребро и особую вершину склеить с крайней обычной вершиной  $3b$ -цепи:



Ниже рассматривается первый случай, остальные аналогично.

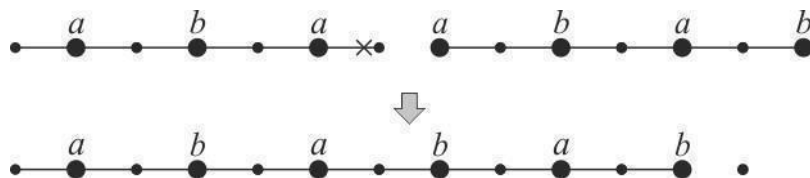
4.4'.  $3a+3b'=3$ .

4.5.  $2a+3=1a, 2b+3=1b$ . В  $3$ -цепи расклеить внешнее  $b$ -ребро и особую вершину склеить с крайней особой вершиной  $2a$ -цепи:



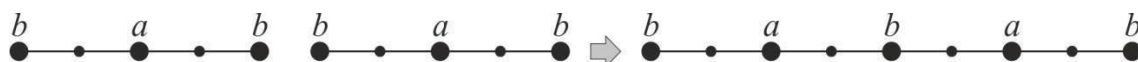
4.5'.  $2a'+3=1a$ .

4.6.  $3a+2=1a, 3b+2=1b$ . В  $3a$ -цепи расклеить внешнее ребро и особую вершину склеить с крайней особой вершиной  $2$ -цепи:



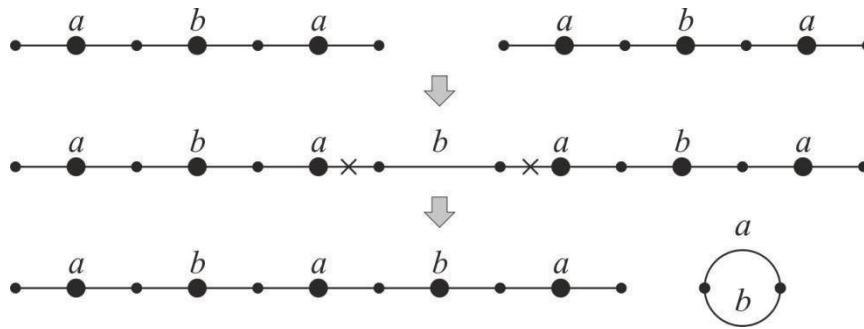
4.6'.  $3a+2'=1a, 3b'+2=1b$ .

4.7.  $2a+2a=2a, 2b+2b=2b$ . Склеить крайние особые вершины двух цепей:



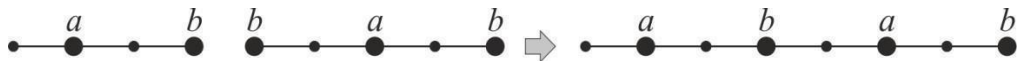
4.7'.  $2a'+2a=2a$ .

4.8.  $3a+3a=3a$ ,  $3b+3b=3b$ . Две крайние обычные вершины цепей соединить обычным ребром с последующим его вырезанием:



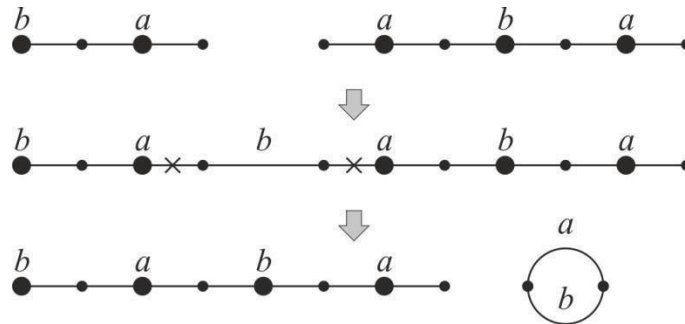
4.8'.  $3b'+3b=3b$ .

4.9.  $1a+2a=1a$ ,  $1b+2b=1b$ . Склеить крайние особые вершины двух цепей:



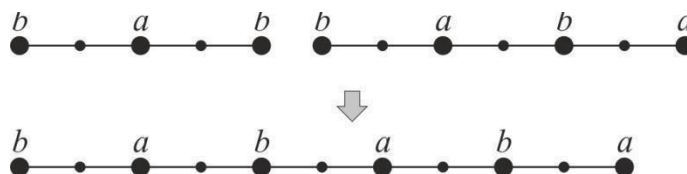
4.9'.  $1a+2a'=1a$ .

4.10.  $1a+3a=1a$ ,  $1b+3b=1b$ . Две крайние обычные вершины цепей соединить обычным ребром с последующим его вырезанием:



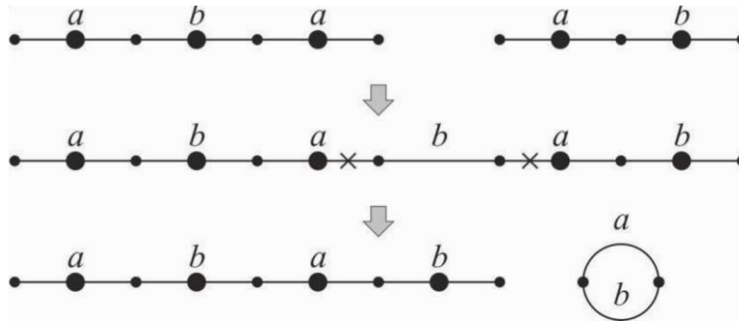
4.10'.  $1b+3b'=1b$ .

4.11.  $2a+2=2$ ,  $2b+2=2$ . Склеить крайние особые вершины двух цепей:



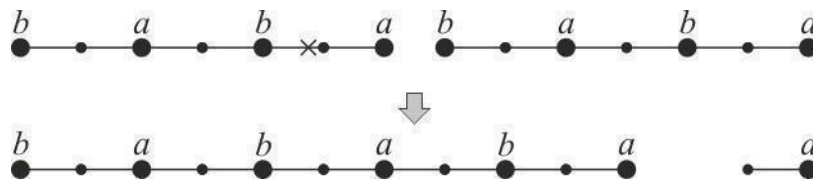
4.11'.  $2a'+2=2$ ,  $2a+2'=2$ ,  $2b+2'=2$ .

4.12.  $3a+3=3$ ,  $3b+3=3$ . Две крайние обычные вершины цепей соединить обычным ребром с последующим его вырезанием:



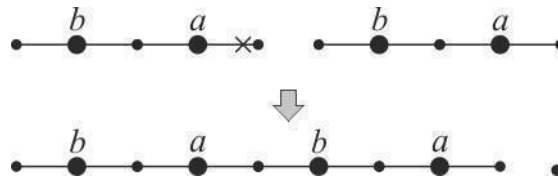
4.12'.  $3b'+3=3$ .

