

Markers and phylogenetics  
oooooooooooooooooooo

Phylogenomics intro  
oooooooooooo

Single copy orthologs  
oooooooooooo

Reduced representation  
oooooo

Alignment-free methods  
oooo

Outro  
o

# Молекулярная филогенетика и филогеномика

Полина Дроздова

9–10 октября 2023



[https://github.com/drozdovapb/ZIN\\_school\\_2023](https://github.com/drozdovapb/ZIN_school_2023)

## Терминология

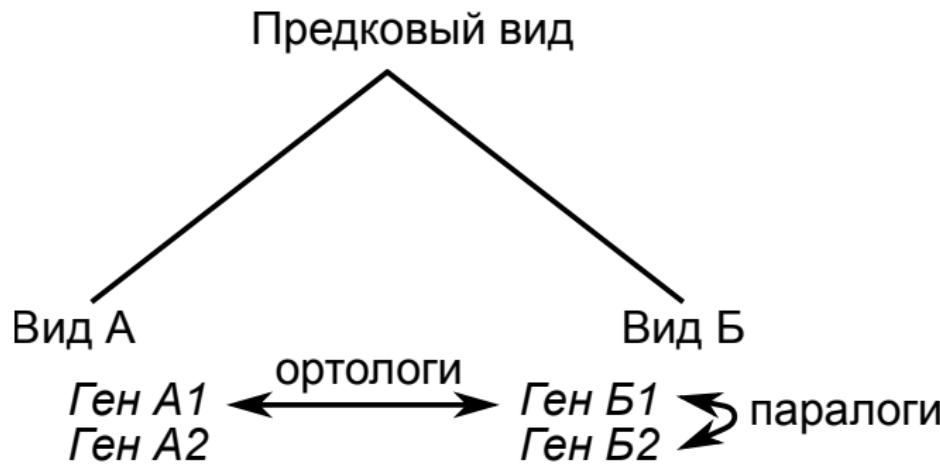
**Гомология (homology)** — сходство признаков, обусловленное общим происхождением.

**Дивергенция (divergence)** — расхождение признаков в ходе эволюции.

**Гомоплазия (homoplasy)** — сходство признаков, не обусловленное общим происхождением. Часто обусловлена **параллельной (конвергентной)** эволюцией.

**Гомологи** = ортологи (orthologs/orthologues) + паралоги (paralogs/paralogues).

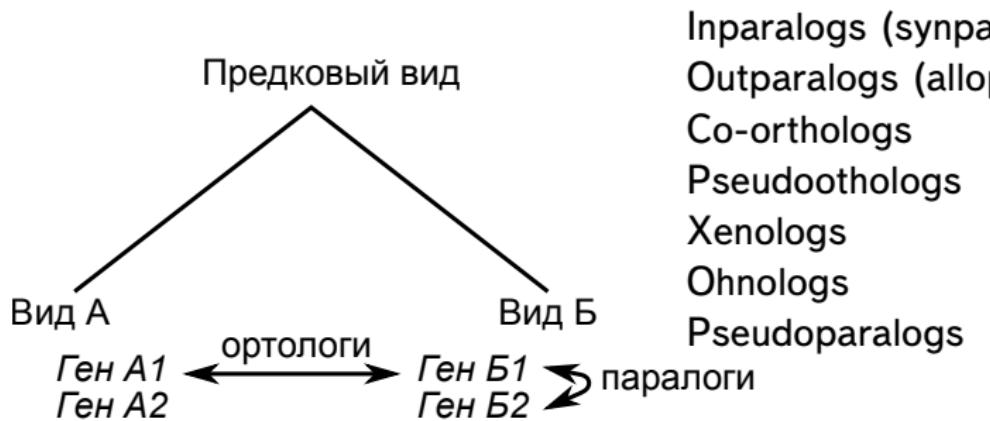
## Варианты гомологии



Термины paralogous / orthologous:

Fitch WM. Distinguishing homologous from analogous proteins.  
Systematic Biology. 1970 Jun;19(2):99-113.

# Варианты гомологии



см. Koonin EV. Orthologs, paralogs, and evolutionary genomics  
1. Annu. Rev. Genet. 2005;39:309-38.

Markers and phylogenetics  
oooooooooooooooo

Phylogenomics intro  
oooooooooooo

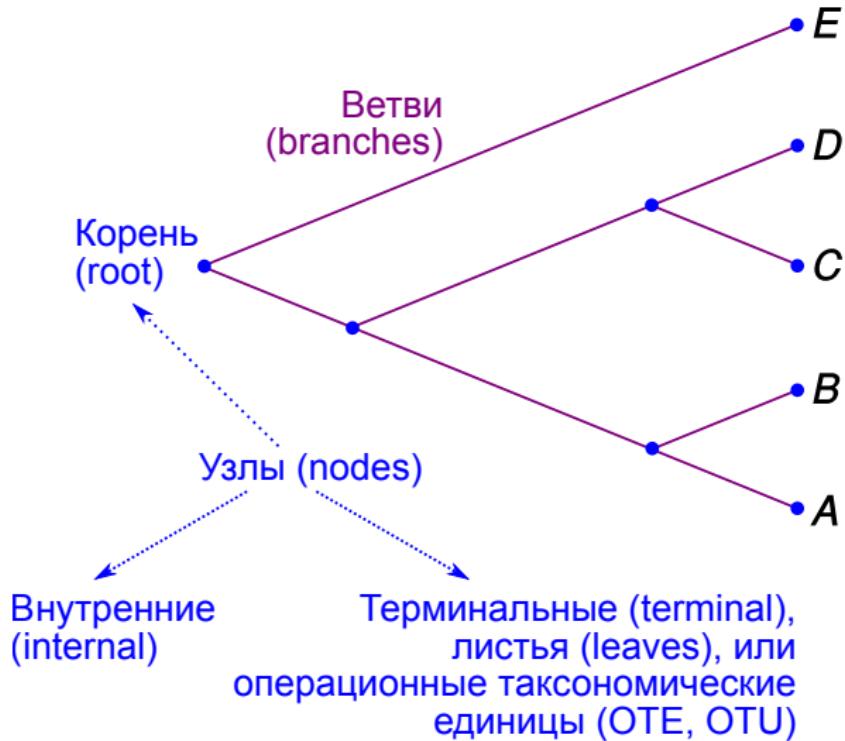
Single copy orthologs  
oooooooooooo

Reduced representation  
oooooo

Alignment-free methods  
oooo

Outro  
o

## Названия частей дерева



Markers and phylogenetics  
oooooooooooooooo

Phylogenomics intro  
oooooooooooo

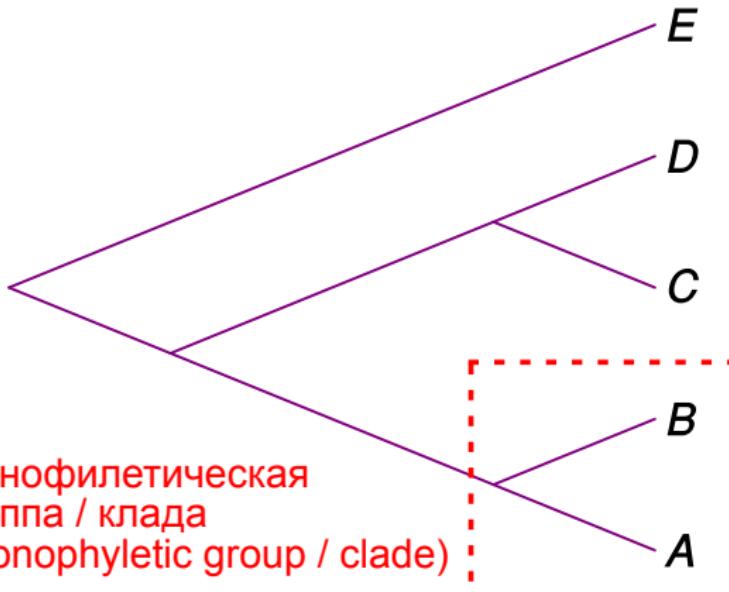
Single copy orthologs  
oooooooooooo

Reduced representation  
oooo

Alignment-free methods  
oooo

Outro  
o

# Монофилия, полифилия и парафилия



Markers and phylogenetics  
oooooooooooooooooooo

Phylogenomics intro  
oooooooooooo

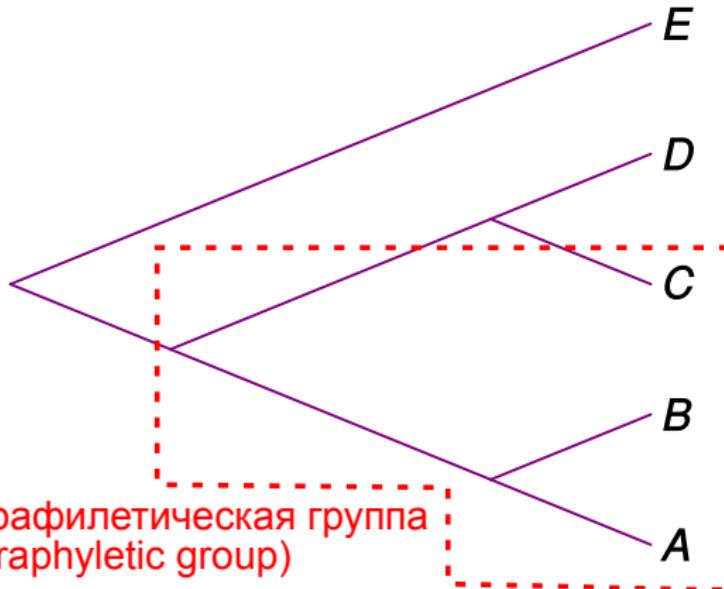
Single copy orthologs  
oooooooooooo

Reduced representation  
oooo

Alignment-free methods  
oooo

Outro  
o

# Монофилия, полифилия и парафилия



Парафилетическая группа  
(paraphyletic group)

