



**Выравнивание  
последовательностей.  
Выбор оптимальной модели  
эволюции нуклеотидов**

# Интерфейс MEGA



- ▶ ALIGN — выравнивание последовательностей
- ▶ DATA — загрузка и редактирование данных
- ▶ MODELS — выбор оптимальной модели эволюции нуклеотидов
- ▶ DISTANCE — расчёт матрицы расстояний по выравниванию
- ▶ PHYLOGENY — построение филогенетического дерева

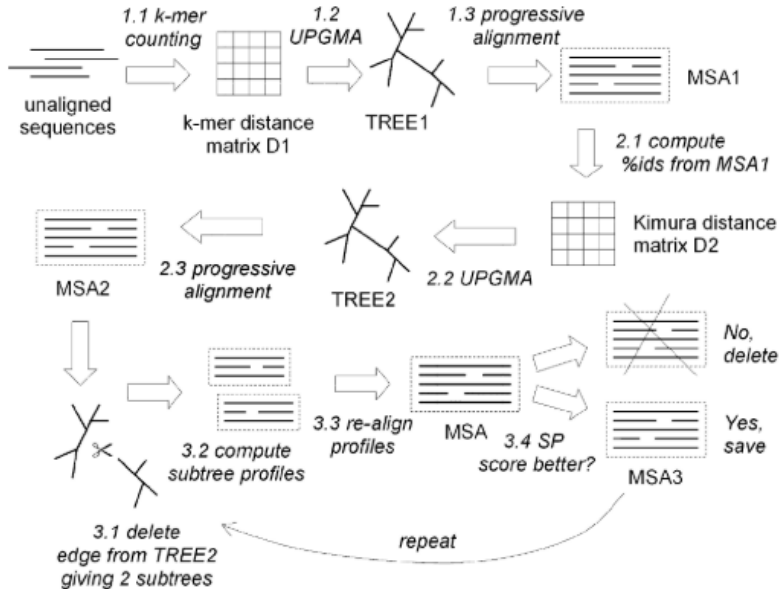
# Выравнивание последовательностей

Откройте файл `CoV_task.fa` для выравнивания.

1. DATA
2. Open File/Session
3. Выберите файл
4. Analyze or Align File? — Align

Species/Abbrv	* * *
1. Bat Alphacoronavirus DQ648792.1	A T G A A A T A T A C A C T T T T A T T T T G T G T A G T G T T T G C T A C G G T G T C T

# Алгоритм MUSCLE



# MUSCLE



- ▶ **Align DNA**

- ▶ Align Codons

Select all sequences

# MUSCLE



- ▶ **Align DNA**
- ▶ **Align Codons**

Select all sequences

Option	Setting	
GAP PENALTIES		
Gap Open	<input checked="" type="checkbox"/>	-400.00
Gap Extend	<input checked="" type="checkbox"/>	0.00
MEMORY/ITERATIONS		
Max Memory in MB	<input checked="" type="checkbox"/>	2048
Max Iterations	<input checked="" type="checkbox"/>	16
ADVANCED OPTIONS		
Cluster Method (Iterations 1,2)	<input checked="" type="checkbox"/>	UPGMA
Cluster Method (Other Iterations)	<input checked="" type="checkbox"/>	UPGMA
Min Diag Length (Lambda)	<input checked="" type="checkbox"/>	24

**Задание 1.** Выполните выравнивание последовательностей с помощью алгоритма MUSCLE. Какова длина полученного выравнивания (порядковый номер последней колонки нуклеотидов, Site)?

**Задание 1.** Выполните выравнивание последовательностей с помощью алгоритма MUSCLE. Какова длина полученного выравнивания (порядковый номер последней колонки нуклеотидов, Site)?