4.13.  $2+2=2+1'$ . Полупортная переклейка с отрезанием двух вершин 2-цепи (крайней особой  $a$ -вершины и соседней обычной) и склейкой образовавшегося края с крайней  $b$ -вершиной другой 2-цепи:



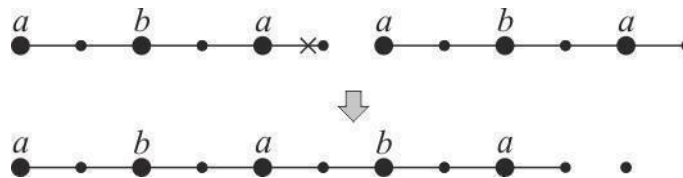
4.13'.  $2'+2=2+1'$ .

4.14.  $3+3=3$ . В 3-цепи расклеить внешнее  $a$ -ребро и образовавшийся край этой цепи склеить с  $b$ -краем другой 3-цепи:



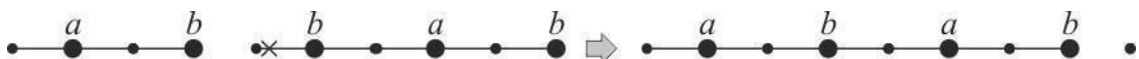
4.14'. Пустое действие.

4.15.  $1_a+1_a=1_a$ ,  $1_b+1_b=1_b$ ,  $1_b+1_c=1_b$  (определяется  $c=b$ ). В  $1_a$ -цепи расклеить внешнее ребро и особую вершину склеить с крайней особой вершиной другой  $1_a$ -цепи:



4.15'.  $1''+1_b=1_b$ ,  $1''+1_c=1_b$  (определяется  $c=b$ ).

4.16.  $1_a+1_b=1_a$ ,  $1_b+1_a=1_b$ ,  $1_a+1_c=1_a$  (определяется  $c=b$ ). В  $1_b$ -цепи расклеить внешнее ребро и особую вершину склеить с крайней особой вершиной  $1_a$ -цепи:



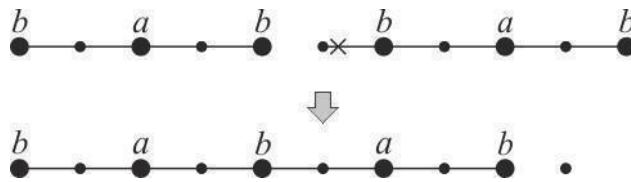
4.16'.  $1a+1''=1a$ .

4.17.  $1a+1_a=1a$ ,  $1b+1_b=1b$ ,  $1b+1_c=1b$  (определяется  $c=b$ ). В  $1a$ -цепи расклеить внешнее ребро и особую вершину склеить с крайней особой вершиной  $1a$ -цепи:



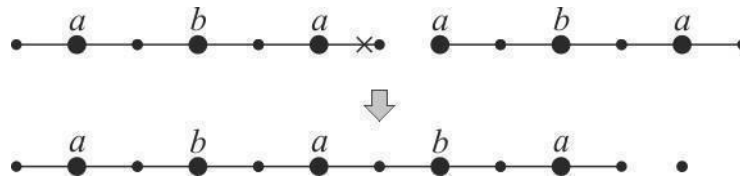
4.17'.  $1b+1''=1b$ .

4.18.  $2a+1_b=2a$ ,  $2b+1_a=2b$ ,  $2a+1_c=2a$  (определяется  $c=b$ ). В  $1_b$ -цепи расклеить внешнее ребро и особую вершину склеить с крайней особой вершиной  $2a$ -цепи:



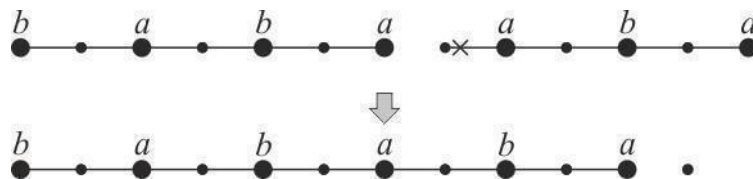
4.18'.  $2a'+1_b=2a$ ,  $2a+1''=2a$ ,  $2a'+1_c=2a$  (определяется  $c=b$ ).

4.19.  $3a+1_a=3a$ ,  $3b+1_b=3b$ ,  $3b+1_c=3b$  (определяется  $c=b$ ). В  $3a$ -цепи расклеить внешнее ребро и особую вершину склеить с крайней особой вершиной  $1a$ -цепи:



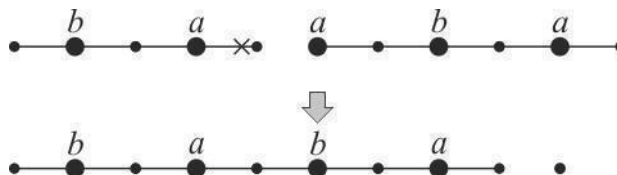
4.19'.  $3b'+1_b=3b$ ,  $3b+1''=3b$ ,  $3b'+1_c=3b$  (определяется  $c=b$ ).

4.20.  $2+1_a=2$ ,  $2+1_b=2$ ,  $2+1_c=2$  (определяется  $c=b$ ). В  $1_a$ -цепи расклеить внешнее ребро и особую вершину склеить с крайней особой вершиной  $2$ -цепи:



4.20'.  $2'+1_a=2$ ,  $2'+1_b=2$ ,  $2+1''=2$ ,  $2'+1_c=2$  (определяется  $c=b$ ).

4.21.  $3+1_a=3$ ,  $3+1_b=3$ ,  $3+1_c=3$  (определяется  $c=b$ ). В  $3$ -цепи расклеить внешнее ребро и особую вершину склеить с крайней особой вершиной  $1_a$ -цепи:



4.21'.  $3+1''=3$ .

4.22.  $1_a+1_c=1_a$  (определяется  $c=a$ ),  $1b+1_c=1b$  (определяется  $c=a$ ),  $1a+1_c=1a$  (определяется  $c=a$ ),  $2b+1_c=2b$  (определяется  $c=a$ ),  $3a+1_c=3a$  (определяется  $c=a$ ).

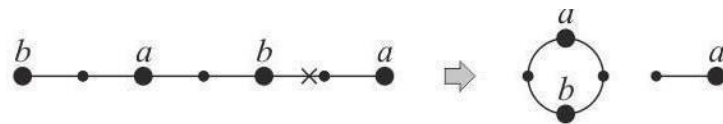
4.22'. Пустое действие.

4.23. Для оставшихся цепей типа  $1_c$  определяем  $c=b$  и выполняем  $1b+1b=1b$ .