Markers and phylogenetics  
oooooooooooooooooooo

Phylogenomics intro  
oooooooooooo

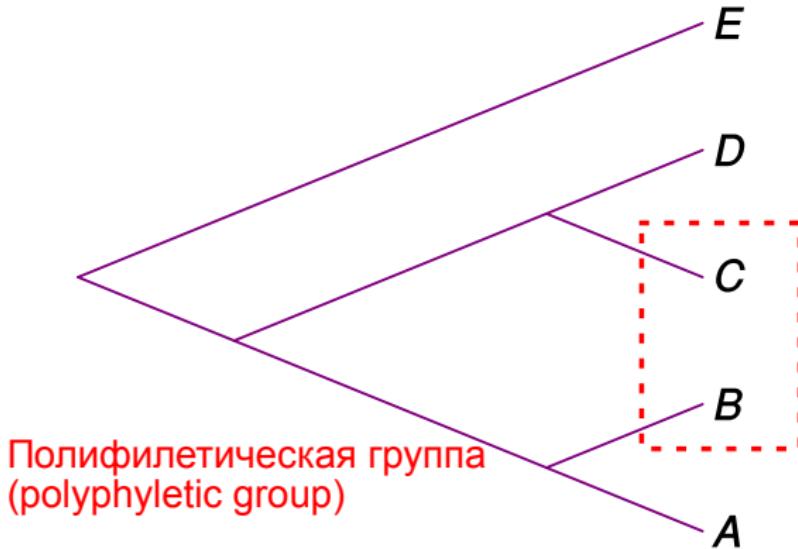
Single copy orthologs  
oooooooooooo

Reduced representation  
oooo

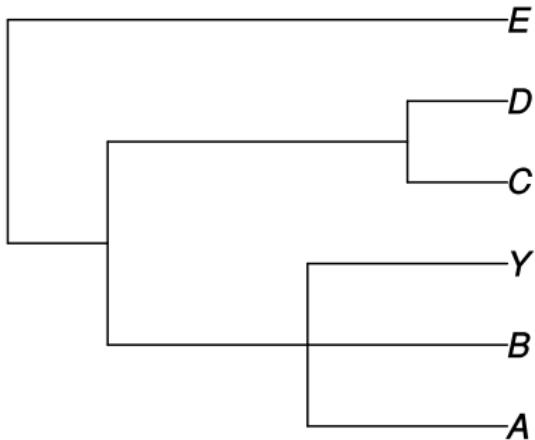
Alignment-free methods  
oooo

Outro  
o

# Монофилия, полифилия и парафилия



# Политомия

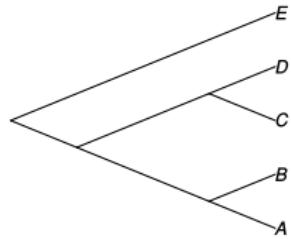


**Soft polytomy** (из-за недостатка данных)  
*vs.*

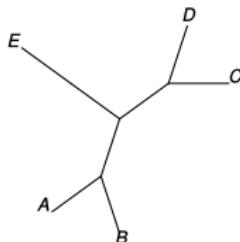
**Hard polytomy** (очень быстрое видообразование)

# Типы дендрограмм

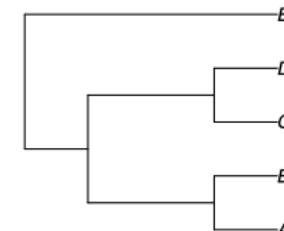
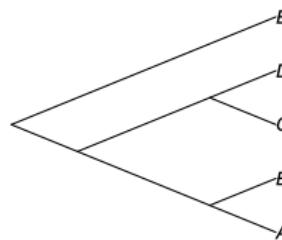
## Укоренённое (rooted) дерево



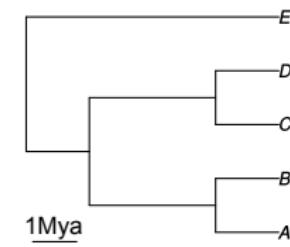
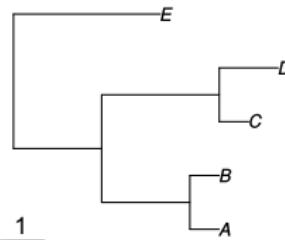
## Неукоренённое (unrooted)



## Кладограммы (длина ветвей не означает ничего)



## Филограммы (длина ветвей что-то означает)



**Аддитивное (additive):**  
длина ветви — число замен

**Ультраметрическое (ultrametric):**  
равная длина ветвей  
(напр., длина ветви — время)

# Число возможных топологий для дерева

Число ОТЕ	Возможных неукоренённых деревьев	Возможных укоренённых деревьев
3	1	3
4	3	15
...	...	...
22	0,5 моля	
...	...	...
n	$(2n - 5)!!$	$(2n - 3)!!$
n	$\frac{(2n-5)!}{2^{n-3} \cdot (n-3)!}$	$\frac{(2n-3)!}{2^{n-2} \cdot (n-2)!}$

Markers and phylogenetics  
oooooooooooooooooooo

Phylogenomics intro  
oooooooooooo

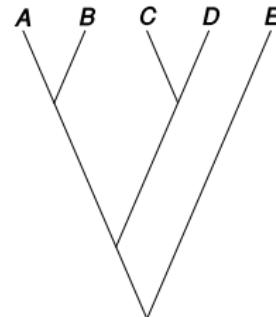
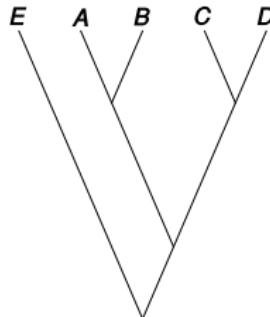
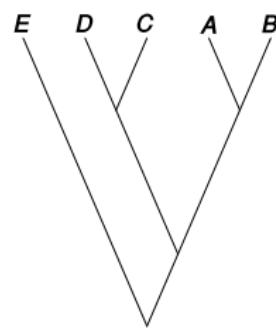
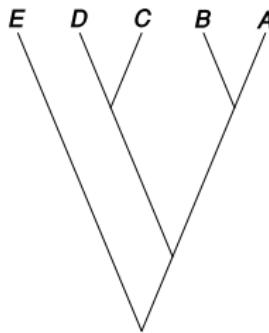
Single copy orthologs  
oooooooooooo

Reduced representation  
oooooo

Alignment-free methods  
oooo

Outro  
o

## Разные варианты отображения топологии



Это одно и то же дерево!

# Почему молекулярные данные?

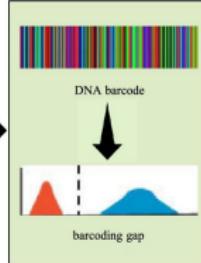
- ✓ Биополимеры есть у всех живых организмов и состоят из одинаковых элементов.
- ✓ Общая методология для разных организмов => проще обучить специалистов;
- ✓ Иногда последовательность ДНК — это единственное, что мы знаем об организме.
- ✓ Нужно мало материала / подходит повреждённый материал.
- ✓ Можно комбинировать с морфологическими данными.

# Баркодинг

(a)

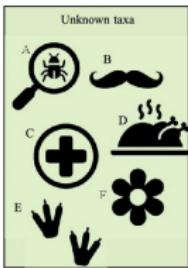


PCR and Sanger sequencing



Reference Libraries

(b)



PCR and Sanger sequencing

Sequence D1 Sequence E1
Sequence A Sequence B1 Sequence F
Sequence C1 Sequence D2 Sequence E2
Sequence B2 Sequence C3 Sequence D3 Sequence E3

OTUs



<https://link.springer.com/article/10.1007/s12686-022-01291-2/>

## Необходимые свойства маркера:

1. **универсальность** => наличие ортологов, желательно без проблем с паралогами;



(Например, компоненты рибосомы)

2. подходящий уровень **разнообразия** (например, межвидовые различия больше внутривидовых), но при этом редкие множественные замены;



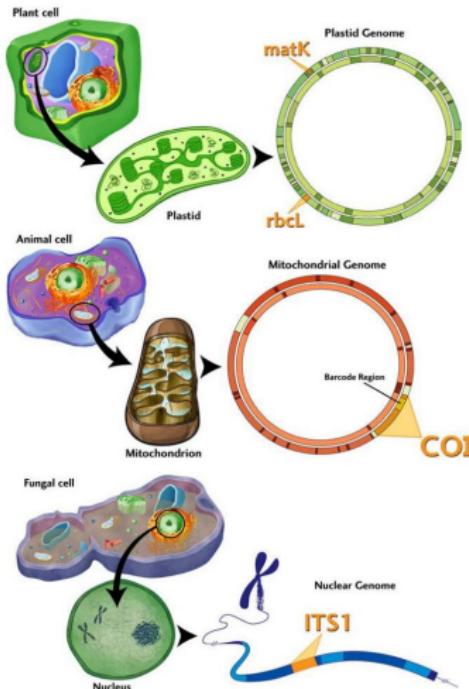
EPIC = exon primed intron crossing

3. желательна селективная нейтральность.

