



АЛГОРИТЪМ CLUSTALW ЗА МНОЖЕСТВЕНО ПОДРАВНЯВАНЕ НА БИОЛОГИЧНИ СЕКВЕНЦИИ

ПРОФ. ПЛАМЕНКА БОРОВСКА

АЛГОРИТЪМ CLUSTALW ЗА МНОЖЕСТВЕНО ПОДРАВНЯВАНЕ НА БИОЛОГИЧНИ СЕКВЕНЦИИ

- Един от най-често използваните алгоритми за множествоно подравняване, който е базиран на прогресивното подравняване
- Основната идея - първоначално подравняване на най-тясно свързаните секвенции
- ClustalW е първият алгоритъм, при който се комбинира подравняването по двойки на секвенциите с глобалното подравняване с цел намаляване на изчислителното време.
- *Алгоритъмът ClustalW има сложност $O(N^2)$* поради използването на метода на присъединяване на съседите (neighbor-joining method).

АЛГОРИТЪМ CLUSTALW ЗА МНОЖЕСТВЕНО ПОДРАВНЯВАНЕ НА БИОЛОГИЧНИ СЕКВЕНЦИИ

- Методът *Clustal V* (предложен от Higgins и Sharp, 1989 г.) групира секвенциите в клъстери на основата на анализ на еволюционните дистанции на всяка двойка секвенции.
- Клъстерите се подреждат по двойки, след това колективно като групи секвенции с цел получаване на глобално подравняване.
- За полученото подравняване на секвенциите се прилага методът за присъединяване на съсед за реконструкция на филогенетичното дърво.
- Предназначението на метода *Clustal W* (Thompson, 1994 г.) е да генерира по-точни подреждания от метода *Clustal V* за секвенции с големи еволюционните дистанции.

CLUSTALW

- Той е евристичен алгоритъм и осигурява оптимално или най-често близко до оптималното подравняване.
- Въпреки недостатъците за качеството на решение, той е широко използван поради неговите предимства, а именно:
 - (1) приемливо време за подравняване на големи набори дълги секвенции, и
 - (2) направляващото дърво намалява чувствителността на подравняването към зашумеността на данните

CLUSTALW

- Clustalw е широко използван софтуерен продукт за подравняване на множество ДНК, РНК или протеинови секвенции.
- Резултатите показват общия произход в еволюцията на тези секвенции.
- ClustalW, в сравнение с други MSA алгоритми, се изпълнява като един от най-бързите, като все още поддържа приемливо ниво на точност.
- Точността за ClustalW е по-ниска при тестване срещу MAFFT, T-Coffee, Clustal Omega и други MSA софтуерни инструменти, но при него изискванията за RAM памет са най-малки.
- Има актуализации и подобрения на алгоритъма, напр., в ClustalW2, които водят до повишаване на точността, при запазване на висока скорост на изчисленията.

АЛГОРИТЪМЪТ CLUSTALW СЕ СЪСТОИ ОТ ТРИ ОСНОВНИ СЪПКИ

- 1. Подравняване по двойки
(pairwise alignment)*
- 2. Конструирание на направляващо дърво
(guide tree)*
- 3. Множествено подравняване*

CLUSTALW – СЪПКА 1: ПОДРАВНЯВАНЕ ПО ДВОЙКИ (PAIRWISE ALIGNMENT)

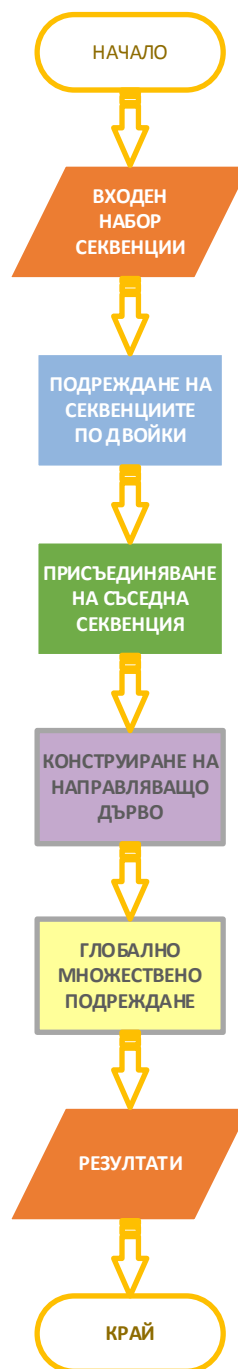
- Изчисляват се разликите между двойките секвенции чрез подравняване на секвенциите по двойки.
- Пресмята се дистанцията (разликите) между всеки две секвенции.
- За всяка двойка секвенции се извършва сравнение на двойките символи, при които няма празна позиция.
- Разстоянието се пресмята, като се раздели броят на двойките, при които съществува различие, на общия брой двойки.

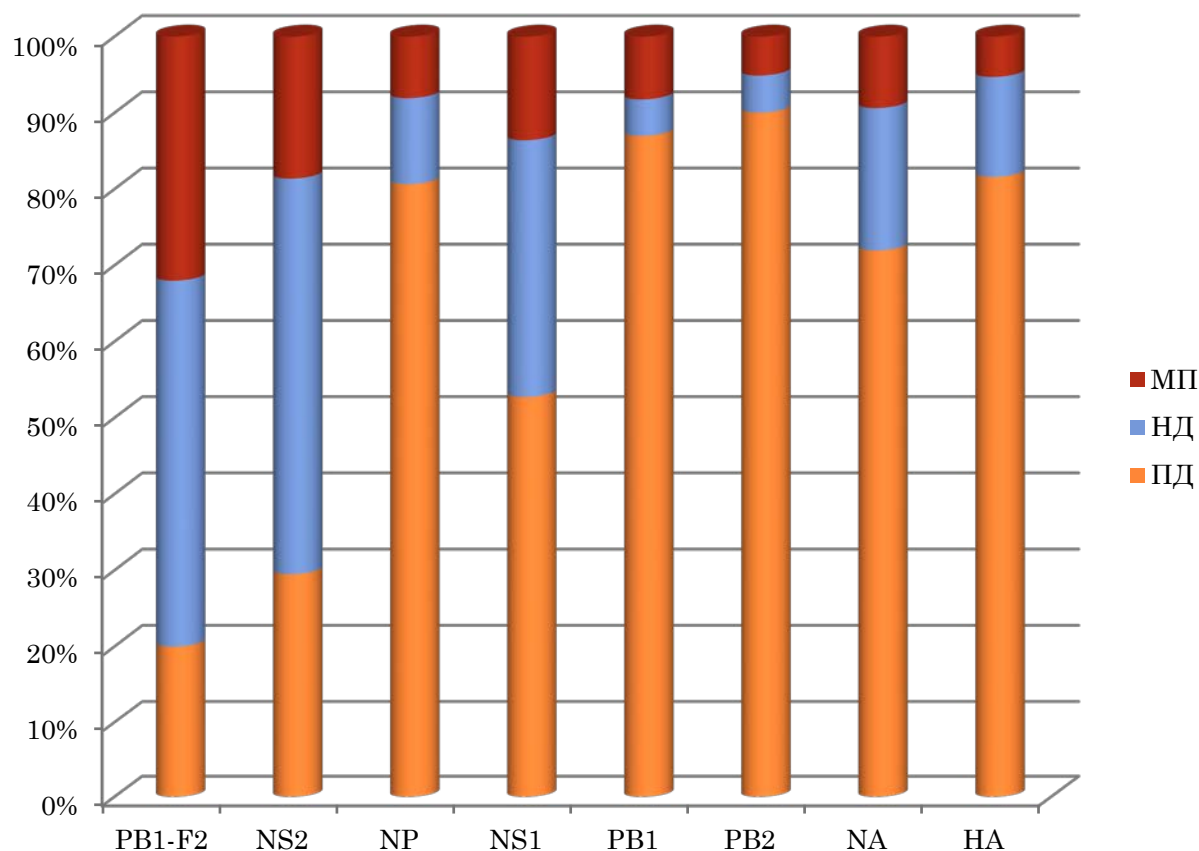
CLUSTALW – СТЬПКА 2: КОНСТРУИРАНЕ НА НАПРАВЛЯВАЩО ДЪРВО

- На базата на матрицата на разстоянията се построява направляващо дърво.
- При метода ClustalW се използва матрицата на дистанциите и методът на присъединяване на съседна секвенция (*Neighbor Joining*) за да се изгради направляващото дърво (guide tree).
- Комбинират се резултатите от подравняването по двойки, като се започне от най-близко свързаните секвенции и се завърши с най-отдалечените.
- Процесът започва от върховете до достигане на корена на дървото.