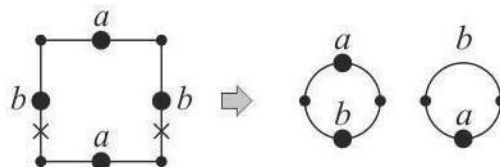
4.23'. Пустое действие.

4.24. Цепи, имеющие невисячее ребро, замыкаем в циклы склейкой (цепи типа  $2a$ ,  $2b$ ,  $3a$ ,  $3b$ ), полуторной переклейкой со слиянием особых вершин (цепи типа  $1a$ ,  $1b$ ,  $1c$ ,  $2$ ) или без слияния (цепи типа  $1a$ ,  $1b$ ,  $3$ ). При замыкании цепи типа  $1_c$  определяем  $c=b$ . При замыкании цепи типа  $2$  выбираем вариант со слиянием двух  $b$ -вершин, из образовавшейся цепи типа  $1'$  удаляем  $a$ -вершину, см. рисунок ниже. Из циклов, получившихся при замыкании цепей типа  $3a$  или  $3b$ , вырезаем обычные ребра. Затем снова применяем шаг 4.2.

4.24'. То же, что и 4.24.



**Шаг 5.** Удаляем изолированные особые вершины и петли. Из оставшихся цепей удаляем особые вершины. Из циклов длины больше 2 вырезаем 2-циклы так, чтобы происходило слияние двух  $b$ -вершин (соответственно, в 2-цикл включается  $a$ -вершина), см. рисунок ниже. Из 2-циклов удаляем особые вершины.



Конец описания алгоритма.

### Алгоритм построения дерева эволюции хромосомных структур

Построение эволюционного дерева, оптимально согласованного с полученной матрицей расстояний для данного множества хромосомных структур, выполняется простыми алгоритмами, подобными алгоритмам иерархической кластеризации UPGMA или NJ (neighbor joining).

## Установка и запуск в среде Windows программы ChromoGGL

Программа ChromoGGL написана на языке C++ и имеет интерфейс командной строки. Для скачивания доступна 32-битная версия исполняемого модуля. Она не требует специальной установки.

Для использования программы достаточно: скачать архив с программой *chromoggl.zip* и распаковать его в выбранную папку (например, *f:/Chromo*); запустить командный процессор Windows (Пуск > Выполнить > *cmd*) и перейти в папку с содержимым архива:

```
f:
```

```
cd f:/Chromo;
```

```
выполнить команду: chromoggl -h.
```

В случае успешного запуска на экране появляется краткая инструкция по запуску программы. В ином случае на экране появляется информация о возникшей ошибке. В случае успешного запуска рекомендуется протестировать программу на имеющихся в архиве примерах с помощью команд *run\_example\_X.bat*, где *X* пробегает значения 1 – 5.

## Входные данные программы ChromoGGL

### Файл структур

На вход программы ChromoGGL подается текстовый **файл хромосомных структур**, содержащий одну или более строк со следующими данными: *число* хромосомных структур; *имя* вида/штамма; *число* хромосом, включённых в структуру из этого вида, т.е. число путей и циклов в структуре; *пометка* (*L*), если хромосома линейная, или (*C*), если она циклическая; *число* генов в хромосоме; *последовательность* генов в хромосоме, записанная как имена генов со знаками «+» или «-», которые указывают на транскрибируемую цепь. Например, для следующих двух хромосомных структур их файл приведён ниже.

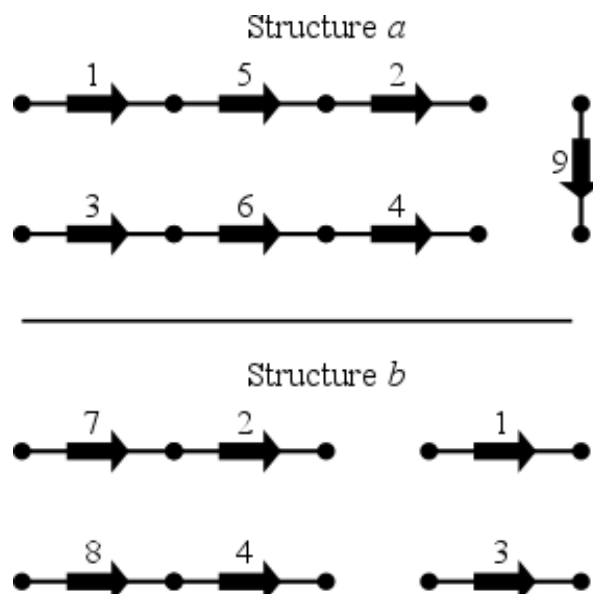




Рис. 1

Файл структур, показанных на рис. 1:

2; Structure\_a; 3; L3: +1+5+2; L3: +3+6+4; L1: +9; Structure\_b; 4; L2: +7+2; L2: +8+4; L1: +1;  
L1: +3

Общий граф этих двух структур показан на рис. 2, где точки увеличенного размера отмечают особые вершины, а обычного размера – обычные вершины:

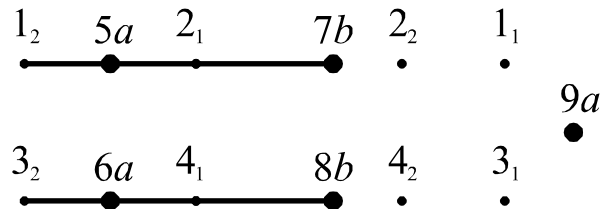


Рис. 2

Напомним: *обычной* называется вершина, соответствующая краю гена, представленного в обеих исходных структурах; *особой* называется вершина, соответствующая максимальному участку генов, представленных лишь в одной из исходных структур (сам участок называется *блоком*). У особой вершины имя структуры относится и к соседним рёбрам. Если особой вершине соответствуют несколько генов, то они выписываются подряд.

### Файл с ценами операций

Цены выписываются в строку со следующим форматом, где в качестве их значений могут использоваться любые положительные рациональные числа:

DP = ..., SP = ..., J = ..., C = ..., aD = ..., bD = ...

Пример файла – строки цен операций:

DP = 1.2, SP = 1.1, J = 1, C = 0.9, aD = 0.8, bD = 1.5

### Файл с весами расстояний разных типов

В этом файле указываются значения весов, с которыми при вычислении расстояния суммируются специальное (т.е. обычное) и брекпойнтовое расстояния между двумя структурами.

Пример файла:

```
[weights]
special=1
breakpoint=0.1
```

## Параметры программы ChromoGGL

В командной строке приводятся параметры программы, смысл которых поясняется ниже. Каждый параметр задаётся своим ключом (однобуквенным или длинным) и следующим после него аргументом. Обязательны параметры -c, -o, -m; для прочих параметров есть значение аргумента по умолчанию. Краткая справка о параметрах командной строки выдаётся при запуске программы с ключом -h (--help).