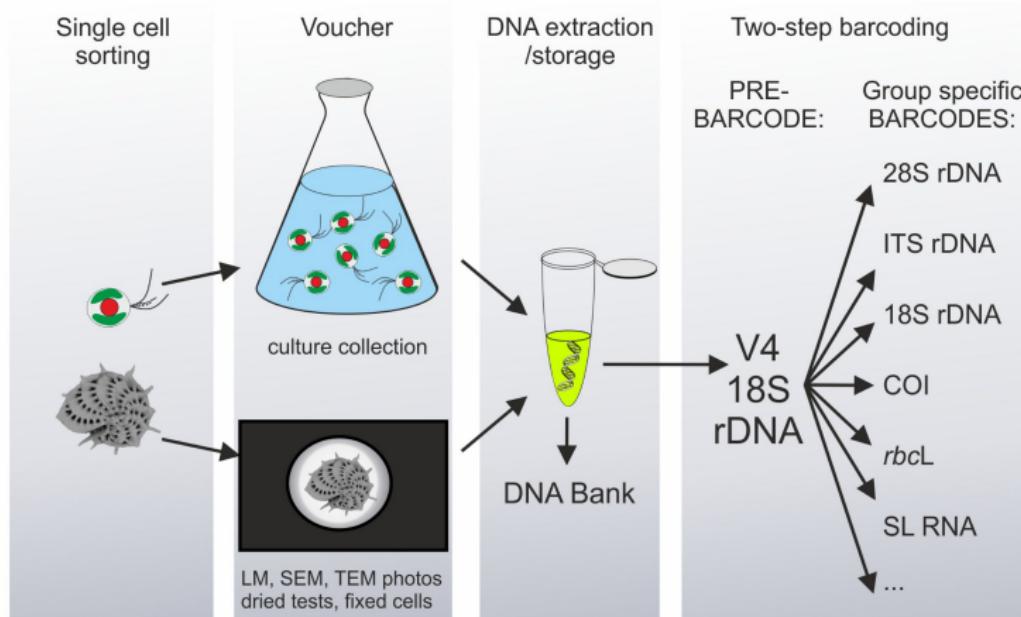
# Примеры часто используемых маркеров

Table: key of tree of life

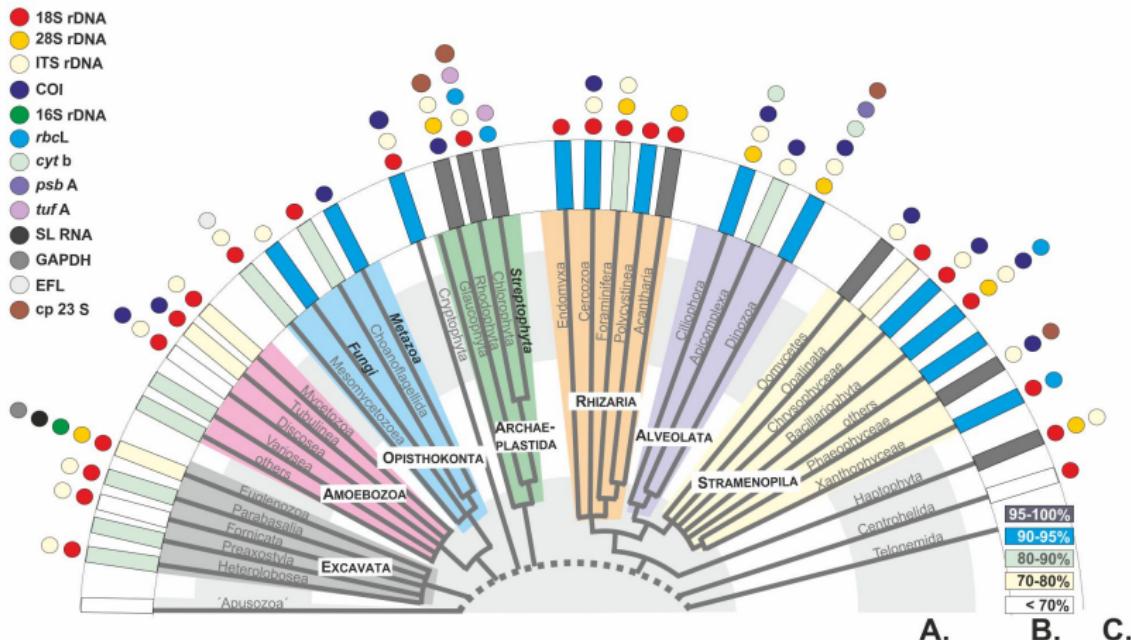
Color	Clade	Primary barcode	Secondary barcode
Red	Animals	COI	COI, 16S
Blue	Fungi	ITS	LSU D1/D2
Green	Green algae	tufA	LSU D2/D3
Purple	Land plants	rbcL/matK	psbA-trnH/ITS
Yellow	Algae	COI-5P	LSU D2/D3



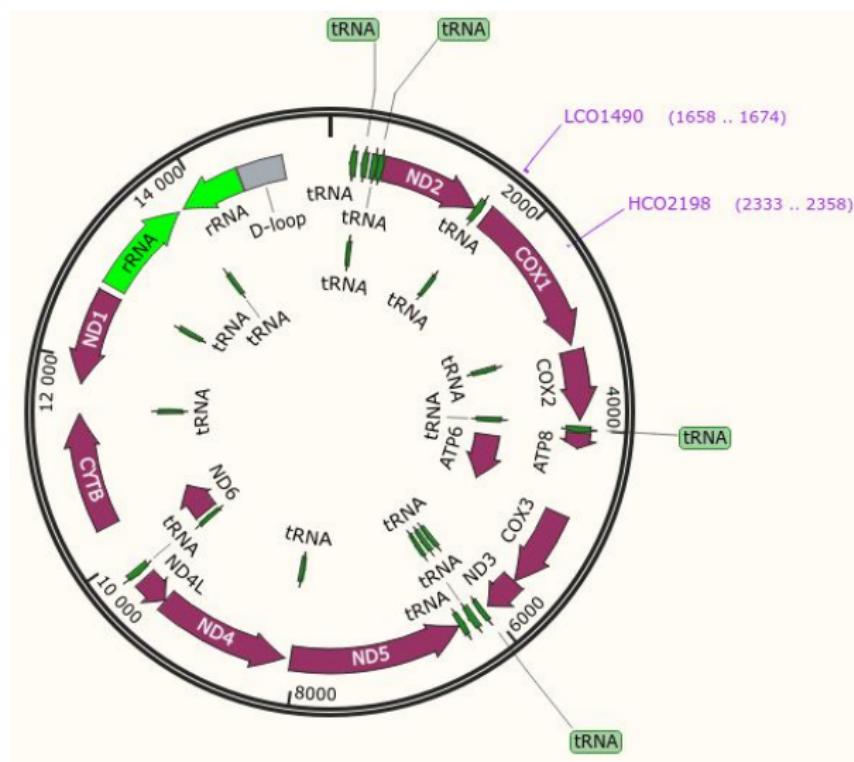
## Примеры часто используемых маркеров: протисты



## Примеры часто используемых маркеров: протисты



# Примеры часто используемых маркеров: животные



# Структура локуса рДНК

Bacteria / Archaea:



Eukaryota:

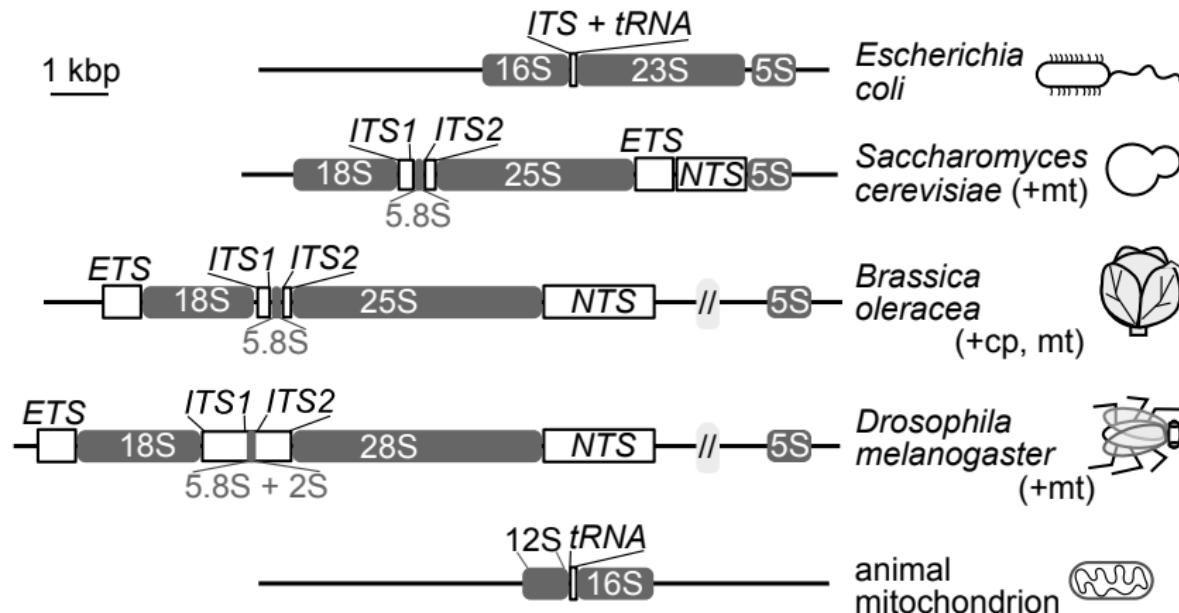


ITS = internal transcribed spacer

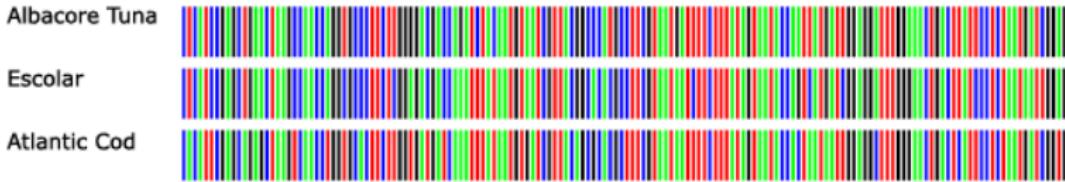
ETS = external transcribed spacer

IGS / NTS = intergenic / non-transcribed spacer

# Структура локуса рДНК у разных организмов

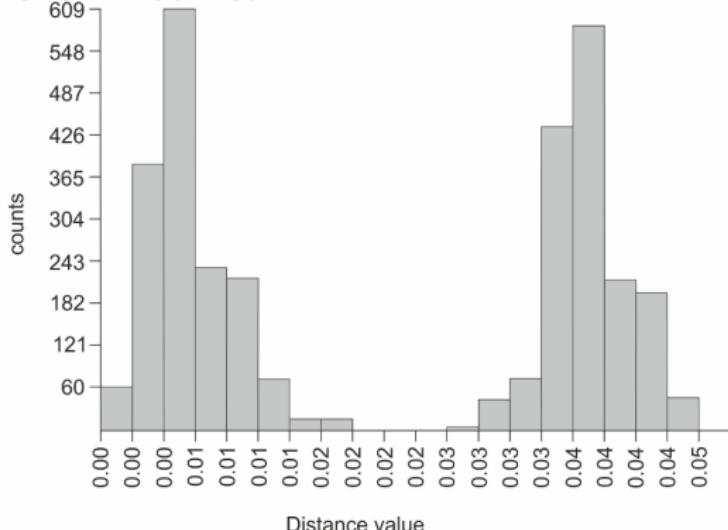


# Штрихкодирование ДНК (DNA barcoding)



<http://www.ibolproject.org/research.php>

Barcode gap — разрыв между внутривидовой и межвидовой изменчивостью.