CLUSTALW – СТЬПКА 3: МНОЖЕСТВЕНО ПОДРАВНЯВАНЕ

- Секвенциите се подреждат на базата на направляващото дърво.
- Започва се с подравняването на двойките, които имат най-голямо сходство.





ОТНОСИТЕЛНИ ДЯЛОВЕ НА ВРЕМЕНАТА ЗА ПАРАЛЕЛНО ИЗПЪЛНЕНИЕ НА РАЗЛИЧНИТЕ ФАЗИ ПРИ МНОЖЕСТВЕНО ПОДРАВНЯВАНЕ С АЛГОРИТЪМА CLUSTALW НА 8-ТЕ СЕГМЕНТА НА РНК НА ВИРУСА НА СВИНСКИЯ ГРИП АН1/Н1 НА 32 ЯДРА НА КОМПЮТЪРЕН КЛЪСТЕР:
 ПД - ПОДРАВНЯВАНЕ ПО ДВОЙКИ, НД – НАСОЧВАЩО ДЪРВО,
 МП – МНОЖЕСТВЕНО ПОДРАВНЯВАНЕ

ВИЗУАЛИЗАЦИЯТА НА РЕЗУЛТАТИТЕ

За улесняване на интерпретацията на резултатите, получени от приложението на алгоритъма ClustalW е необходимо изходните данни да съдържат следното:

- Визуализация на подреденото множество секвенции;
- Конструираното филогенетично дърво;
- Таблица на дистанциите между секвенциите;
- Визуализация и хистограма на консенсусните региони в секвенциите;
- Възможност за редактиране на секвенциите.

ВИЗУАЛИЗАЦИЯТА НА КОНСЕНСУСНИТЕ РЕГИОНИ В ПОДРАВНЕНОТО МНОЖЕСТВО СЕКВЕНЦИИ ПОСРЕДСТВОМ МНОГОЦВЕТНОТО ИМ МАРКИРАНЕ ПО ДИАПАЗОНИ

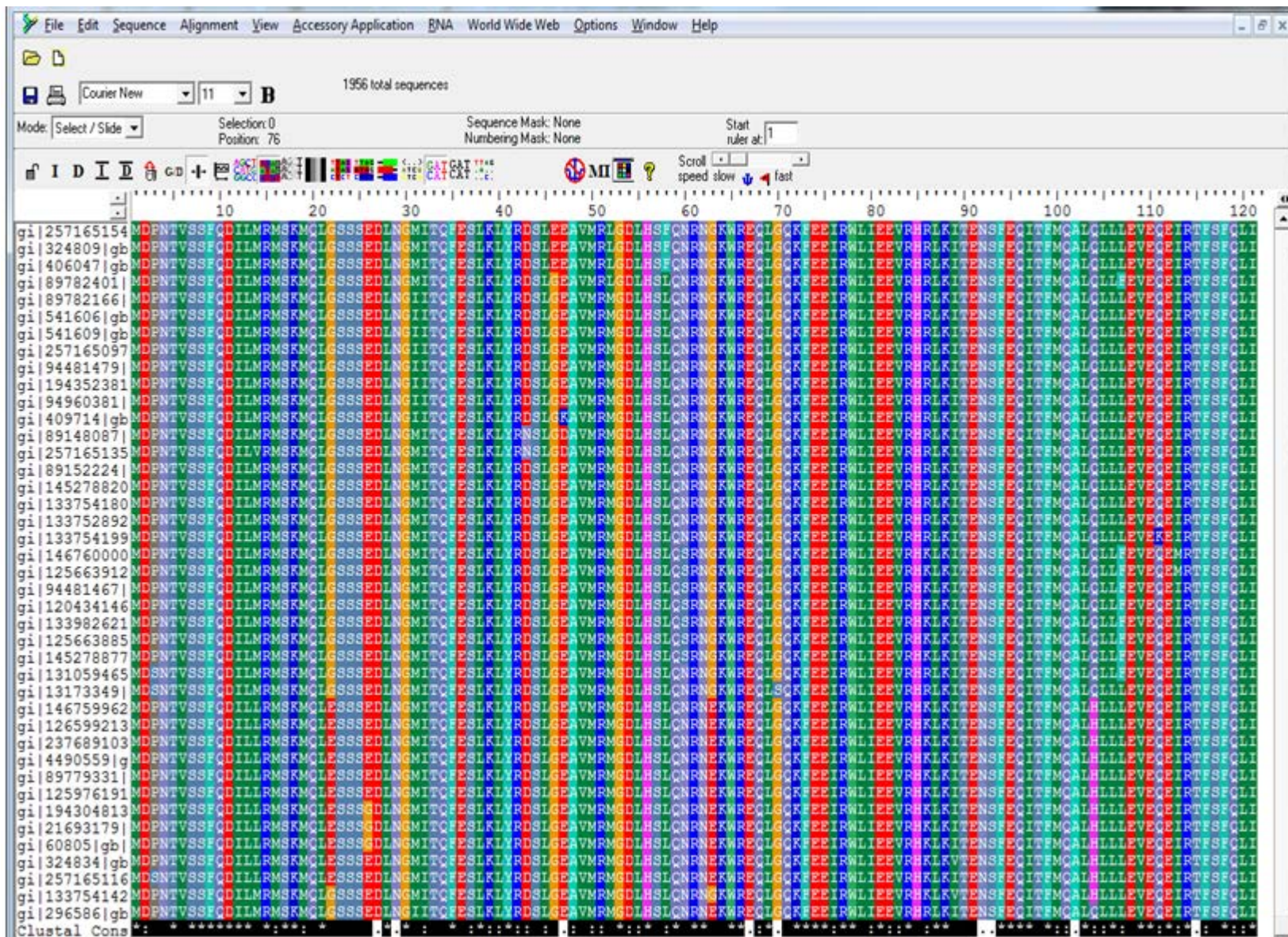
- Таблицата на дистанциите между секвенциите съдържа *процента на сходство за всяка двойка секвенции*, който се определя по следния начин:

Процент на сходство = (брой на съвпаденията \times 100) / общ брой на символите в региона (вкл. вътрешните празни позиции, не се отчитат тези в края на секвенциите)

- За визуализацията на консенсусните региони в подреденото множество секвенции в общия случай се използва многоцветното им маркиране по диапазони

ВИЗУАЛИЗАЦИЯТА НА КОНСЕНСУСНИТЕ РЕГИОНИ В ПОДРАВНЕНОТО МНОЖЕСТВО СЕКВЕНЦИИ ПОСРЕДСТВОМ МНОГОЦВЕТНОТО ИМ МАРКИРАНЕ ПО ДИАПАЗОНИ

Степен на сходство	Цвят
<1%	Черен
1-20%	Тъмно син
21-40%	Син
41-60%	Зелен
61-80%	Жълт
81-100%	Червен



проф. Боровска

ДОКУМЕНТИРАНЕ НА КОДА НА CLUSTALW С DOXYGEN

- документиране на кода в HTML формат посредством стандартния софтуерен инструмент за документиране на софтуер от сорс файлове Doxygen
- Doxygen е стандартен инструмент за генериране на документация от аотирани източници на C ++, но също така поддържа други популярни езици за програмиране като C, C #, PHP, Java, Python, и др.
- Doxygen е разработен под операционните системи Mac OS X и Linux, но също така е преносим.