### Командная строка

- -c (--chrom) – имя файла хромосомных структур (может включать путь);
- -o (--oper) – имя файла с ценами отдельных операций (может включать путь);
- -m (--mode) – режим запуска, со значением аргумента dist, если решается задача 1, или tree, если решается задача 2;
- -r (--res) – путь к папке с результатами работы программы. По умолчанию используется текущая папка;
- -d (--coeff) – имя файла с ценами отдельных операций (может включать путь); по умолчанию берётся файл distance\_weights.ini в текущей папке;
- -p (--method) – алгоритм, применяемый для построения дерева, urgma или nj. По умолчанию используется алгоритм urgma;
- -f – вывод матрицы расстояний дополнительно в формате phylip, если аргумент равен 1. По умолчанию не выполняется.

### Выходные данные программы ChromoGGL в задаче 1

В задаче 1 выдаются два файла *joint\_graph* – описание общего графа и *shortest\_sequence* – описание искомой кратчайшей последовательности.

*Первый файл* имеет вид: заголовок *joint\_graph*, заголовок cycles (...), который включает указание числа циклов в общем графе, затем они перечисляются; заголовок paths (...) с указанием числа путей, которые затем перечисляются. Каждый цикл или путь записывается как линейная последовательность его вершин: общая вершина описывается как имя гена, сопровождаемое цифрой 1, если она – начало гена, и 2, если она – его конец. Особая вершина содержит имена генов, входящих в блок, и в конце имя структуры, к которой относится блок. Если разделение имён генов не является однозначным, то имена генов заключаются в круглые скобки. Ребро между обычными вершинами записывается знаком «\_», после него указывается имя структуры, к которой относится ребро. Ребро к особой вершине записывается тем же знаком; имя структуры опускается. В конце файла

указывается число  $a$ - и  $b$ -особых вершин в общем графе, сокращенно записываемых как “spnodes”.

*Пример первого файла* для общего графа, показанного на рис. 2:

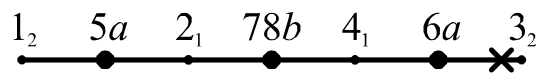
*Joint\_graph* cycles (0): paths (7):1.1; 2.2; 3.1; 4.2; 9a; 1.2\_5a\_2.1\_7b;  
3.2\_6a\_4.1\_8b  $a$ -spnodes: 3,  $b$ -spnodes: 2

*Второй файл* содержит кратчайшую последовательность операций, преобразующую общий граф к финальному виду. В его первой строке указаны: *число* операций в последовательности; её *суммарная цена*. Затем перечисляются сами операции. Каждая из них записывается в виде: компоненты, к которым применяются одна или несколько последовательных операций; сокращенные имена операций; полученный результат.

*Пример второго файла* для общего графа, показанного на рис. 2:

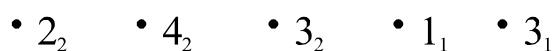
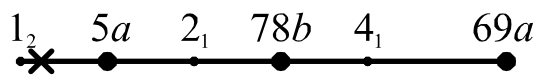
5; 5.4;  
[1.2\_5a\_2.1\_7b (L); 3.2\_6a\_4.1\_8b (L)] C [1.2\_5a\_2.1\_78b\_4.1\_6a\_3.2 (L)]

Здесь особые вершины  $7b$  и  $8b$  на двух путях объединяются в особую вершину  $78b$  на одном пути. Пометки (L) и (C) означают «цепь» и «цикл». Ниже изображен результат.



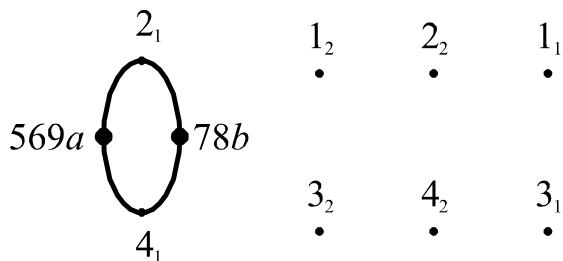
[1.2\_5a\_2.1\_78b\_4.1\_6a\_3.2 (L); 9a (L)] SP [3.2; 1.2\_5a\_2.1\_78b\_4.1\_6a (L)]

Отклейка вершины  $3.2$  и объединение особых вершин  $6a$  и  $9a$  в вершину  $69a$ . Результат ниже.



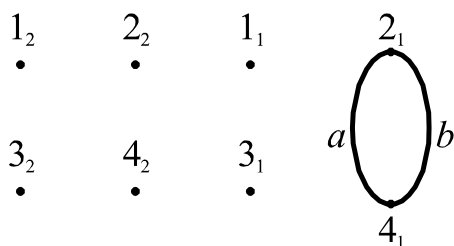
[1.2\_5a\_2.1\_78b\_4.1\_6a (L)] SP [2.1\_78b\_4.1\_569a\_2.1 (C); 1.2]

Отклейка вершины  $1.2$  и объединение особых вершин  $5a$  и  $69a$  в вершину  $569a$ . Ниже результат.



[2.1\_78b\_4.1\_569a\_2.1 (C)] aD, bD [2.1\_a4.1\_b2.1 (C)]

Удаление особых вершин 569a и 78b. Ниже результат.



Получен граф финального вида, алгоритм закончил работу.

## Примеры решения задачи 1

В архиве с программой находятся несколько примеров решения задачи 1, в том числе искусственный пример №0 – пара структур, показанная на рис. 1, и биологический – пара структур от хромосомы *Salmonella enterica* к хромосоме *Escherichia coli*. Для запуска программы на этих примерах необходимо выполнить соответствующие скрипты *run\_example\_1.bat* и *run\_example\_2.bat*, также лежащие в архиве.