Markers and phylogenetics

oooooooooooo●oooooooo

Phylogenomics intro

oooooooooooo

Single copy orthologs

oooooooooooo

Reduced representation

ooooo

Alignment-free methods

oooo

Outro

o

# Проект Barcode of Life



The image shows the homepage of the Barcode of Life Data System (BOLD). The header features the text "BOLD SYSTEMS" in white and orange, with a magnifying glass icon and a menu icon. The main title "BARCODE OF LIFE DATA SYSTEM v4" is displayed prominently in white over a background of a world map filled with silhouettes of various animals. Below the title, the tagline "Advancing biodiversity science through DNA-based species identification." is written in white. A large orange button at the bottom center contains the text "EXPLORE THE DATA".

<https://boldsystems.org/>

Markers and phylogenetics  
oooooooooooo●oooo

Phylogenomics intro  
oooooooooooo

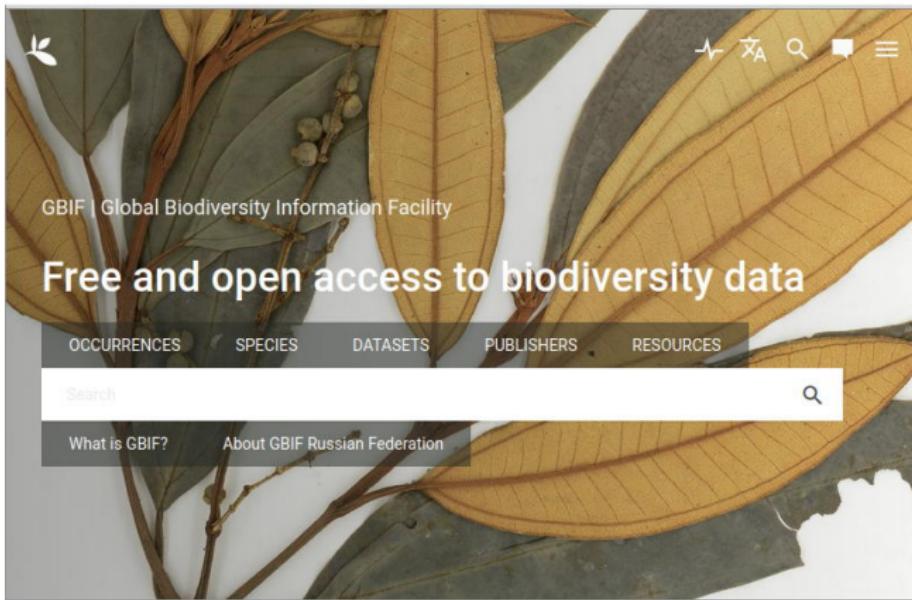
Single copy orthologs  
oooooooooooo

Reduced representation  
oooooo

Alignment-free methods  
oooo

Outro  
o

# База GBIF

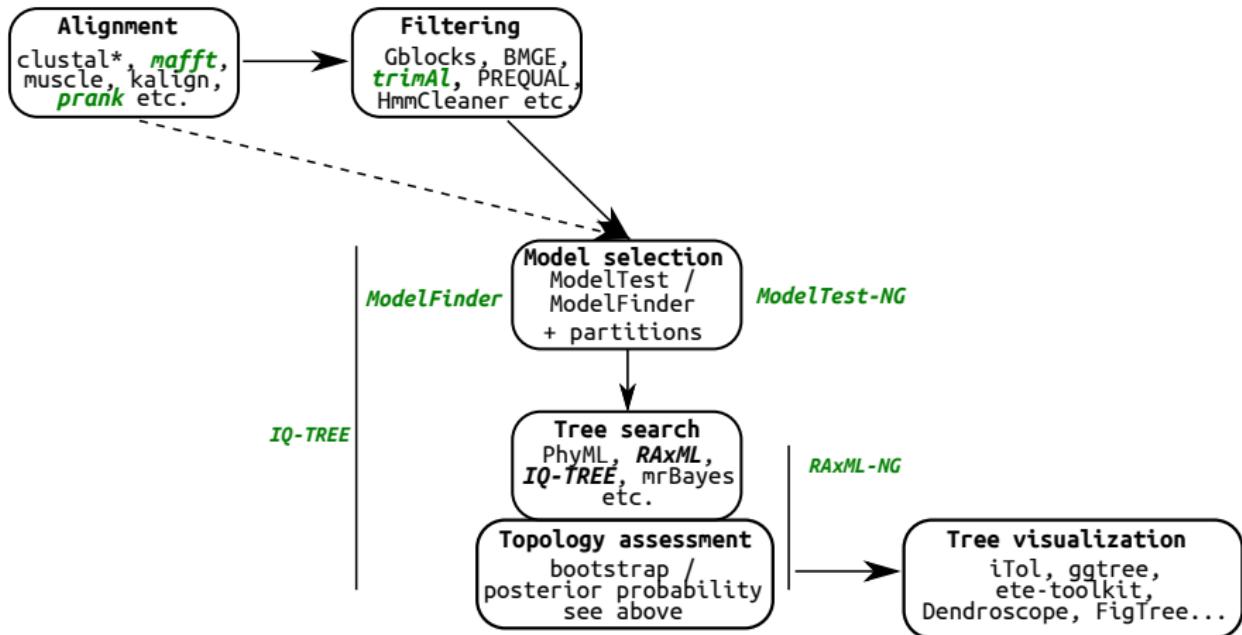


iNaturalist

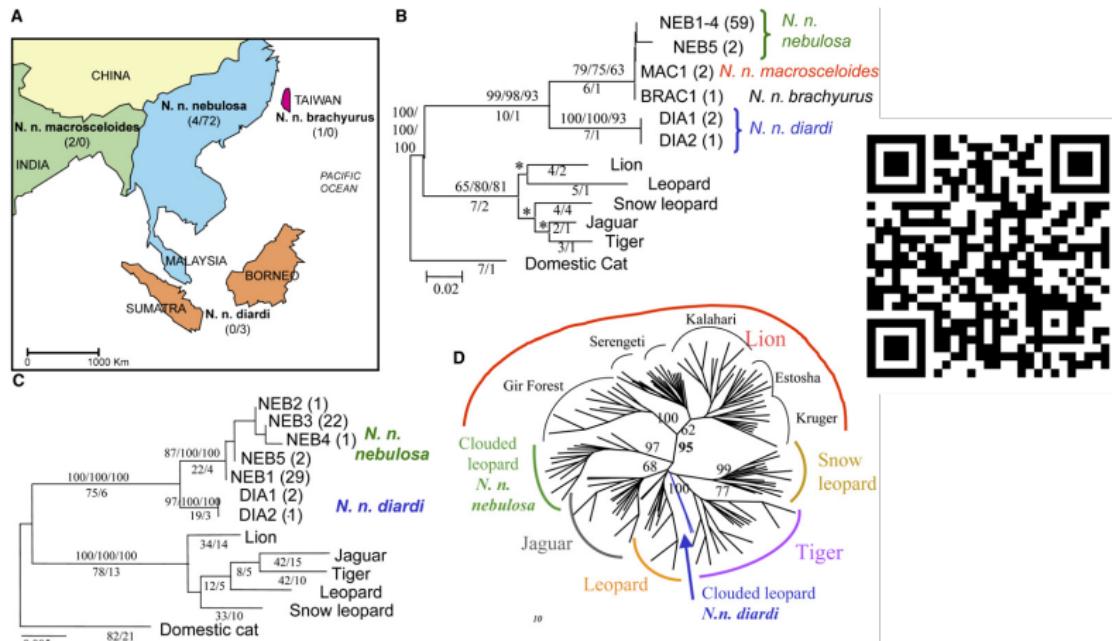
BOLD SYSTEMS

<https://www.gbif.org/>

# Общая схема филогенетического анализа



# Упражнение №1: восстановление филогении по фрагменту гена (*cytb*)



# Упражнения: филогения 12 видов млекопитающих

- ✓ *Physeter catodon*
- ✓ *Ovis aries*
- ✓ *Equus caballus*
- ✓ *Felis catus*
- ✓ *Canis lupus familiaris*
- ✓ *Myotis lucifugus*
- ✓ *Mus spretus*
- ✓ *Mus musculus*
- ✓ *Marmota marmota*
- ✓ *Macaca mulatta*
- ✓ *Ornithorhynchus anatinus*
- ✓ *Vombatus ursinus*

# Упражнения: филогения 12 видов млекопитающих



Markers and phylogenetics

oooooooooooooo●

Phylogenomics intro

oooooooooo

Single copy orthologs

oooooooooo

Reduced representation

oooo

Alignment-free methods

oooo

Outro

o

## Упражнение №2: филогения 12 видов млекопитающих по одному гену (*cytb*)

Markers and phylogenetics  
oooooooooooooooooooo

Phylogenomics intro  
●oooooooooooo

Single copy orthologs  
oooooooooooo

Reduced representation  
oooooo

Alignment-free methods  
oooo

Outro  
o

# Филогенетика и филогеномика

Филогенетика: один или несколько генов.

Филогеномика: много генов (полный геном / транскриптом или его существенная часть)

# Филогенетика и филогеномика

Филогенетика: один или несколько генов.

Филогеномика: много генов (полный геном / транскриптом или его существенная часть)

Зачем использовать много генов?

- ✓ Иногда дерево генов не соответствует дереву видов...