ФУНКЦИОНАЛНОСТ НА DOXYGEN

- Може да генерира онлайн браузър за документация (в HTML) и/или офлайн справочно ръководство от набор от документиранни изходни файлове.
- Има и поддръжка за генериране на изходни данни в RTF (MS-Word), PostScript, хипервръзка PDF, компресиран HTML.
- Документацията се извлича директно от сорс кода, което много улеснява поддържането на актуална документацията в съответствие с изходния код.

ФУНКЦИОНАЛНОСТ НА DOXYGEN

- Доxygen може да се конфигурира за извличане на структура на кода от недокументирани файлове на сорс кода.
- Доxygen може също така да визуализира релациите между различните компоненти на кода чрез генериране на графики на зависимост, диаграми за наследяване и колаборативни диаграми, които се генерират автоматично.
- Доxygen може също така да се използва за създаване на документация на кода.
- В документацията е включен и програмния код, като по този начин е възможно винаги да се поддържа актуална версия на кода по отношение на направени промени.

ФУНКЦИОНАЛНОСТ НА DOXYGEN

- Поддържат се връзки между отделните документиранни сорс файлове и компонентите на софтуера – функции, структури от данни, променливи и др.
- Препратките позволяват да се проследяват връзките, извикванията на отделните функции, структурите данни, променливите, дефинициите и др. с цел улесняване на анализа на кода.
- Включени са метаданни за описание на кода.
- Също така е включена е възможност да се визуализират отношенията между различните елементи на програмната система чрез диаграми, представящи графи на зависимостите и връзките.

СТРУКТУРА, СЪДЪРЖАНИЕ И РЕФЕРЕНЦИИ НА ОСНОВНИЯ ФАЙЛ CLUSTALW.C

ClustalW: src/ClustalW.c File Reference - Opera

Open Save Print Find Home Tile Cascade Voice

Opera ClustalW: src/Clust... x

localhost/F:/Vessy/Documents%20%20Vessy/PRACE%202/Clustalw-mpi/ClustalDoc/ClustalW_8c.html

Kayak eBay Amazon My Opera Communi... Shopping Vienna Public Trans... Wikipedia

Home Index Contents Search Glossary Help First Previous Next Last Up Copyright Author

Find in page Find Next Voice Author Mode Show Images Fit to Width 100%

ClustalW 1.82

Main Page Data Structures Files

File List Globals

src/ClustalW.c File Reference

#include <stdio.h>
#include <string.h>
#include <stdlib.h>
#include <ctype.h>
#include <stdarg.h>
#include <signal.h>
#include <setjmp.h>
#include "clustalw.h"
#include "mpi.h"

Include dependency graph for ClustalW.c:

```
graph TD; src[src/ClustalW.c] --> stdio[stdio.h]; src --> string[string.h]; src --> stdlib[stdlib.h]; src --> ctype[ctype.h]; src --> stdarg[stdarg.h]; src --> signal[signal.h]; src --> setjmp[setjmp.h]; src --> clustalw[clustalw.h]; src --> mpi[mpi.h]; clustalw --> general[general.h];
```

Go to the source code of this file.

Functions

```
void * ckalloc (size_t)
void init_matrix (void)
void init_interface (void)
void fill_chartab (void)
int main (int argc, char **argv)
void fatal (char *msg,...)
void error (char *msg,...)
void warning (char *msg,...)
void info (char *msg,...)
char prompt_for_yes_no (char *title, char *prompt)
```

Variables

```
double ** tmat
char revision_level[] = "W (1.82)"
```

ПРОГРАМНИТЕ ФАЙЛОВЕ НА СОФТУЕРА CLUSTALW

clustalw.h - Main header file for ClustalW

general.h - General purpose header file

dayhoff.h - Table of estimated PAMS matrices

matrices.h – Scoring matrices

param.h – List of parameters

clustalw-mpi.c – function main()

parallel_compare.c – parallel pairwise alignment, calls function pair_align()

pairalign_new.c - parallel pairwise alignment

trees.c - phylogenetics tree calculating functions

malign_mpi_progressive.c -(function malign_mpi_progressive()) - full progressive multiple alignment

malign_mpi_pdiff.c (function malign_mpi_pdiff()) - full progressive multiple alignment

malign.c (function malign()) - sequential full progressive alignment

alnscore.c - calculation overall score for the alignment

calcgapcoeff.c - calculation gap penalty

calctree.c - Calculation guide tree

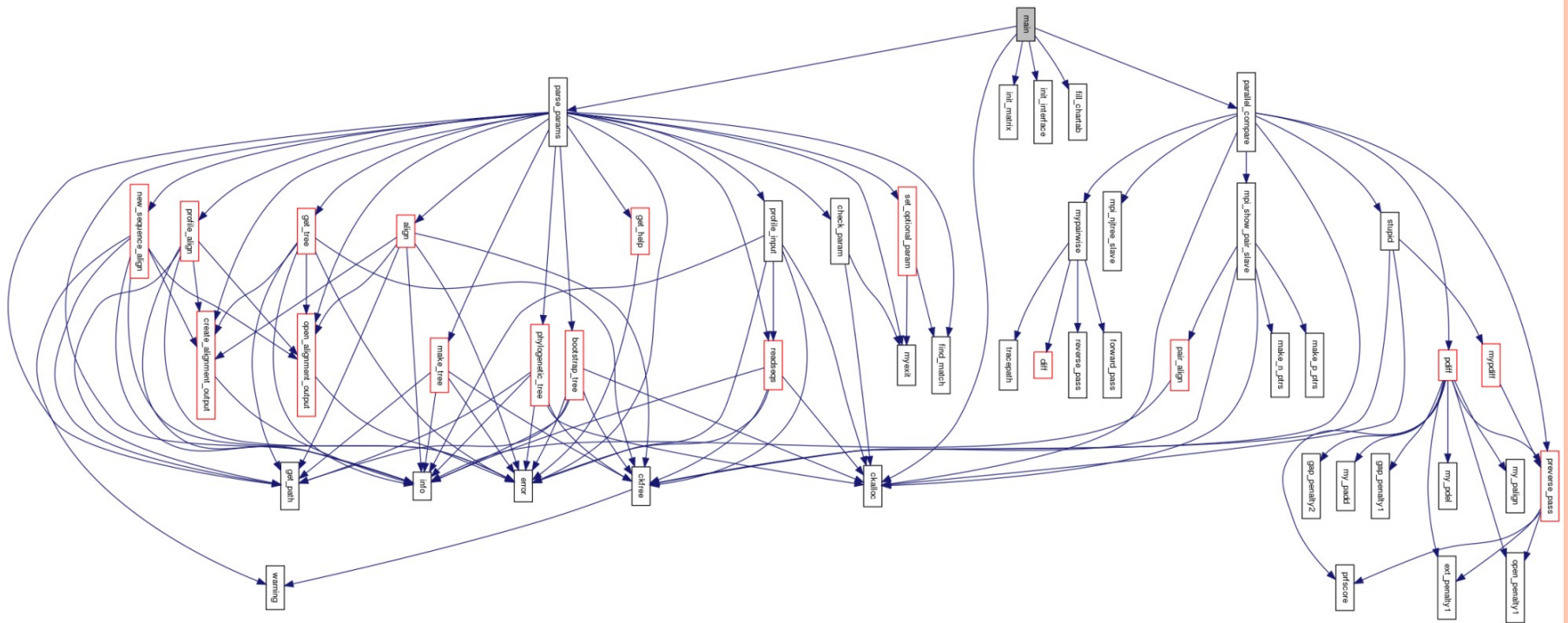
interface.c - Command line interface

random.c – Random number generators

readmat.c – Matrix initialization

sequence.c – Read the sequences and structure data

ГРАФ НА ИЗВИКВАНИЯТА НА ГЛАВНАТА ФУНКЦИЯ MAIN()



РЕДАКТИРАНЕ НА РЕГИОНИ

- След финализирането на множественото подравняване, има възможност за ръчно редактиране на региони от секвенциите посредством софтуерните инструменти на палетата „*Straighten Columns*“, „*Shuffle Right*“ и „*Shuffle Left*“.
- При „изправянето“ на колоните (*Straighten Columns*) на подравняването в ляво за маркирани региони от избрани секвенции се вмъкват празнини.
- При разместването (*Shuffling*) елементите на секвенциите се преместват от единия край на празнина до другия край.