### Искусственный пример 0

Исходные хромосомные структуры (файл data/artificial\_0.txt):

```
2;
Structure_a; 3; L3: +1+5+2; L3: +3+6+4; L1: +9;
Structure_b; 4; L2: +7+2; L2: +8+4; L1: +1; L1: +3;
```

Исходный общий граф для этих структур имеет вид:

```
Joint graph
Cycles (0):
Paths (7): 1.1; b_2.1_a_1.2; 2.2; 3.1; b_4.1_a_3.2; 4.2; a;
a-spnodes: 3, b-spnodes: 2
```

Найденные кратчайшие последовательности для разных наборов цен операций:

- для DP = 0.9, SP = 1.2, J = 1.1, C = 1, aD = 0.8, bD = 1.5 (файл цен config\_cir.ini):

```
5; 5.7;
[3.2_6a_4.1_8b(L), 1.2_5a_2.1_7b(L), 9a(L)] C, SP [3.2_6a_4.1_87b_2.1_59a(L), 1.2(L)]
[3.2_6a_4.1_87b_2.1_59a(L)] SP [2.1_659a_4.1_87b_2.1(C), 3.2(L)]
[2.1_659a_4.1_87b_2.1(C)] aD [2.1_a4.1_87b_2.1(C)]
[2.1_a4.1_87b_2.1(C)] bD [2.1_a4.1_b2.1(C)]
```

- для DP = 1, SP = 1, J = 1, C = 1, aD = 1, bD = 1 (файл цен config\_equal.ini):

5; 5;

[3.2\_6a\_4.1\_8b(L), 1.2\_5a\_2.1\_7b(L), 9a(L)] C, SP [3.2\_6a\_4.1\_87b\_2.1\_59a(L), 1.2(L)]

[3.2\_6a\_4.1\_87b\_2.1\_59a(L)] SP [2.1\_659a\_4.1\_87b\_2.1(C), 3.2(L)]

[2.1\_659a\_4.1\_87b\_2.1(C)] aD [2.1\_a4.1\_87b\_2.1(C)]

[2.1\_a4.1\_87b\_2.1(C)] bD [2.1\_a4.1\_b2.1(C)]

- для DP = 1.2, SP = 1.1, J = 1, C = 0.9, aD = 0.8, bD = 1.5 (файл цен config\_lin.ini):

5; 5.4;

[3.2\_6a\_4.1\_8b(L), 1.2\_5a\_2.1\_7b(L), 9a(L)] C, SP [3.2\_6a\_4.1\_87b\_2.1\_59a(L), 1.2(L)]

[3.2\_6a\_4.1\_87b\_2.1\_59a(L)] SP [2.1\_659a\_4.1\_87b\_2.1(C), 3.2(L)]

[2.1\_659a\_4.1\_87b\_2.1(C)] aD [2.1\_a4.1\_87b\_2.1(C)]

[2.1\_a4.1\_87b\_2.1(C)] bD [2.1\_a4.1\_b2.1(C)]

Приведенные кратчайшие последовательности лежат вместе с общим графом в папке results/artificial\_0 в подпапках cir, equal и lin, соответственно.

### *Другие искусственные примеры*

Также программа тестировалась на наборах искусственных данных №№ 1–8, которые содержатся в дистрибутивном архиве. Для запуска программы на любом из этих примеров скрипт run\_example\_1.bat запускается с параметром – номером примера. Например, для вычисления примера №4 следует ввести команду:

```
run_example_1.bat 4.
```

Исходные хромосомные структуры для примера с номером X содержатся в файле data/artificial\_X.txt. Найденные кратчайшие последовательности для тех же трёх наборов цен операций лежат вместе с общим графом в папке results/artificial\_X в подпапках cir, equal и lin, соответственно и здесь для краткости не приводятся.

### *Биологический пример*

Исходные хромосомные структуры (файл data/Salmonella\_Escherichia.txt):

2;

Salmonella\_enterica; 1; C47: -rpl25-rpsV-rplT-rpsA-rplS-rpsP+rpsU-rplU-rps1-rplM-rplQ-rpsD-rps11-rps13-rplO-rpsE-rplR-rplF-rpsH-rpsN-rplE-rplX-rpl14-rpsQ-rpl29-rplP-rpsC-rplV-rpsS-rplB-rplW-rplD-rpl3-tufA-rps7-rps12+rpoZ+tufA+rpl11+rplA+rplJ+rplL+rpoB+rpos1+rpsF+rpsR+rpl1;

Escherichia\_coli; 1; C56: -rpsT+rpsB+rpsA-rpsV-rplT+rplY-rpoE-rplS-rpsP-rpoS+rpsU+rpoD-rpsO-rplU+rpoN+rpoN-rps1-rplM-rplQ-rpoA-rpsD-rpsK-rpsM-rplO-rpsE-rplR-rplF-rpsH-rpsN-rplE-rplX-rplN-rpsQ-rplP-rpsC-rplV-rpsS-rplB-rplW-rplD-rplC-rpsJ-tufA-rpsG-rpsL-rpoH-rpoZ+rplK+rplA+rplJ+rplL+rpoB+rpoC+rpsF+rpsR+rpl1;

Исходный общий граф для этих структур:

Joint graph

Cycles (33): rpsV.1\_brplT.2\_arpsV.1;

rpsV.2\_brpsA.2\_arplT.1\_rplYrpoEb\_rplS.2\_arpsA.1\_rpsBrpsTb\_rpl1.2\_rpl25a\_rpsV.2; rplS.1\_brpsP.2\_arplS.1;

rpsP.1\_rpoSb\_rpsU.1\_arpsP.1; rpsU.2\_rpoDrpsOb\_rplU.2\_arpsU.2; rplU.1\_rpoNrpoNb\_rps1.2\_arplU.1;

rps1.1\_brplM.2\_arps1.1; rplM.1\_brplQ.2\_arplM.1; rplQ.1\_rpoAb\_rpsD.2\_arplQ.1;

rpsD.1\_rpsKrpsMb\_rplO.2\_rps13rps11a\_rpsD.1; rplO.1\_brpsE.2\_arplO.1; rpsE.1\_brplR.2\_arpsE.1;

rplR.1\_brplF.2\_arplR.1; rplF.1\_brpsH.2\_arplF.1; rpsH.1\_brpsN.2\_arpsH.1; rpsN.1\_brplE.2\_arpsN.1;

rplE.1\_brplX.2\_arplE.1; rplX.1\_rplNb\_rpsQ.2\_rpl14a\_rplX.1; rpsQ.1\_brplP.2\_rpl29a\_rpsQ.1;  
 rplP.1\_brpsC.2\_arplP.1; rpsC.1\_brplV.2\_arpsC.1; rplV.1\_brpsS.2\_arplV.1; rpsS.1\_brplB.2\_arpsS.1;  
 rplB.1\_brplW.2\_arplB.1; rplW.1\_brplD.2\_arplW.1; rplD.1\_rplCrpsJb\_tufA.2\_rpl3a\_rplD.1;  
 tufA.1\_rpsGrpsLrpoHb\_rpoZ.2\_tufArpl11a\_rplA.1\_rplKb\_rpoZ.1\_rps12rps7a\_tufA.1; rplA.2\_brplJ.1\_arplA.2;  
 rplJ.2\_brplL.1\_arplJ.2; rplL.2\_brpoB.1\_arplL.2; rpoB.2\_rpoCb\_rpsF.1\_rpoC1a\_rpoB.2; rpsF.2\_brpsR.1\_arpsF.2;  
 rpsR.2\_brpl1.1\_arpsR.2;  
 Paths (0):  
 a-spnodes: 8, b-spnodes: 12