Markers and phylogenetics  
oooooooooooo

Phylogenomics intro  
○●○○○○○○○○

Single copy orthologs  
○○○○○○○○○○

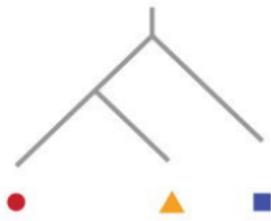
Reduced representation  
○○○○

Alignment-free methods  
○○○○

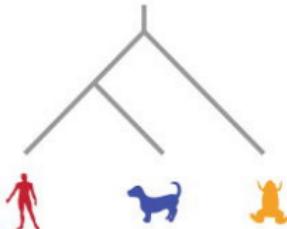
Outro  
○

# Потеря гомологов

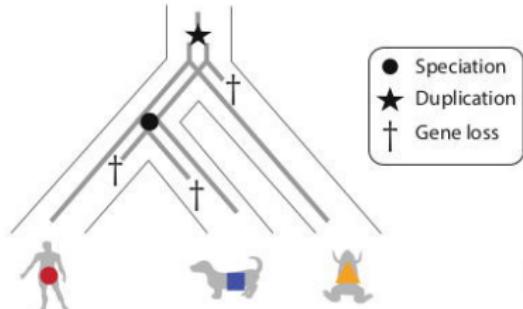
Gene Tree



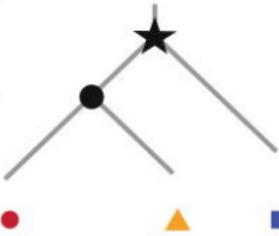
Species Tree



Reconciled Tree  
(Full Representation)



Reconciled Tree  
(Simple Representation)



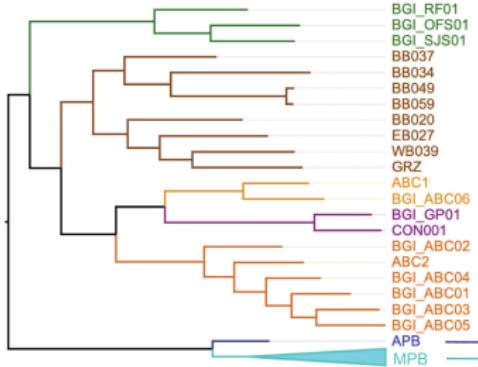
https:

//www.researchgate.net/publication/334265032\_Inferring\_Orthology\_and\_Paralogy

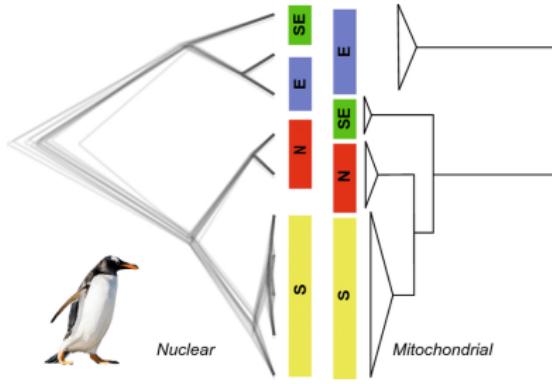
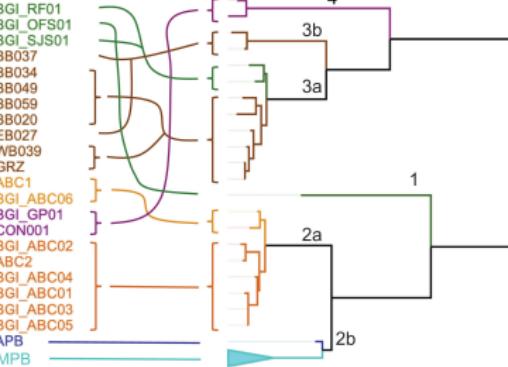


# Mito-nuclear discordance

A Autosomal SNPs



Mitochondrial genomes



Markers and phylogenetics  
○○○○○○○○○○○○○○

Phylogenomics intro  
○○○●○○○○○

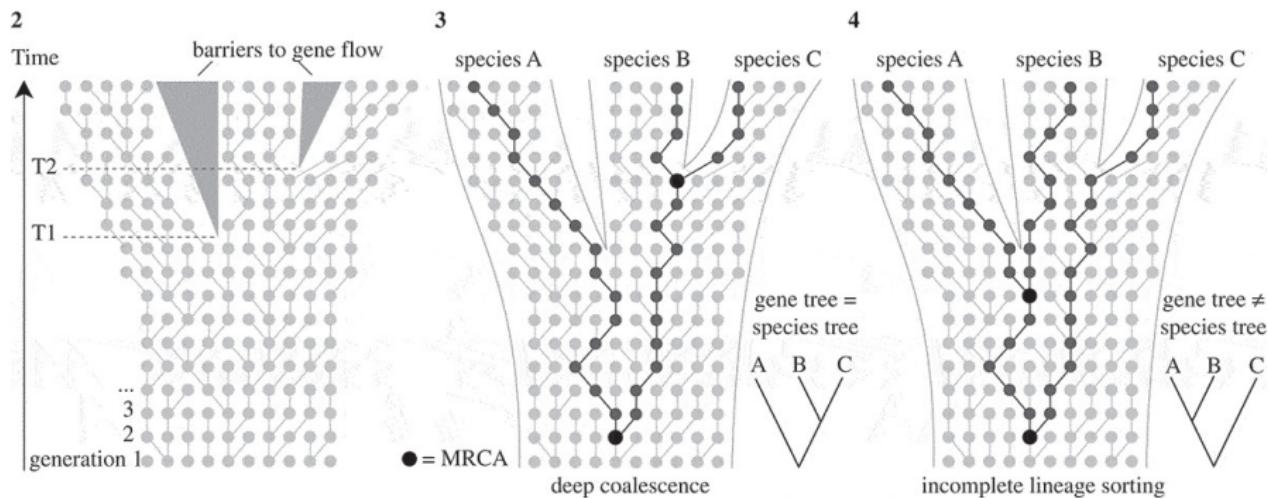
Single copy orthologs  
○○○○○○○○

Reduced representation  
○○○○

Alignment-free methods  
○○○○

Outgroups  
○

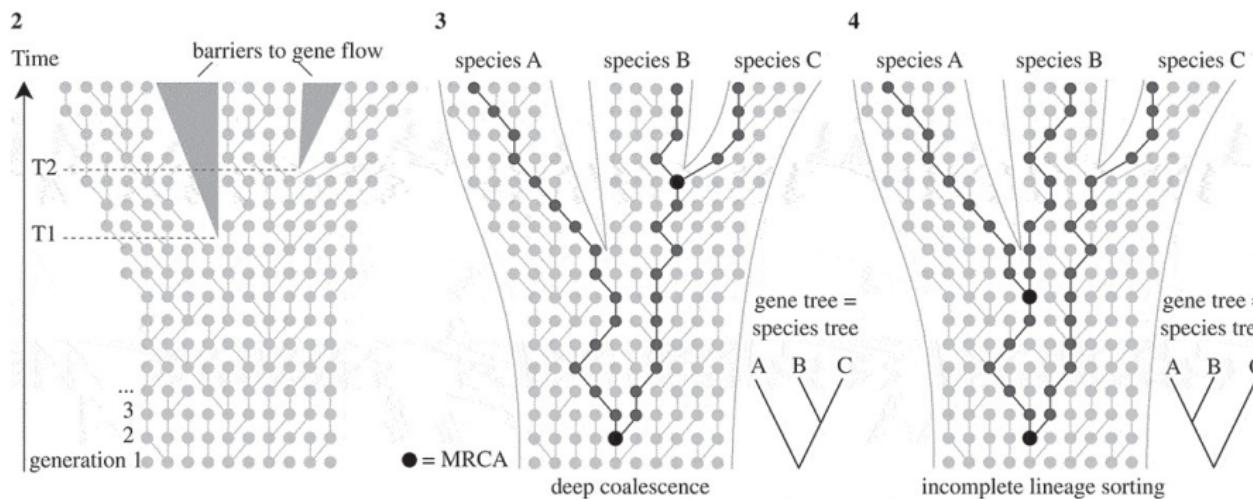
# Incomplete lineage sorting



Leliaert et al., 2014,

<https://dx.doi.org/10.1080/09670262.2014.904524>

# Incomplete lineage sorting



Leliaert et al., 2014,

<https://dx.doi.org/10.1080/09670262.2014.904524>

Mito-nuclear discordance часто именно ILS

Markers and phylogenetics  
oooooooooooooooooooo

Phylogenomics intro  
oooo●oooo

Single copy orthologs  
oooooooooooo

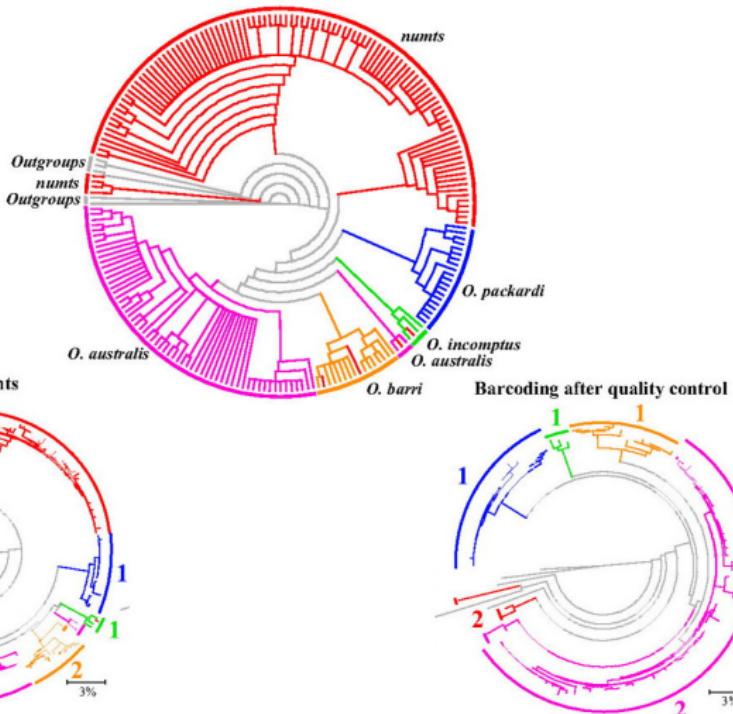
Reduced representation  
ooooo

Alignment-free methods  
oooo

Outro  
o

# NUMTs = nuclear mitochondrial DNA segment

[new might]



Markers and phylogenetics  
oooooooooooooooooooo

Phylogenomics intro  
oooooo●oooo

Single copy orthologs  
oooooooooooo

Reduced representation  
ooooo

Alignment-free methods  
oooo

Outro  
o

# Филогенетика и филогеномика

Филогенетика: один или несколько генов.

Филогеномика: много генов (полный геном /  
транскриптом или его существенная часть)

# Филогенетика и филогеномика

Филогенетика: один или несколько генов.