ВИРУСНА БИОИНФОРМАТИКА

- Вирусите са обект на еволюцията и на естествения отбор също като клетъчните форми на живот. Много от вирусите еволюират бързо.
- Когато два вируса инфектират една и съща клетка едновременно, те могат да си обменят генетичен материал и така да се получат нови, "смесени" вируси с уникални свойства. По този начин могат да се появят различни щамове грип.
- РНК вирусите имат високи нива на мутации, което им позволява да еволюират изключително бързо.
- Пример за това е ХИВ, който еволюира във форми, резистентни към лекарства

РЕКОМБИНАЦИЯ

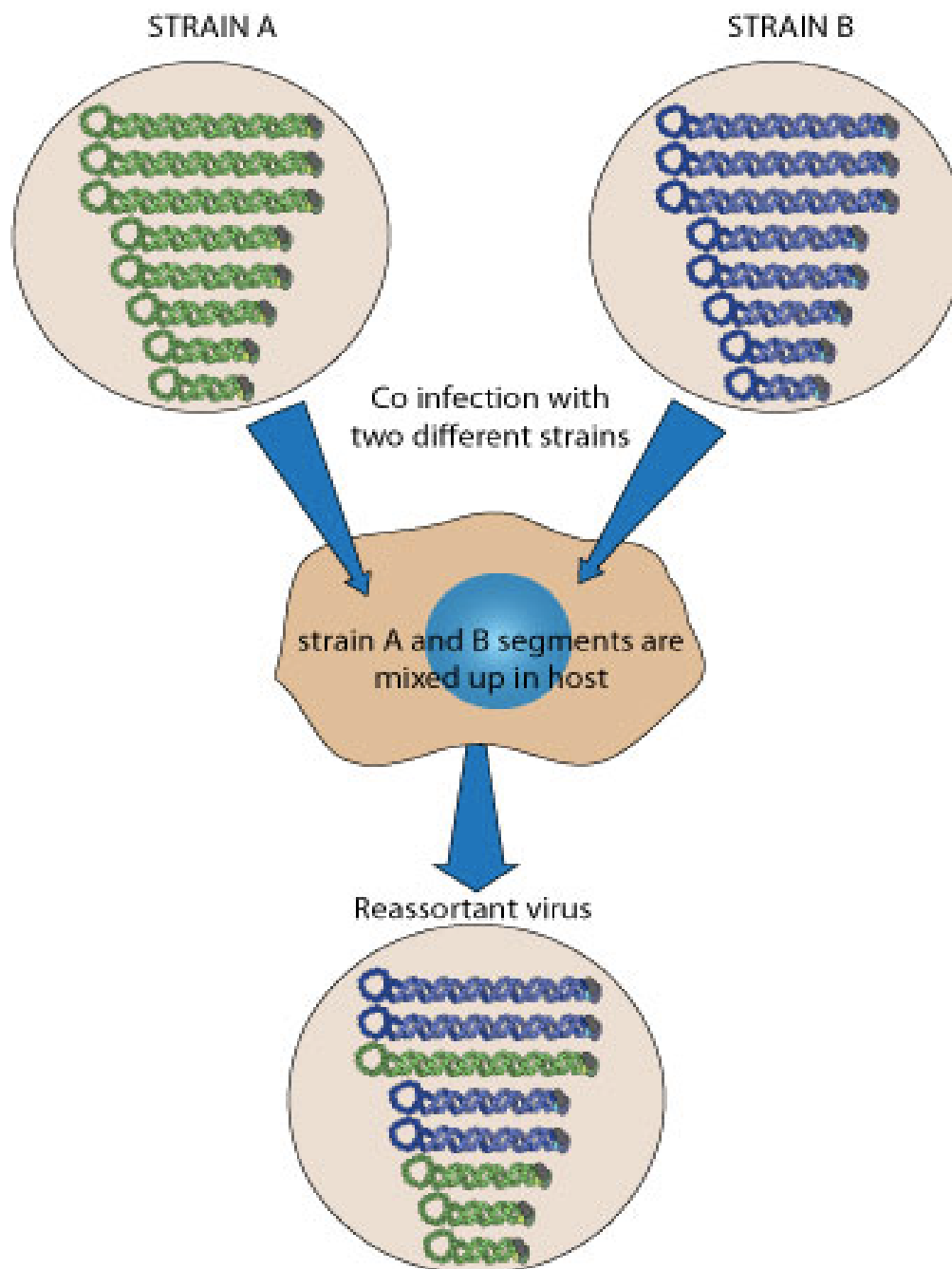
- Вирусите не просто еволюират, те еволюират по-бързо от гостоприемниците си, например човека.
- Това прави еволюцията на вирусите важна тема не само за биолозите, които изучават вируси, но също и за лекарите, сестрите и останалите хора, които работят в здравната система, както и за всеки, който може да бъде изложен на вируси.
- Вирусите си разменят ДНК или РНК чрез процес на име **рекомбинация**.
- Най-често рекомбинацията се случва, когато два вируса заразят едновременно една и съща клетка.
- Тъй като и двата вируса използват клетката, за да произведат нови вирусни частици, във вътрешността ѝ ще има много плаващи вирусни елементи по едно и също време.

РЕКОМБИНАЦИЯ

- Рекомбинацията е процесът, при който две отделни молекули от ДНК или РНК обменят участъци от своя геном.
- Обменът често е в хомоложни области на генома и се извършва както в едноверижни, така и в двуверижни молекули на ДНК и РНК.
- По този начин рекомбинацията е основен механизъм, движещ еволюционната промяна.

РЕКОМБИНАЦИЯ

- При тези обстоятелства рекомбинацията може да протече по два начина.
- Първият е подобни региони от вирусните геноми да се сдвоят и и да си обменят части от тях като физически разрушат и след това свържат отново ДНК или РНК.
- При втория начин вируси с различни сегменти, подобно на малки хромозоми, могат да си разменят някои от тези сегменти чрез процес, наречен ресортиране.



РЕКОМБИНАЦИЯ И ИНФЛУЕНЦА

- Вирусите инфлуенца (вирусите на грипа) са “майстори” на ресортирането.
- Те имат осем сегмента РНК, всеки от които носи по един или по няколко гена.
- Когато два вируса инфлуенца заразят една и съща клетка по едно и също време, някои от новите вирусни частици, произведени в клетката могат да имат микс от сегменти (например сегменти 1-4 от щам А и сегменти 5-8 от щам В).

КАК ЖИВОТИНСКИТЕ ВИРУСИ ЗАРАЗЯВАТ КЛЕТКИ?

- Животинските вируси, също като останалите вируси, зависят от клетка гостоприемник, за да изпълнят жизнения си цикъл.
- За да се възпроизведе, вирусът трябва да зарази клетка гостоприемник и да я програмира да произвежда нови вирусни частици.
- Първата ключова стъпка от инфектирането е разпознаване.
- Животинските вируси имат специални молекули на повърхността си, които им позволяват да се свържат с рецептори от мембраната на клетката гостоприемник.
- Щом се прикрепят за клетката гостоприемник, животинските вируси могат да навлязат в нея по различни начини

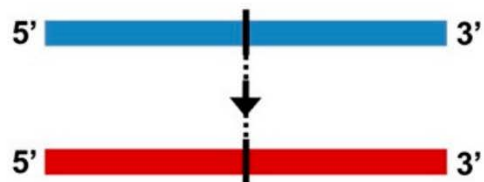
РЕКОМБИНАЦИЯ

- Прасетата са добре известни "смесителни съдове" за вирусите инфлуенца
- Клетките на прасетата могат да бъдат разпознати и заразени от човешки и от птичи вируси на инфлуенца (както и от свински вируси, разбира се).
- Ако дадена клетка в прасето е инфектирана от два вида вирус по едно и също време, тя може да освободи нови вируси, които съдържат смес от генетичните материали на човешкия и на птичия вирус.