Найденные кратчайшие последовательности для тех же трёх наборов цен операций лежат вместе с общим графом в папке results/Salmonella\_Escherichia в подпапках *cir*, *equal* и *lin*, соответственно. Например, для первого набора цен ( $DP = 0.9$ ,  $SP = 1.2$ ,  $J = 1.1$ ,  $C = 1$ ,  $aD = 0.8$ ,  $bD = 1.5$ ) построена следующая кратчайшая последовательность:

21; 21.6;  
 [rpsV.2\_brpsA.2\_arplT.1\_b\_rplS.2\_arpsA.1\_b\_rpl1.2\_a\_rpsV.2(C)] DP [rpsV.2\_brpsA.2\_arpsV.2(C),  
 rplT.1\_b\_rplS.2\_arpsA.1\_b\_rpl1.2\_a\_rplT.1(C)] , step:  
 [rplT.1\_b\_rplS.2\_arpsA.1\_b\_rpl1.2\_a\_rplT.1(C)] DP [rplS.2\_arpsA.1\_brplS.2(C), rplT.1\_b\_rpl1.2\_a\_rplT.1(C)] , step:  
 [rpoB.2\_a\_rpsF.1\_b\_rpoB.2(C), tufA.1\_a\_rpoZ.1\_b\_rplA.1\_a\_rpoZ.2\_b\_tufA.1(C)] DP, DP  
 [tufA.1\_a\_rpoB.2\_b\_rplA.1\_a\_rpoZ.2\_b\_tufA.1(C), rpoZ.1\_brpsF.1\_arpoZ.1(C)] , step: 4.2  
 [tufA.1\_a\_rpoB.2\_b\_rplA.1\_a\_rpoZ.2\_b\_tufA.1(C), rplD.1\_a\_tufA.2\_b\_rplD.1(C)] DP, DP  
 [rplD.1\_a\_tufA.1\_b\_rpoZ.2\_a\_rplA.1\_b\_rplD.1(C), tufA.2\_brpoB.2\_atufA.2(C)] , step:  
 [rplD.1\_a\_tufA.1\_b\_rpoZ.2\_a\_rplA.1\_b\_rplD.1(C), rplX.1\_a\_rpsQ.2\_b\_rplX.1(C)] DP, DP  
 [rplX.1\_a\_rplD.1\_b\_rplA.1\_a\_rpoZ.2\_b\_rplX.1(C), rpsQ.2\_btufA.1\_arpsQ.2(C)] , step:  
 [rplX.1\_a\_rplD.1\_b\_rplA.1\_a\_rpoZ.2\_b\_rplX.1(C), rpsD.1\_a\_rplO.2\_b\_rpsD.1(C)] DP, DP  
 [rpsD.1\_a\_rplX.1\_b\_rpoZ.2\_a\_rplA.1\_b\_rpsD.1(C), rplO.2\_brplD.1\_arplO.2(C)] , step:  
 [rpsD.1\_a\_rplX.1\_b\_rpoZ.2\_a\_rplA.1\_b\_rpsD.1(C), rplT.1\_b\_rpl1.2\_a\_rplT.1(C)] DP, DP  
 [rpsD.1\_b\_rplA.1\_a\_rpoZ.2\_b\_rpl1.2\_a\_rpsD.1(C), rplT.1\_brplX.1\_arplT.1(C)] , step:  
 [rpsD.1\_b\_rplA.1\_a\_rpoZ.2\_b\_rpl1.2\_a\_rpsD.1(C)] DP [rpsD.1\_b\_rpl1.2\_a\_rpsD.1(C), rplA.1\_a\_rpoZ.2\_brplA.1(C)]  
 , step:  
 [rpsP.1\_arpsU.1\_b\_rpsP.1(C)] bD [rpsP.1\_arpsU.1\_brpsP.1(C)] , step:  
 [rpsU.2\_arplU.2\_b\_rpsU.2(C)] bD [rpsU.2\_arplU.2\_brpsU.2(C)] , step:  
 [rplU.1\_arps1.2\_b\_rplU.1(C)] bD [rplU.1\_arps1.2\_brplU.1(C)] , step:  
 [rplQ.1\_arpsD.2\_b\_rplQ.1(C)] bD [rplQ.1\_arpsD.2\_brplQ.1(C)] , step:  
 [rpsQ.1\_a\_rplP.2\_brpsQ.1(C)] aD [rpsQ.1\_arplP.2\_brpsQ.1(C)] , step:  
 [rplA.1\_a\_rpoZ.2\_brplA.1(C)] aD [rplA.1\_arpoZ.2\_brplA.1(C)] , step:  
 [rpsD.1\_b\_rpl1.2\_a\_rpsD.1(C)] bD [rpsD.1\_brpl1.2\_a\_rpsD.1(C)] , step:  
 [rpsD.1\_brpl1.2\_a\_rpsD.1(C)] aD [rpsD.1\_brpl1.2\_arpsD.1(C)] , step:

## Выходные данные программы ChromoGGL в задаче 2

В задаче 2 алгоритм выдает в папку для записи результатов следующие файлы: *distance.tsv* – матрица попарных расстояний между хромосомными структурами, в формате TSV (значения, разделённые символом табуляции) для импорта в Excel; *tree.tre* – дерево эволюции хромосомных структур, оптимально согласованное с матрицей, в скобочном формате Newick, в том числе с указанием длин рёбер (в текущей версии – только при построении дерева методом neighbor joining). Также при значении 1 параметра -f выдается еще два файла: *infile* – матрица попарных расстояний в формате phylip, предназначенная для использования в пакетах phylip, например, построения дерева с помощью программы Neighbor (<http://evolution.genetics.washington.edu/phylip/doc/neighbor.html>). В этой матрице