5256

Сохраните полученное выравнивание в формате fasta:  
Data/Export Alignment/FASTA format



# Расчёт матрицы расстояний

- ▶ Выбор файла с выравниванием
- ▶ Настройка импорта данных



Compute  
Pairwise  
Distance

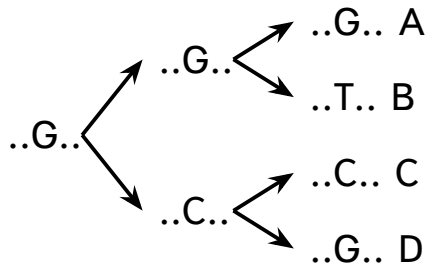
- ▶ Уточнение в выравнивании белок кодирующая последовательность или нет (выбор таблицы генетического кода)

# Расчёт матрицы расстояний

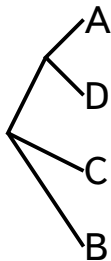
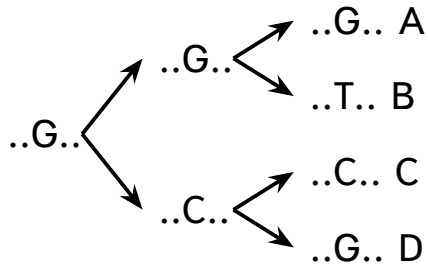


ANALYSIS	
Scope	→ Pairs of taxa
ESTIMATE VARIANCE	
Variance Estimation Method	→ None
No. of Bootstrap Replications	→ Not Applicable
SUBSTITUTION MODEL	
Substitutions Type	→ Nucleotide
Genetic Code Table	→ Not Applicable
Model/Method	→ No. of differences
Fixed Transition/Transversion Ratio	→ Not Applicable
Substitutions to Include	→ d: Transitions + Transversions
RATES AND PATTERNS	
Rates among Sites	→ Uniform Rates
Gamma Parameter	→ Not Applicable
Pattern among Lineages	→ Same (Homogeneous)
DATA SUBSET TO USE	
Gaps/Missing Data Treatment	→ Pairwise deletion
Site Coverage Cutoff (%)	→ Not Applicable
Select Codon Positions	→ <input checked="" type="checkbox"/> 1st <input checked="" type="checkbox"/> 2nd <input checked="" type="checkbox"/> 3rd <input checked="" type="checkbox"/> Noncoding Sites

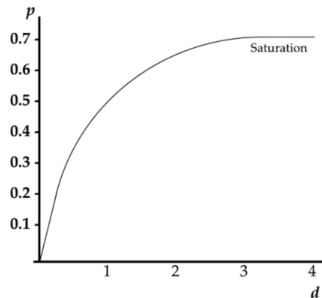
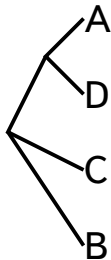
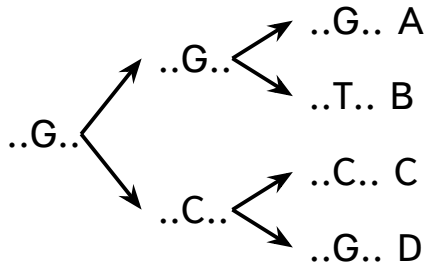
# Несоответствие наблюдаемого и действительного расстояния



# Несоответствие наблюдаемого и действительного расстояния



# Несоответствие наблюдаемого и действительного расстояния



Lemey, 2009

# Способы расчёта расстояний

- ▶ No. of differences ( $n_d$ )
- ▶ p.distance  $p = n_d/L$
- ▶ Jukes-Cantor distance (JC69):  $d = -\frac{3}{4}\log_e(1 - \frac{4}{3}p)$
- ▶ Kimura 2-parameter distance (K80):  
 $d = -0.5\log_e(w_1) - 0.25\log_e(w_2)$   
 $w_1 = 1 - 2P - Q$   
 $w_2 = 1 - 2Q$

**Задание 2.** Рассчитайте матрицу попарных расстояний для получившегося выравнивания. Какое максимальное расстояние между последовательностями (No. of differences)?

# Способы расчёта расстояний

- ▶ No. of differences ( $n_d$ )
- ▶ p.distance  $p = n_d/L$
- ▶ Jukes-Cantor distance (JC69):  $d = -\frac{3}{4}\log_e(1 - \frac{4}{3}p)$
- ▶ Kimura 2-parameter distance (K80):  
 $d = -0.5\log_e(w_1) - 0.25\log_e(w_2)$   
 $w_1 = 1 - 2P - Q$   
 $w_2 = 1 - 2Q$

**Задание 2.** Рассчитайте матрицу попарных расстояний для получившегося выравнивания. Какое максимальное расстояние между последовательностями (No. of differences)?