РЕКОМБИНАЦИЯ

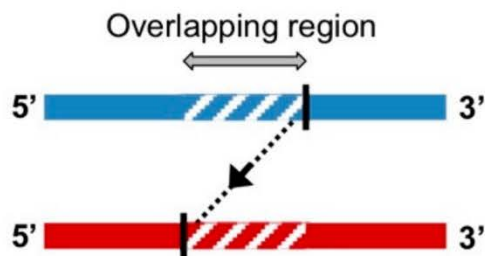
- Този вид обмяна е много често срещана в природата при вирусите инфлуенца.
- Например, щама H1N1 (свински грип), който предизвика пандемия през 2009 г.
- H1N1 съдържаеше сегмент от човешки и птичи вируси, както и свински вируси от Северна Америка и Азия.
- Тази комбинация отразява поредица от ресортирания, които са се случили стъпка по стъпка в продължение на много години, за да се стигне до щама H1N1

Homologous recombination



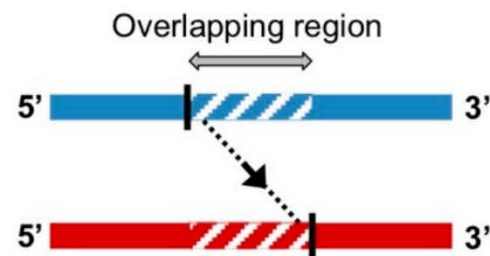
(a)

Nonhomologous recombination



Duplicated sequences

(b)



Deleted sequence

(c)

ЧЕСТОТА НА МУТАЦИИ НА ВИРУСИТЕ

- Някои вируси имат много висока честота на мутации, което им помага да еволюират бързо, осигурявайки им по-голямо генетично разнообразие.
- Два други фактора, които допринасят за бързата еволюция са големият размер на популациите и краткият жизнен цикъл.
- Колкото по-голяма е популацията, толкова по-голям е шансът в нея да има вирус, който носи случайна мутация, която ще бъде селектирана от естествения отбор, например мутация, която прави вируса резистентен към лекарства или усилва вирулентността (инфекциозността) му.
- Освен това вирусите се възпроизвеждат бързо, така че популациите им еволюират по-бързо от тези на гостоприемниците им. Например вирусът ХИВ завършва жизнения си цикъл само за 52 часа, което е много по-кратко време в сравнение с жизнения цикъл на човека.

РЕКОМБИНАЦИЯ НА ВИРУСНИ РНК-И

- Рекомбинация на РНК има по-голяма вероятност да възникне на места (сайтове), с висока степен на хомология между секвенциите на двата родителски вируса, а също и по-лесно се осъществява между тясно свързани вирусни РНК-и.
- Предполагаемите рекомбинационни горещи точки (hot spots) съответстват на хипервариантния регион, в който се появяват чести делеции (изтривания) след преминаване на вируса в тъканна култура или животни
- рекомбинация между суперинфектираща вирусна РНК и РНК

RAT (RECOMBINATION ANALYSIS TOOL)

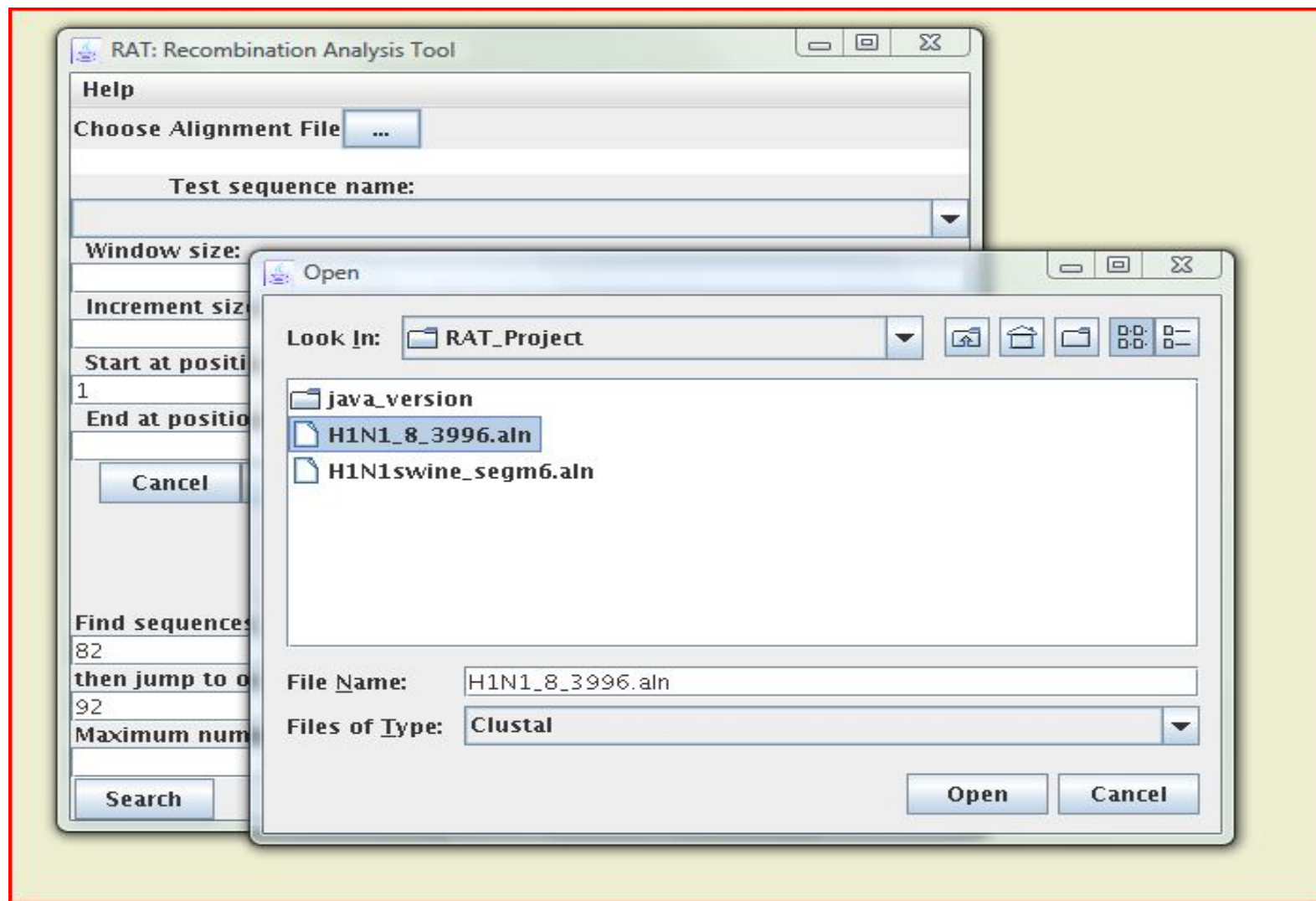
- RAT (Recombination Analysis Tool) е кросплатформено Java-базирано приложение за изследване на рекомбинацията в произволен брой подравнени секвенции (протеини или ДНК) с произволна дължина (къси секвенции на вируси до пълни геноми).
- Позволява анализа на данни от файлове с до 7 формата
- Използва метода на дистанциите за откриване на рекомбинация.
- Всички операции при RAT се изпълняват през графичния потребителски интерфейс GUI.