Филогеномика: много генов (полный геном / транскриптом или его существенная часть)

Зачем использовать много генов?

- ✓ Иногда дерево генов не соответствует дереву видов...
- ✓ Лучшее разрешение в близких таксонах.
- ✓ Лучшее разрешение в макротаксономии
- ✓ (Эти данные у нас уже есть...)

# Упражнения: филогения 12 видов млекопитающих



Markers and phylogenetics  
oooooooooooooooo

Phylogenomics intro  
oooooooo●○

Single copy orthologs  
oooooooo

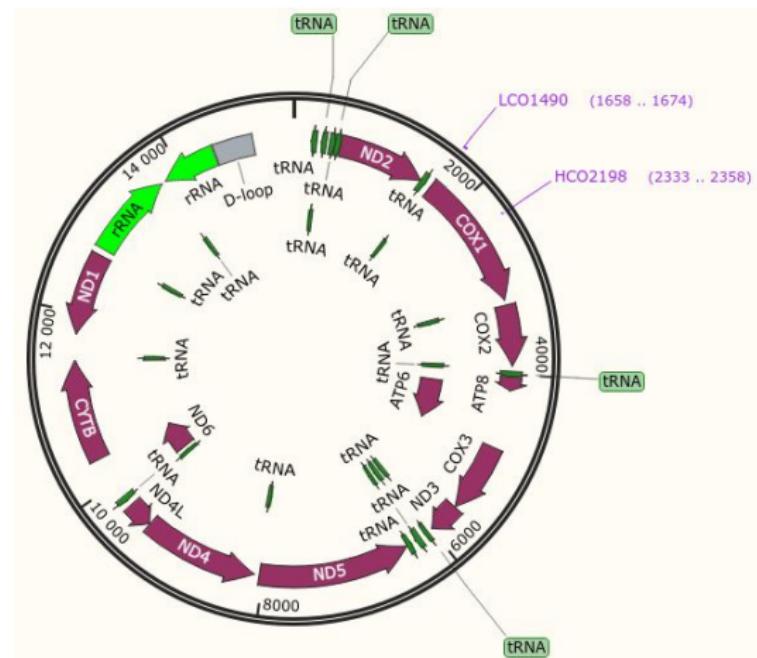
Reduced representation  
oooo

Alignment-free methods  
oooo

Outro  
o

## Упражнение №2: восстановление филогении по 13 митохондриальным белкам

# Митохондриальный геном животных



1. *ATP6*
2. *ATP8*
3. *COX1*
4. *COX2*
5. *COX3*
6. *CYTB*
7. *ND1*
8. *ND2*
9. *ND3*
10. *ND4*
11. *ND4L*
12. *ND5*
13. *ND6*

Markers and phylogenetics  
ooooooooooooooo

Phylogenomics intro  
oooooooooo

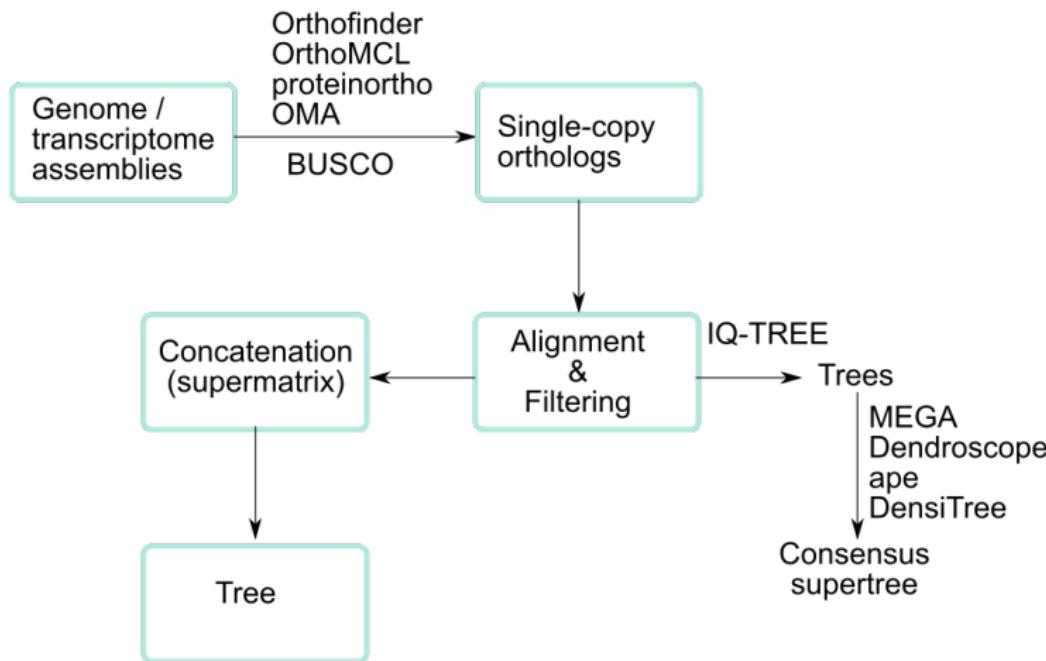
Single copy orthologs  
●oooooooo

Reduced representation  
oooo

Alignment-free methods  
oooo

Outro  
o

## Single copy orthologs



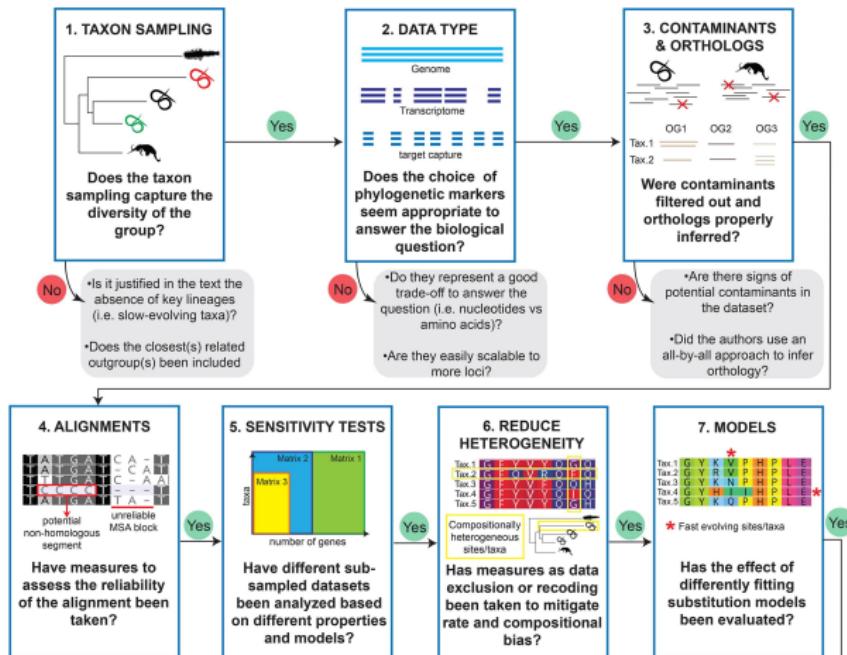
Работает как с геномами, так и с транскриптомами