# Обрезать или не обрезать выравнивание?

Gblocks

**Gblocks Server**

Paste an alignment in NBRF/PIR or FASTA format:

Or upload an alignment file:

all\_seq3\_align.fas

Type of sequence:

DNA ☒ || Protein ☐ || Codons ☐

Options for a less stringent selection:

- ☒ Allow smaller final blocks
- ☒ Allow gap positions within the final blocks
- ☒ Allow less strict flanking positions

Options for a more stringent selection:

- ☐ Do not allow many contiguous nonconserved positions



**Задание 3.** Обрежьте выравнивание с наиболее щадящими настройками.

Перед загрузкой на сервер **Gblocks** сохраните выравнивание в формате fasta: File/Export alignment/FASTA format

```
Processed file: input.fasta  
Number of sequences: 34  
Alignment assumed to be: DNA  
New number of positions: 2462 (selected positions are underlined in blue)
```

Сохраните обрезанное выравнивание (Cured alignment in FASTA Format).

**Задание 3.** Обрежьте выравнивание с наиболее щадящими настройками.

Перед загрузкой на сервер **Gblocks** сохраните выравнивание в формате fasta: File/Export alignment/FASTA format

```
Processed file: input.fasta  
Number of sequences: 34  
Alignment assumed to be: DNA  
New number of positions: 2462 (selected positions are underlined in blue)
```

Сохраните обрезанное выравнивание (Cured alignment in FASTA Format).

**Задание 4.** Повторите процедуру выравнивания для исходного файла, изменив значение Gap Open penalty на -200. Сравните длины выравниваний до и после обрезки Gblocks

**Задание 4.** Повторите процедуру выравнивания для исходного файла, изменив значение Gap Open penalty на -200. Сравните длины выравниваний до и после обрезки Gblocks

6291, 2120

# Форматы выравниваний

## FASTA

```
>A
ATGAAATATACACTTTTATTTT--
>B
ATGTTGGTGATATTGTTAATGTTA
>C
ATGTTTTTGATACTTTTAATTTCC
```

## CLUSTAL

```
CLUSTAL O(1.2.4) multiple sequence
alignment
```

```
A   ATGAAATATACACTTTTATTTT--
B   ATGTTGGTGATATTGTTAATGTTA
C   ATGTTTTTGATACTTTTAATTTCC
    ***          * * * * * *
```

## PHYLIP

```
PHYLIP
3  24
A   ATGAAATATACACTTTTATTTT--
B   ATGTTGGTGATATTGTTAATGTTA
C   ATGTTTTTGATACTTTTAATTTCC
```

## NEXUS

```
begin data;
dimensions ntax=3 nchar=24;
format interleave datatype=DNA missing=N gap=-
;
matrix
A   ATGAAATATACACTTTTATTTT--
B   ATGTTGGTGATATTGTTAATGTTA
C   ATGTTTTTGATACTTTTAATTTCC
;
end;
```

# Модели эволюции нуклеотидов

Модель	Частоты нуклеотидов	Частоты переходов	Свободные параметры
JC69	равные	равные	0
K80	равные	$T_s \neq T_v$	1
T92	$AT \neq GC$	$T_s \neq T_v$	2
F81	неравные	равные	3
F84	неравные	$T_s \neq T_v$	4
GTR	неравные	неравны для всех переходов	8

**Транзиции** ( $T_s$ ) — A-G или C-T. **Трансверзии** ( $T_v$ ) — A-T, A-C, G-C, G-T.

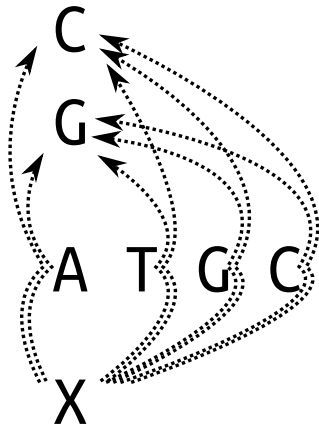
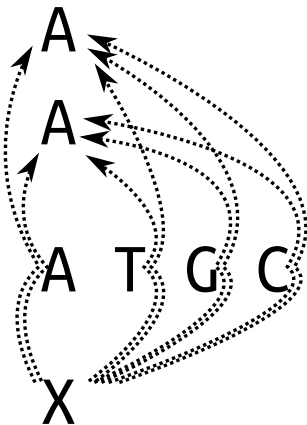
+G — различные частоты замен для сайтов