RAT (RECOMBINATION ANALYSIS TOOL)

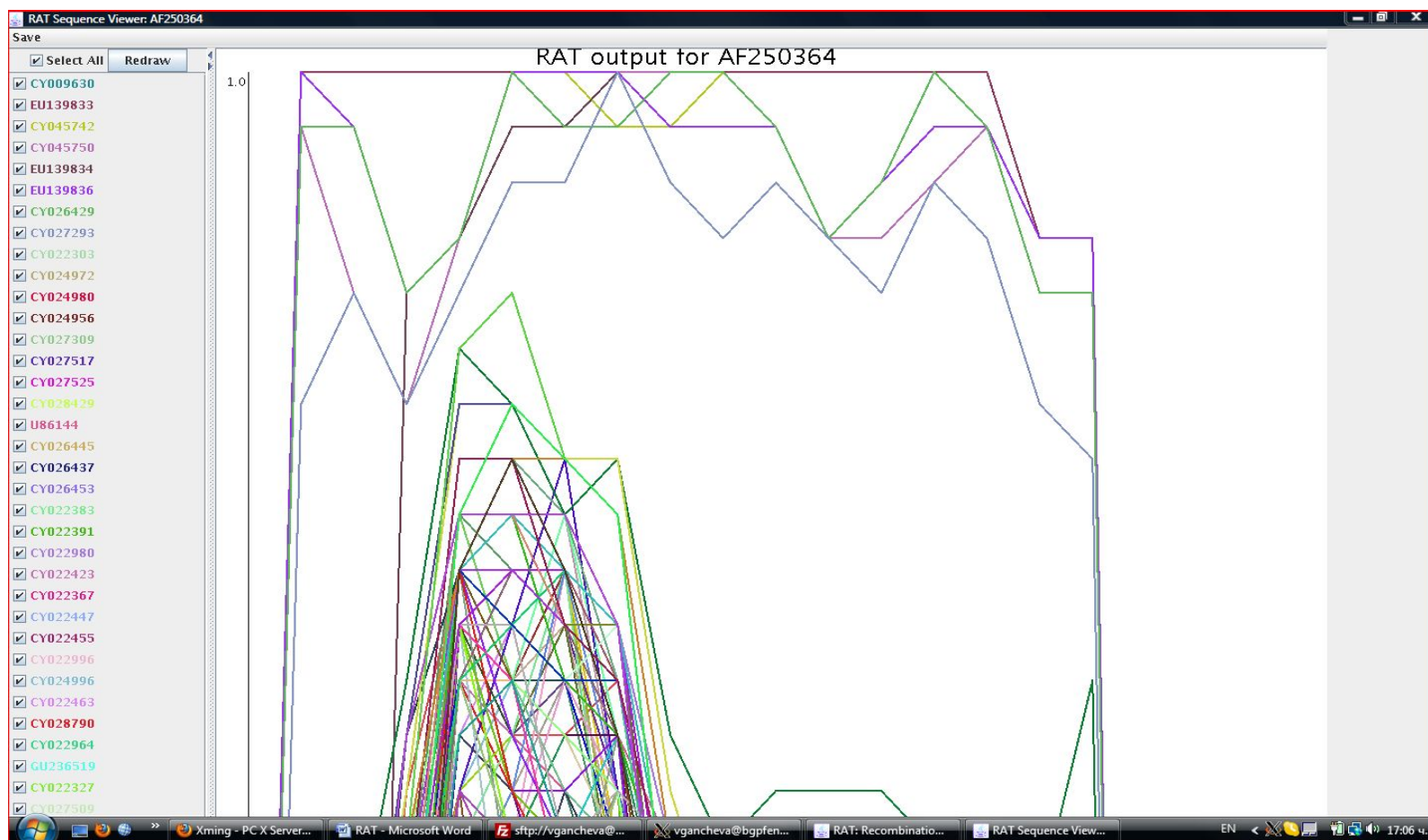
- Изходните данни се съхраняват във файлове с формат .txt, .xls, .csv или като .jpg файлове при графичен изход.
- Стартовата позиция за анализа по подразбиране е първата.
- При желание да се анализира дадена област в рамките на секвенцията, може да се промени стартовата позиция и това ще промени областта за анализ.
- Плъзгащият се прозорец на софтуера RAT се движи по сцените на кратки стъпки и се разглежда пространството с размера на зададения "размер на прозореца".

RAT (RECOMBINATION ANALYSIS TOOL)

- Автоматичното търсене сканира секвенции, които отговарят на параметрите, представени в три полета. Разглеждат се всички секвенции в съответния ред и се сравняват с всички други секвенции.
- *Поле 1:* Търсят се участъци, които започват със сходство, зададени в %.
- *Поле 2:* Търсят се участъци с прилика над зададен %.
- *Поле 3:* Определя се максималният брой секвенции, които да се анализират за рекомбинация. Може да се ограничи броят на секвенциите, които да се включат като полезни в рекомбинацията. По подразбиране това е броят на всички секвенции (без ограничение).



ЕКСПЕРИМЕНТАЛНИТЕ РЕЗУЛТАТИ ЗА ИЗСЛЕДВАНЕ НА ГОРЕЩИ ТОЧКИ НА РЕКОМБИНАЦИЯ НА ГРИПЕН ВИРУС A/H1N1, СЕГМЕНТ 1 (3780 СЕКВЕНЦИИ ОТ РАЗЛИЧНИ ИЗОЛАТИ)



ИЗСЛЕДВАНЕ НА ГОРЕЩИ ТОЧКИ НА РЕКОМБИНАЦИЯ ГРИПЕН ВИРУС ТИП А/Н3N8

The image shows two windows from a software application. The left window is titled 'RAT: Recombination Analysis Tool' and contains various input fields and buttons. The right window is titled 'Auto Search Output' and displays a list of possible recombinant sequences and a detailed breakdown of a specific sequence.

RAT: Recombination Analysis Tool

Help

Choose Alignment File ...

C:\Users\USER\Documents\HorseH3N8\H3N8_horse_nucl1.aln

Test sequence name:

CY028843

Window size:

234

Increment size:

117

Start at position:

1

End at position:

2341

Cancel Execute

Auto Search

Find sequences that start below the following similarity (%):

82

then jump to over (%):

92

Maximum number of contributing sequences

95

Search

Auto Search Output

Save

Possible Recombinant: CY067507

Possible Recombinant: CY067563

Possible Recombinant: DQ222920

Possible Recombinant: EU794548

Possible Recombinant: EU794564

Possible Recombinant: EU794516

Possible Recombinant: EU794492

Possible Recombinant: EU794500

Possible Recombinant: EU794508

Possible Recombinant: EU794540

Possible Recombinant: EU794572

Possible Recombinant: EU794524

Possible Recombinant: GU571147

Possible Recombinant: EU794556

Possible Recombinant: FJ375218

Possible Recombinant: FJ375219

Click here to see

GU051895:

Window 1, at nucleotide position 234.0

Possible Recombinant: CY032952

Window 2, at nucleotide position 351.0

Window 3, at nucleotide position 468.0

Window 4, at nucleotide position 585.0

Window 5, at nucleotide position 702.0

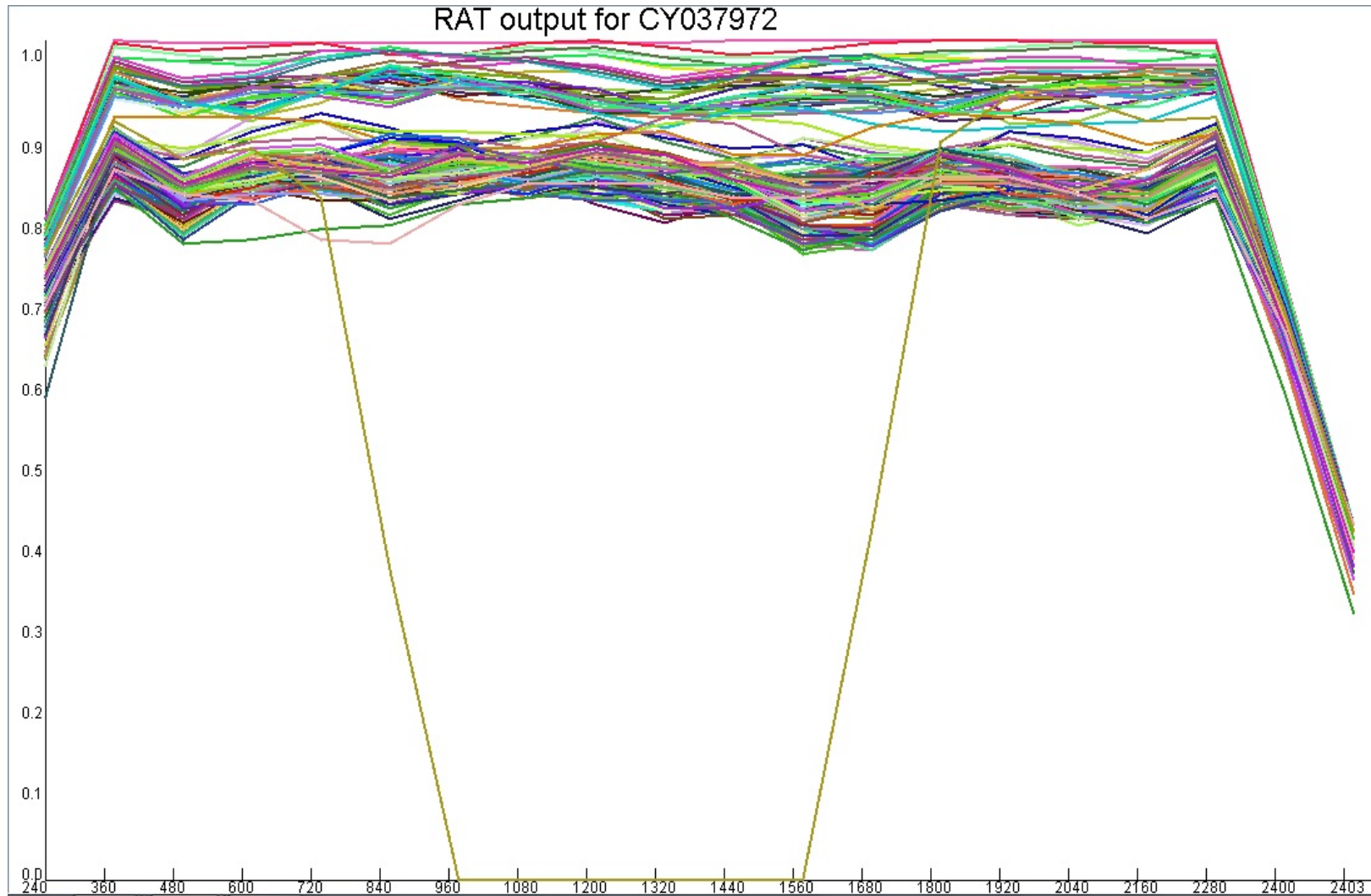
Window 6, at nucleotide position 819.0

Window 7, at nucleotide position 936.0

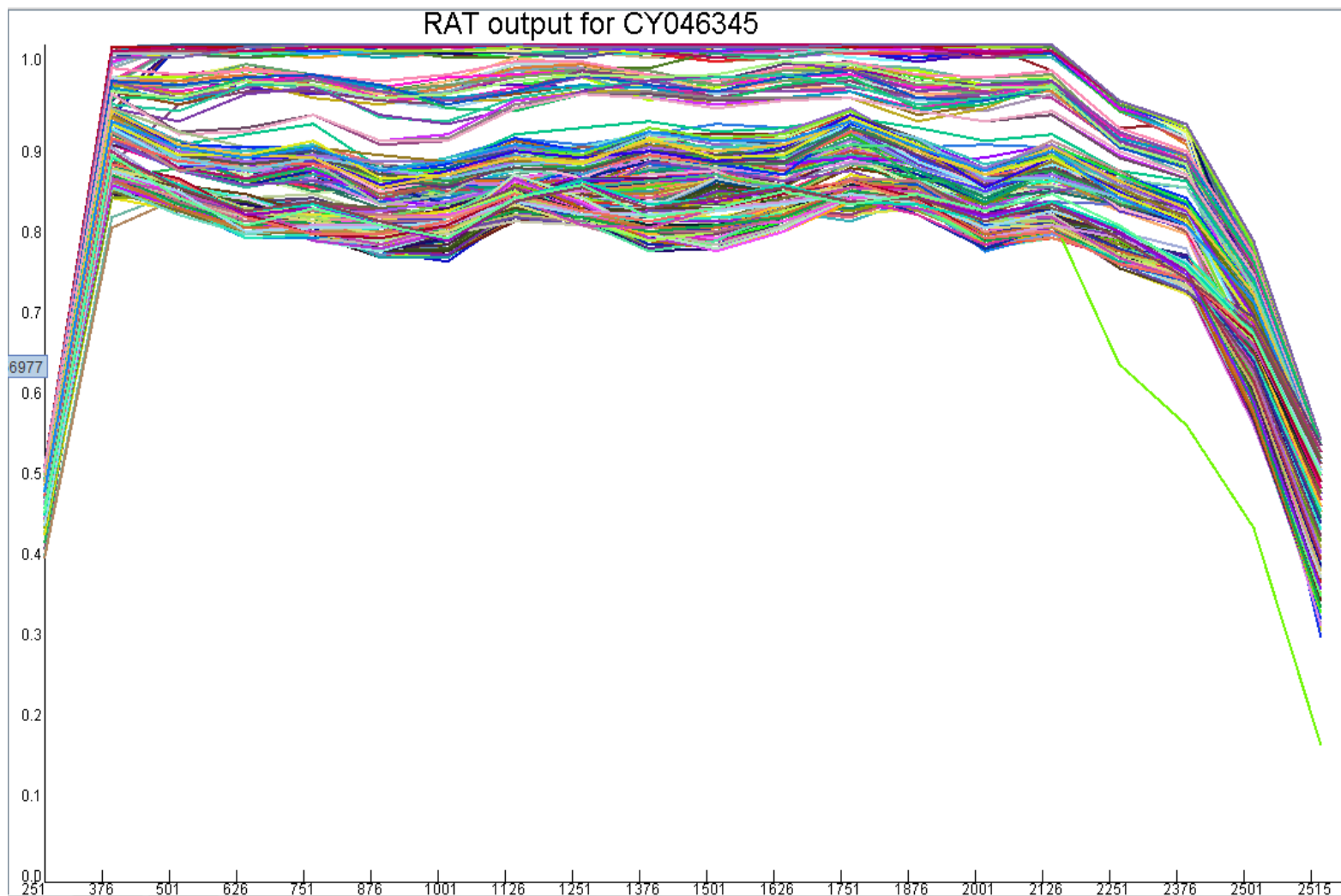
Window 8, at nucleotide position 1053.0