# Recommended!

## A Practical Guide to Design and Assess a Phylogenomic Study

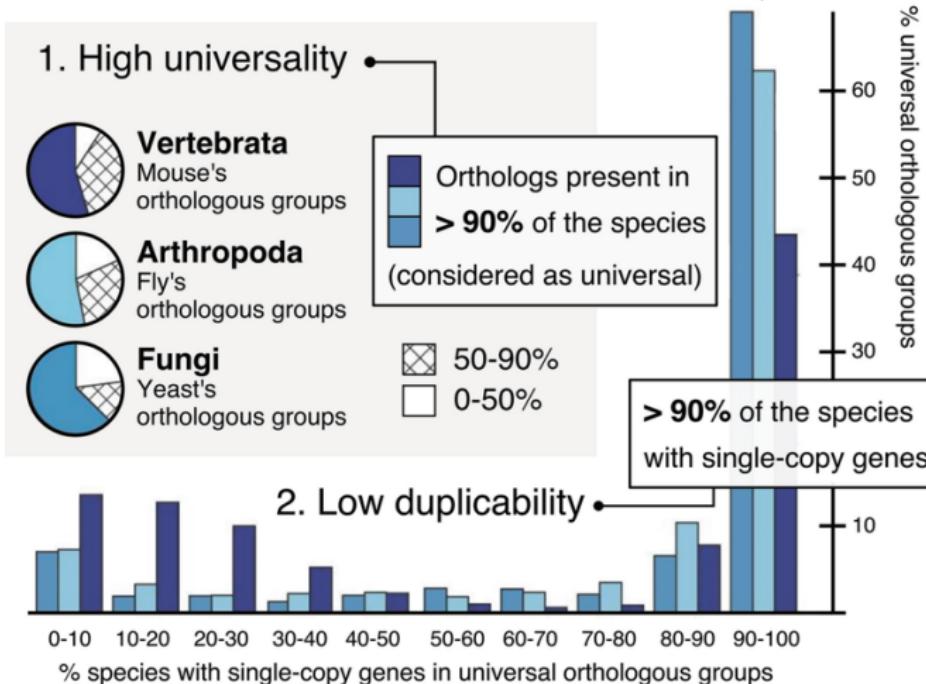
Jesus Lozano-Fernandez ①,2,\*

GBE



# BUSCO

## BUSCO sampling space



Markers and phylogenetics  
oooooooooooooooooooo

Phylogenomics intro  
oooooooooooo

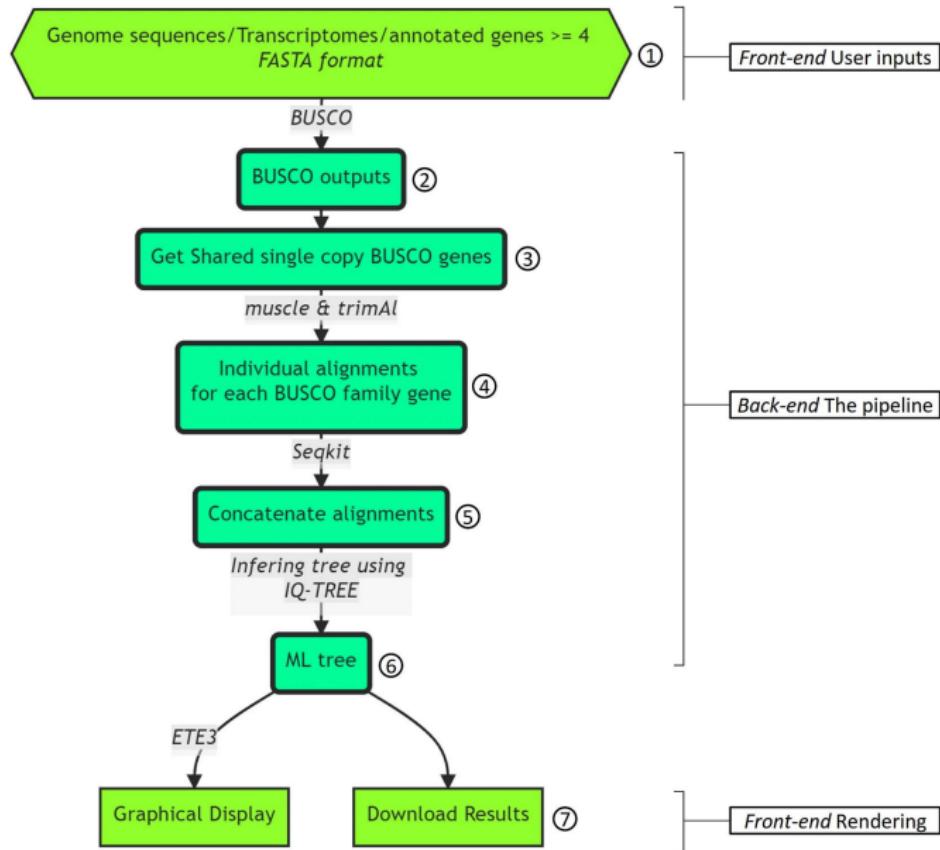
Single copy orthologs  
ooo●oooooo

Reduced representation  
ooooo

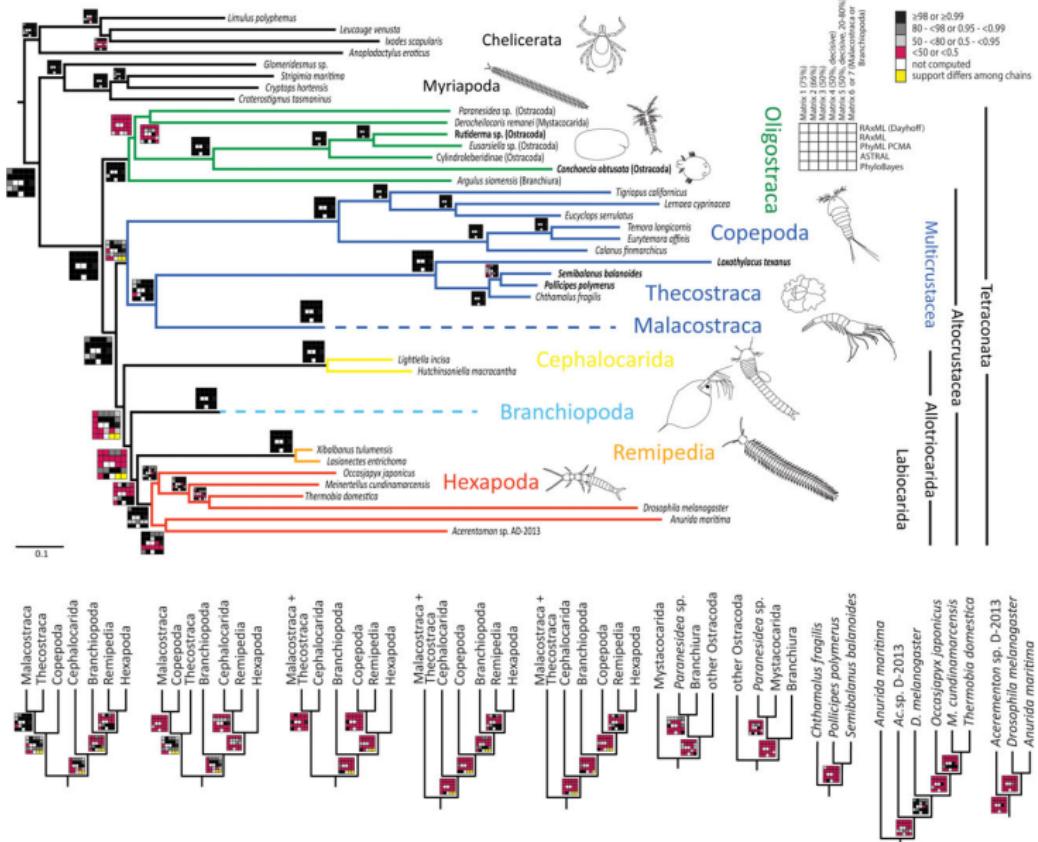
Alignment-free methods  
oooo

Outro  
o

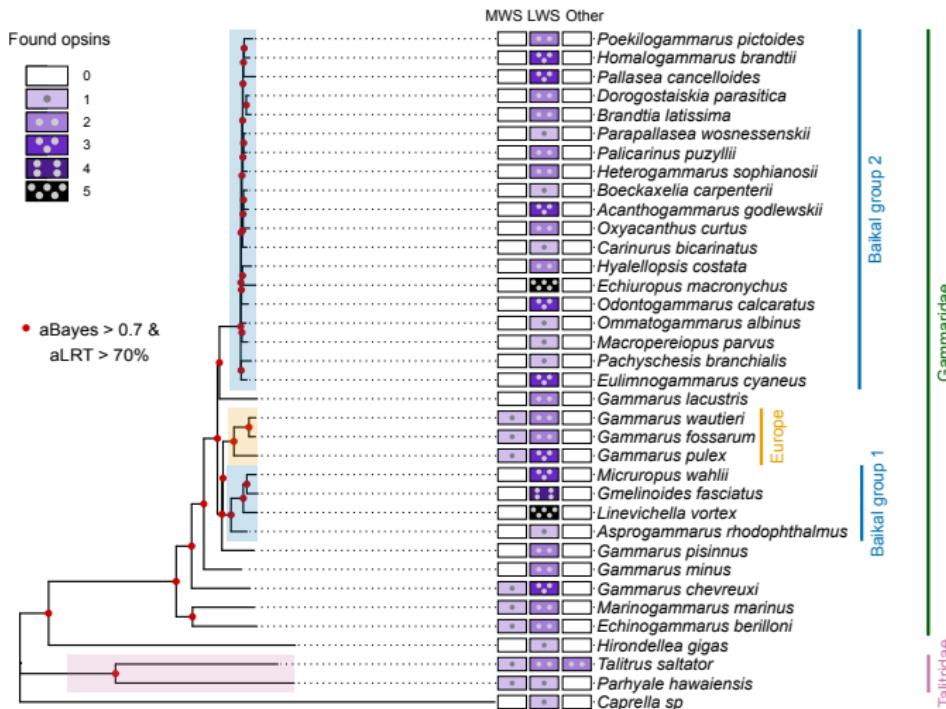
# BuscoPhylo



## Single copy orthologs: examples



# Single copy orthologs: examples



Drozdova et al., 2021, <https://doi.org/10.1186/s12862-021-01806-9>

[https://github.com/AlenaKizenko/pia3\\_amphipod\\_opsins/blob/master/other\\_scripts/2.3\\_species\\_tree.sh](https://github.com/AlenaKizenko/pia3_amphipod_opsins/blob/master/other_scripts/2.3_species_tree.sh), 2.4\_Fig2\_S2B\_species\_tree.R