+I — доля инвариантных сайтов

# Марковские цепи и выбор модели эволюции

Частоты  
переходов

Частоты  
нуклеотидов



# Выбор оптимальной модели эволюции



Find best DNA/protein model

Option	Setting
<b>ANALYSIS</b>	
Tree to Use →	<i>Automatic (Neighbor-joining tree)</i>
User Tree File →	Not Applicable
Statistical Method →	<i>Maximum Likelihood</i>
<b>SUBSTITUTION MODEL</b>	
Substitutions Type →	Nucleotide ▼
Genetic Code Table →	Not Applicable
<b>DATA SUBSET TO USE</b>	
Gaps/Missing Data Treatment →	<i>Use all sites</i>
Site Coverage Cutoff (%) →	Not Applicable
Select Codon Positions →	<input checked="" type="checkbox"/> 1st <input checked="" type="checkbox"/> 2nd <input checked="" type="checkbox"/> 3rd <input checked="" type="checkbox"/> Noncoding Sites
Branch Swap Filter →	<i>None</i>
<b>SYSTEM RESOURCE USAGE</b>	
Number of Threads →	3

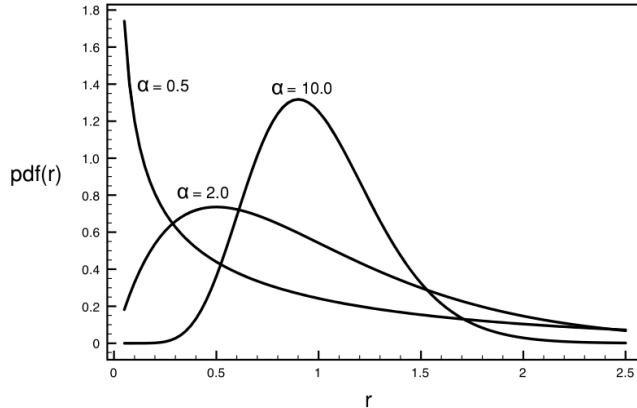


# Просмотр оптимальных параметров

Запустите подбор модели, не меняя настройки по умолчанию. Используйте обрезанное выравнивание ДНК, которое было получено без модификации настроек, его длина 2462.

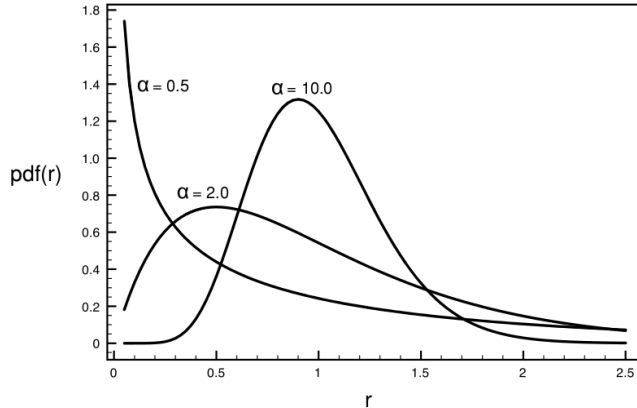
- ▶ Model — название модели
- ▶ Parameters — общее число параметров
- ▶ BIC, AIC — информационные критерии
- ▶ lnL — правдоподобие
- ▶ (+G)/(+I) — различия в частотах замен между сайтами и доля инвариантных сайтов
- ▶ R —  $T_s/T_v$
- ▶ f(N) — частоты нуклеотидов
- ▶ f(NN) — частоты переходов

# Гамма (G) распределение



**Задание 5.** Какая модель наилучшая?

# Гамма (G) распределение



**Задание 5.** Какая модель наилучшая?

GTR+G+I

**Задание 6.** Откройте последовательности из файла DBY\_intron 7\_part.fasta. Постройте выравнивание с помощью алгоритма muscle (учтите, что даны последовательности интронов). Обрежьте его в Gblocks при щадящих настройках. Подберите модель эволюции нуклеотидов.

**Задание 6.** Откройте последовательности из файла DBY\_intron 7\_part.fasta. Постройте выравнивание с помощью алгоритма muscle (учтите, что даны последовательности интронов). Обрежьте его в Gblocks при щадящих настройках. Подберите модель эволюции нуклеотидов.

T92+G