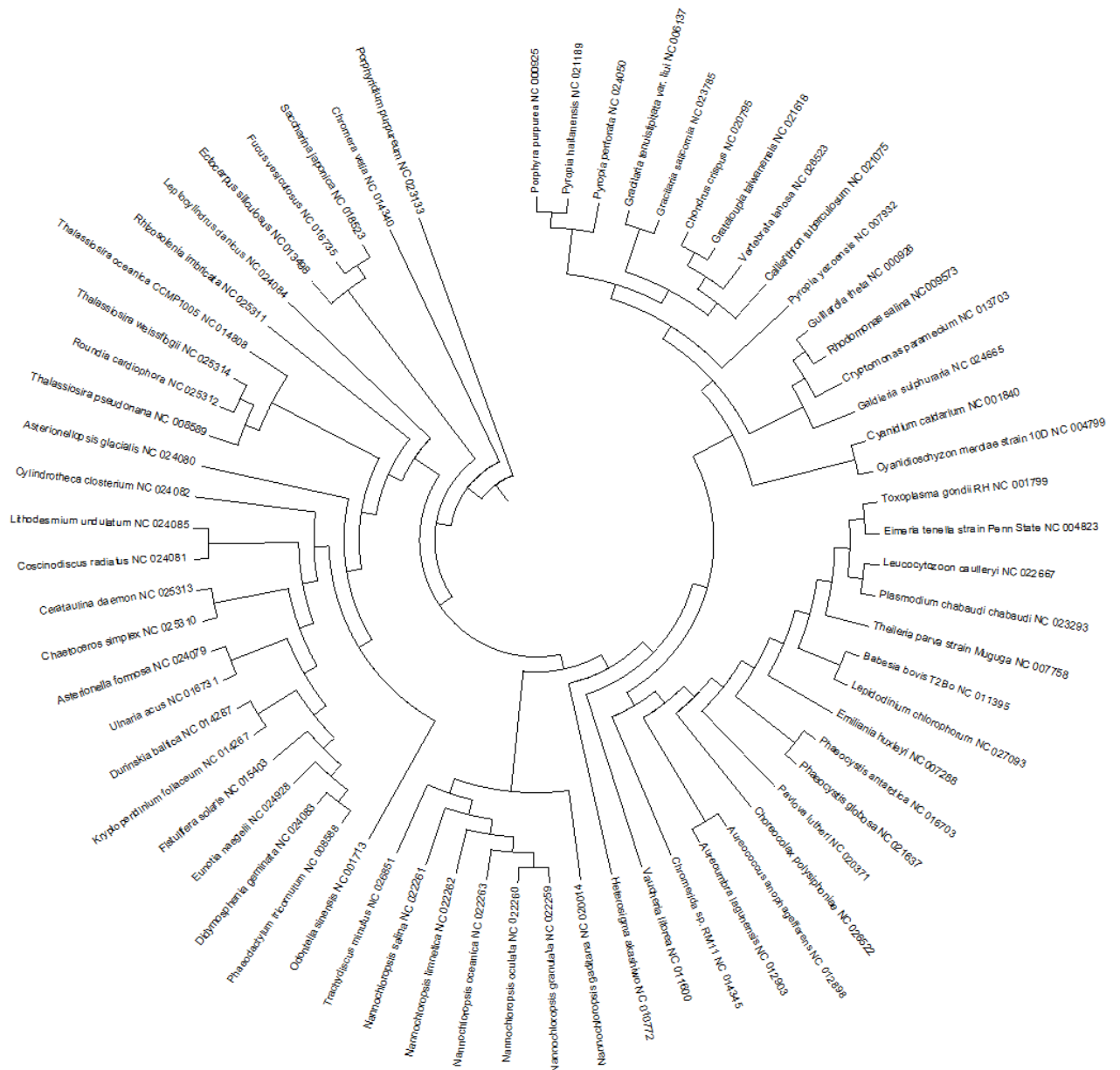
названия видов заменяются номерами из Генбанка, так как rhyUp поддерживает только короткие имена. Соответствие номеров и имен записывается в файл *map.txt*.

## Примеры решения задачи 2

### Пластиды родофитной ветви

Первый пример множества хромосомных структур содержит 66 пластид родофитной ветви. Исходные хромосомные структуры содержатся в файле *data/rhodophytic\_branch.txt*. Для запуска программы на этом примере необходимо выполнить скрипт *run\_example\_3.bat*.

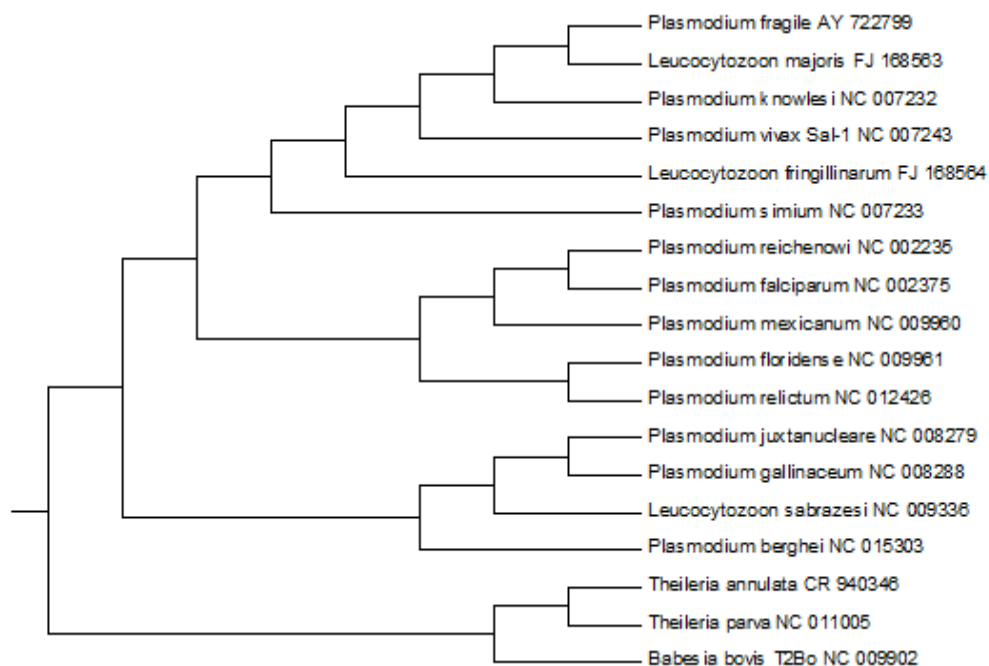
Построенная в результате запуска матрица попарных расстояний содержится в файле *results/rhodophytic\_branch/distance.tsv*, а дерево в формате Newick – в файле *results/rhodophytic\_branch/tree.tre* (ниже приводится его визуализация).



## Митохондрии класса *Aconoidasida*

Второй пример множества хромосомных структур содержит 18 митохондрий класса *Aconoidasida*. В нём представлены как кольцевые, так и линейные хромосомы. Исходные хромосомные структуры содержатся в файле data/mitochondria.txt. Для запуска программы на этом примере необходимо выполнить скрипт run\_example\_4.bat.

Построенная в результате запуска матрица попарных расстояний содержится в файле results/mitochondria/distance.tsv, а дерево в формате Newick – в файле results/mitochondria/tree.tre (ниже приводится его визуализация).



## Вспомогательная утилита gbk\_converter

Утилита gbk\_converter служит для преобразования генома из формата GenBank (.gbk, .gbff) в формат программы ChromoGGL.