Window 9, at nucleotide position 1170.0

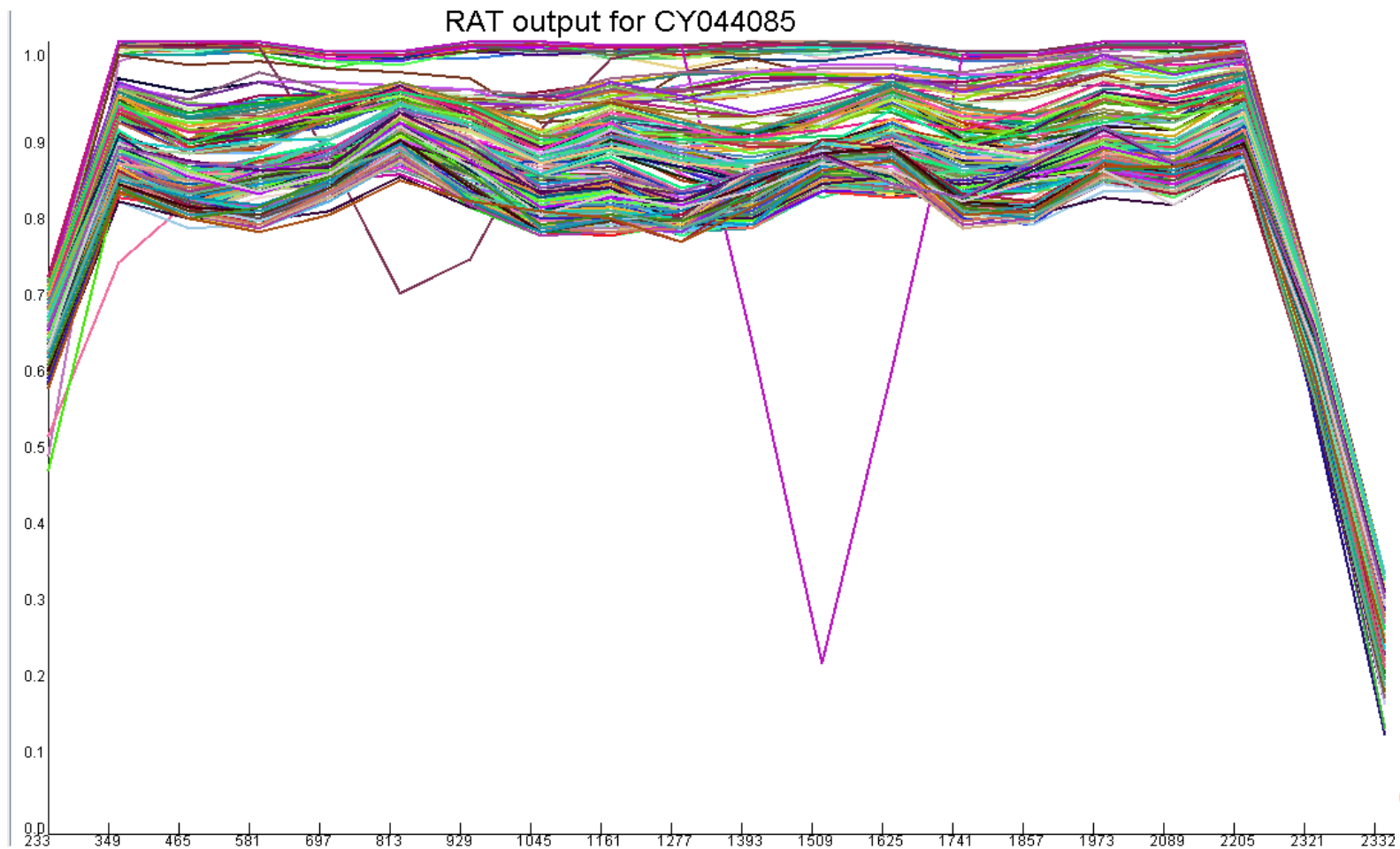
ТОЧКИ НА РЕКОМБИНАЦИЯ ПРИ ГРИПЕН ВИРУС A/H1N1, СЕГМЕНТ 1



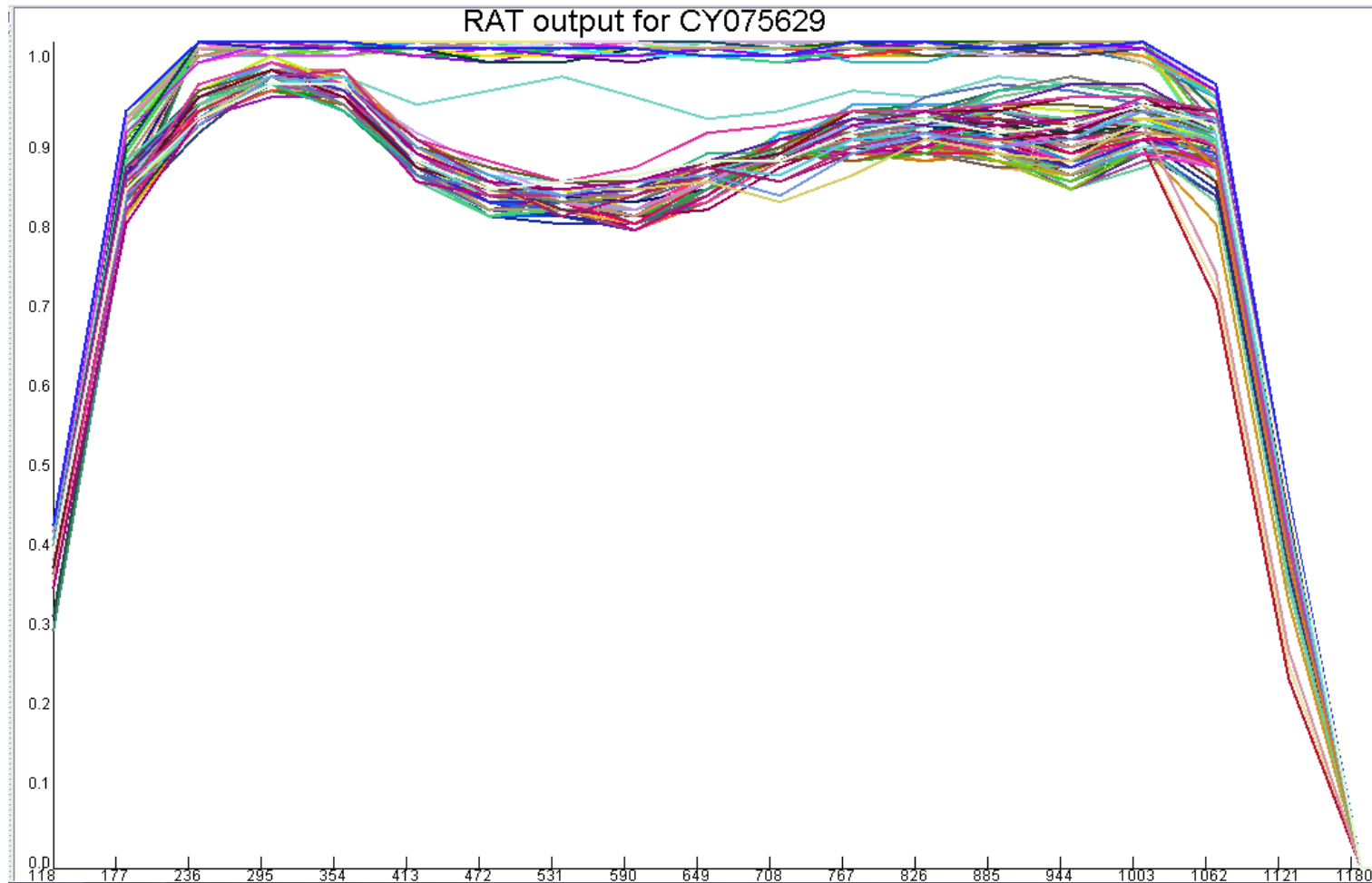
ТОЧКИ НА РЕКОМБИНАЦИЯ ПРИ ГРИПЕН ВИРУС А/Н1N1, СЕГМЕНТ 2



ТОЧКИ НА РЕКОМБИНАЦИЯ ПРИ ГРИПЕН ВИРУС А/Н1N1, СЕГМЕНТ 3



ТОЧКИ НА РЕКОМБИНАЦИЯ ПРИ ГРИПЕН ВИРУС A/H1N1, СЕГМЕНТ 7



ТОЧКИ НА РЕКОМБИНАЦИЯ ПРИ ГРИПЕН ВИРУС А/Н1N1, СЕГМЕНТ 8