### Входные данные gbk\_converter:

1) Список файлов: каждый файл в формате gbff содержит фрагмент генома определенного организма из Генбанка. Пример начальной части файла:

```
LOCUS       CP004015                3837060 bp    DNA     circular BCT 31-JAN-2014
DEFINITION  Rhizobium tropici CIAT 899, complete genome.
ACCESSION   CP004015
VERSION     CP004015.1
DBLINK      BioProject: PRJNA42391
            BioSample: SAMN02604105
KEYWORDS    .
SOURCE      Rhizobium tropici CIAT 899
  ORGANISM  Rhizobium tropici CIAT 899
            Bacteria; Proteobacteria; Alphaproteobacteria; Rhizobiales;
            Rhizobiaceae; Rhizobium/Agrobacterium group; Rhizobium.
REFERENCE   1 (bases 1 to 3837060)
  AUTHORS   Ormeno-Orrillo,E., Menna,P., Almeida,L.G., Ollero,F.J.,
            Nicolas,M.F., Pains Rodrigues,E., Shigeyoshi Nakatani,A., Silva
            Batista,J.S., Oliveira Chueire,L.M., Souza,R.C., Ribeiro
            Vasconcelos,A.T., Megias,M., Hungria,M. and Martinez-Romero,E.
```



```

TITLE      Genomic basis of broad host range and environmental adaptability of
           Rhizobium tropici CIAT 899 and Rhizobium sp. PRF 81 which are used
           in inoculants for common bean (Phaseolus vulgaris L.)
JOURNAL    BMC Genomics 13 (1), 735 (2012)
PUBMED     23270491
REMARK     Publication Status: Online-Only
REFERENCE  2 (bases 1 to 3837060)
AUTHORS    Ormeno-Orrillo,E. and Martinez-Romero,E.
TITLE      Direct Submission
JOURNAL    Submitted (28-DEC-2012) Centro de Ciencias Genomicas, Universidad
           Nacional Autonoma de Mexico, Av. Universidad s/n, Col. Chamilpa,
           Cuernavaca, Morelos 62210, Mexico
COMMENT    Source DNA is available from Esperanza Martinez-Romero
           (emartine@ccg.unam.mx).
FEATURES   Location/Qualifiers
           source          1..3837060
                           /organism="Rhizobium tropici CIAT 899"
                           /mol_type="genomic DNA"
                           /strain="CIAT 899"
                           /isolation_source="Phaseolus vulgaris root nodule"
                           /db_xref="taxon:698761"
                           /country="Colombia"
                           /note="type strain of Rhizobium tropici"
           gene            427..1248
                           /locus_tag="RTCIAT899_CH00005"
           CDS              427..1248
                           /locus_tag="RTCIAT899_CH00005"
                           /codon_start=1
                           /transl_table=11
                           /product="putative phosphotransferase"
                           /protein_id="AGB69434.1"
                           /translation="MENRTSFFHLHLISDSTGETLISAGRAASVQFHASQPIEHVYPL
IRNRKQLLPVLEAIDHSPGIVLYTIVDRELADFIAERCREMGVPSVNVLEPVMNVFQT
YLGASRRRVGAQHVMNADYFARIEALNFTMDHDDGQMPDDYDDADVVIIGISRTSKT
PTSIYLANRGIKTANIPIVHGVPPLPESLARATKPLIVGLVATTDTRISQVRENRI LGTT
PGFDRGGYTDRAAISEELKYARSLCARHNWPIIDVTRRSIEETAAAI VALRPKLR"
           gene            1258..1881
           .....
```

2) файл masks.txt со списком типов генов, которые следует отобразить в структуру. Например:

```

psa
psb
rpl
rps
rpo
```

**Командная строка:** python gbk\_converter.py -i [path to directory with gbff files] -m [path to file with masks] -o [path to file with output sequences]

**Выходные данные gbk\_converter:** Файл result.txt с геномными структурами для каждого вида. Например:

```

Bradyrhizobium japonicum USDA 110 chromosome, complete genome; 1; C53:
-rpoH3+rpsU+rplU+rpsP+rplS-rplT+rpoN2+rpsA-rpsO-rpsT+rpoN1-rplI-rpsR-rpsF-rpsB-rpsI-
rplM+rpoH1-rplQ-rpoA-rpsK-rpsM-rplO-rpsE-rplR-rplF-rpsH-rpsN-rplE-rplX-rplN-rpsQ-rplP-
rpsC-rplV-rpsS-rplB-rplW-rplD-rplC-rpsJ-rpsG-rpsL-rpoC-rpoB-rplL-rplJ-rplA-
rplK+rpsD+rpsU+rpoH2-rplY;
Rhizobium tropici CIAT 899, complete genome; 1; C41: -rpsA-rpsO+rplT+rpsT-rpoNch-rplI-rpsR-
rpsI-rplM+rplK+rplA+rplL+rpoB+rpoC+rpsL+rpsG+rpsJ+rplC+rplD+rplW+rplB+rplV+rpsC+rplP-
+rpsQ+rplN+rplX+rpsH+rplF+rplR+rpsE+rplO+rpsM+rpsK+rpoA+rplQ+rpsB-rpsD+rpoH1-
rpoH2+rplU;
```

## Рекомендуемые стандартные программы

Для просмотра файлов рекомендуется использовать стандартные программы.

Ссылки удобно открывать в новой вкладке, используя, например, браузер Google Chrome.

Microsoft Excel – для просмотра матриц. Они записываются в формате tsv (tab separated values). В качестве разделителя в числах используется точка. Для отображения матриц нужно в настройках задать соответствующий разделитель, который находится в меню Файл>Параметры>Дополнительно>Использовать системные разделители.

TreeViewX (<https://code.google.com/p/treeviewx/>) – для просмотра деревьев, которые используют формат Newick. Любой текстовый редактор для просмотра общих графов, кратчайших последовательностей, файлов с описанием исходных листовых хромосомных структур и оптимальной расстановкой структур по дереву.

## Литература

[1] К.Ю. Горбунов, Р.А. Гершгорин, В.А. Любецкий. Перестройка и реконструкция хромосомных структур. *Молекулярная биология*, 2015, том. 49, №3, стр. 372–383.

[Англ. перевод: K.Yu. Gorbunov, R.A. Gershgorin, V.A. Lyubetsky. Rearrangement and Inference of Chromosome Structures. *Molecular Biology*, 2015, Vol. 49, No. 3, pp. 327–338.]

[2] V.A. Lyubetsky, R.A. Gershgorin, A.V. Seliverstov, K.Yu. Gorbunov. Algorithms for reconstruction of chromosomal structures // *BMC Bioinformatics*, 2016, Vol. 17, Article 40, 23 pp. DOI: 10.1186/s12859-016-0878-z.

[3] V.A. Lyubetsky, R.A. Gershgorin, K.Yu. Gorbunov. Chromosome structures: reduction of certain problems with unequal gene content and gene paralogs to integer linear programming // *BMC Bioinformatics*, 2017, Vol. 18, Article 537, 18 pp. DOI: 10.1186/s12859-017-1944-x.