## Markers and phylogenetics

Phylogenomics intro

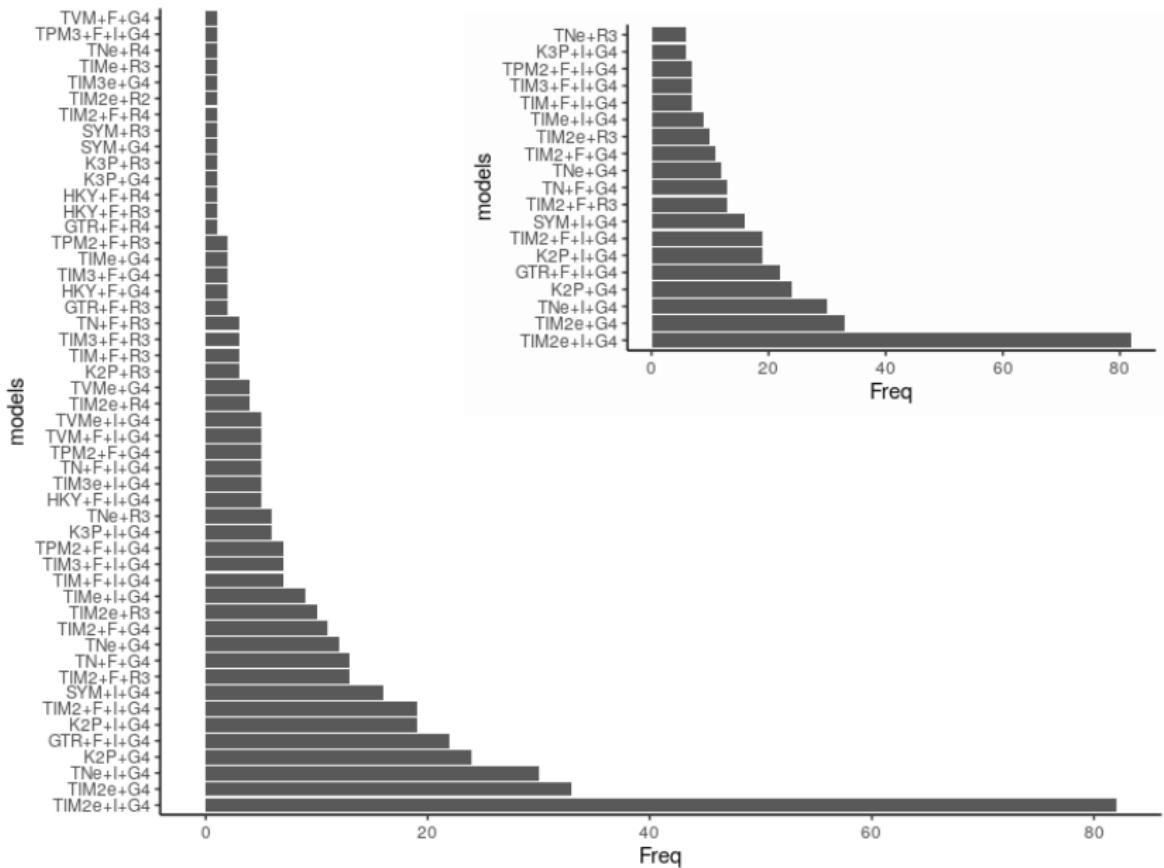
### Single copy orthologs

## Reduced representation

## Alignment-free methods

Outro

## Statistics on best models





Markers and phylogenetics  
oooooooooooooooooooo

Phylogenomics intro  
oooooooooooo

Single copy orthologs  
oooooooo●

Reduced representation  
ooooo

Alignment-free methods  
oooo

Outro  
o

## Упражнение 3: филогения 12 видов млекопитающих по однокопийным ортологам

Markers and phylogenetics  
oooooooooooooooooooo

Phylogenomics intro  
oooooooooooo

Single copy orthologs  
oooooooooooo

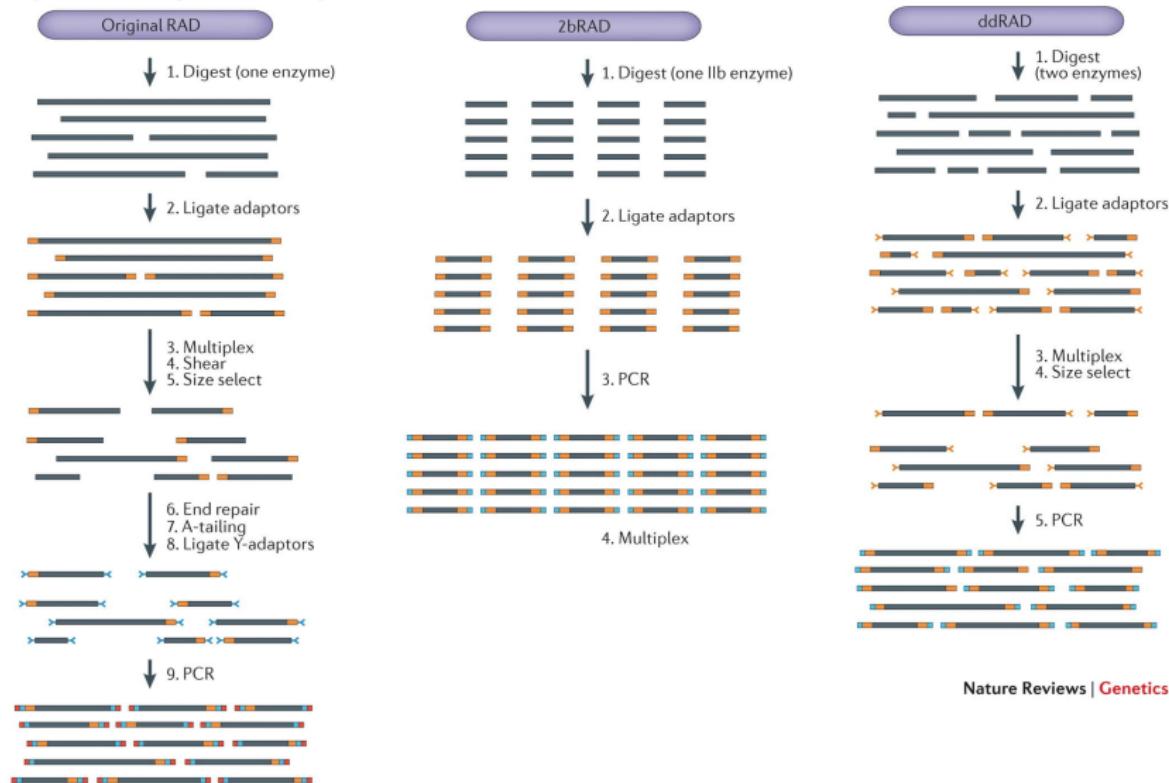
Reduced representation  
●ooooo

Alignment-free methods  
oooo

Outro  
o

# RADseq

Sequence next to single restriction enzyme cut sites



Nature Reviews | Genetics

## Markers and phylogenetics

## Phylogenomics intro

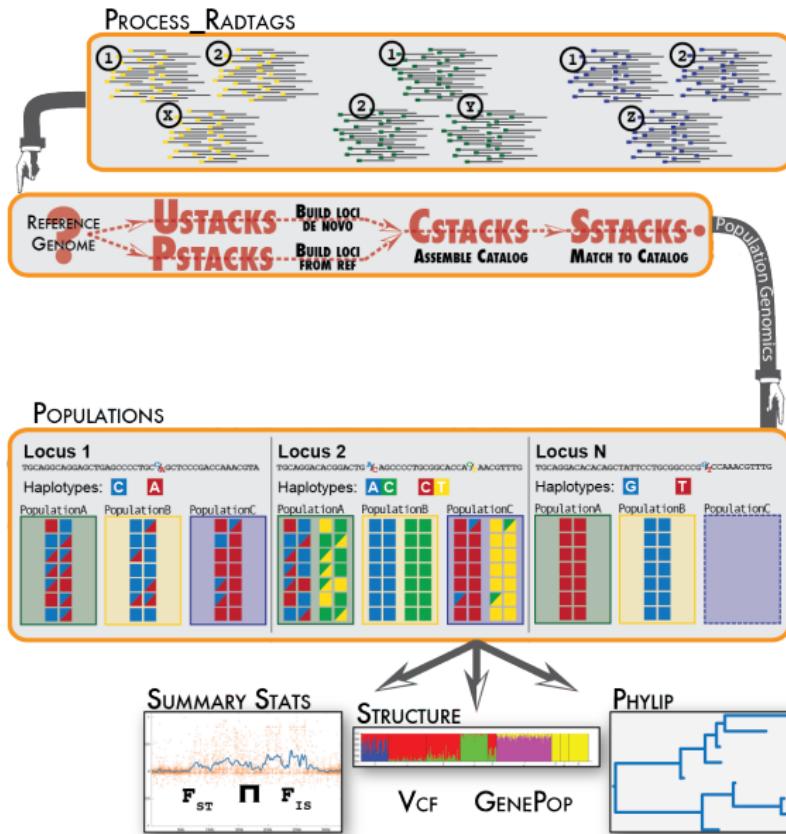
## Single copy orthologs

## Reduced representation

## Alignment-free methods

Outro

## RADseq: analysis



## STEP

### 01 PROBE SELECTION

#### DNA GENOME

Genome sequencing if not available  
Probe selection/design

## STEP

### 02 DNA EXTRACTION

#### SPECIMENS

Pinned specimen (<20 years)  
Specimen in ethanol/low temperature

### 03 UCE LIBRARY PREPARATION

#### DNA EXTRACT

DNA shearing  
Adaptor ligation  
PCR amplification  
Hybridization enrichment  
qPCR quantification  
DNA size selection

### 04 SEQUENCING

#### POOLED SAMPLES

MiSeq for small projects  
HiSeq, NextSeq, and NovaSeq for large projects

### 05 ILLUMINA READ PROCESSING

#### RAW READS

Illumiprocessor /Trim Galore!

### 06 CONTIG ASSEMBLY

#### CLEAN READS

ABySS/Trinity/Velvet/SPAdes

### 07 IDENTIFY TARGETS

#### CONTIGS

Match contigs to probe set  
Extract UCE loci

### 08 a ALIGNMENT

MUSCLE/MAFFT

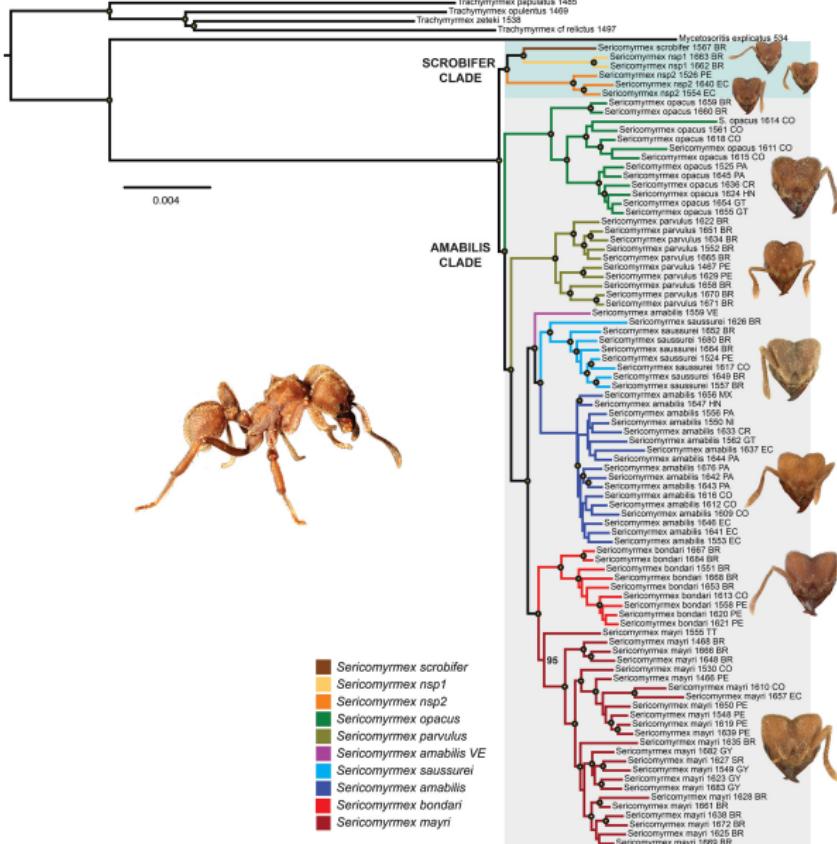
#### SELECTED CONTIGS

### 08 b PHASING

Reference Assembly  
Phase allele



# Ultraconserved elements = UCEs (example)



Markers and phylogenetics  
oooooooooooooooooooo

Phylogenomics intro  
oooooooooooo

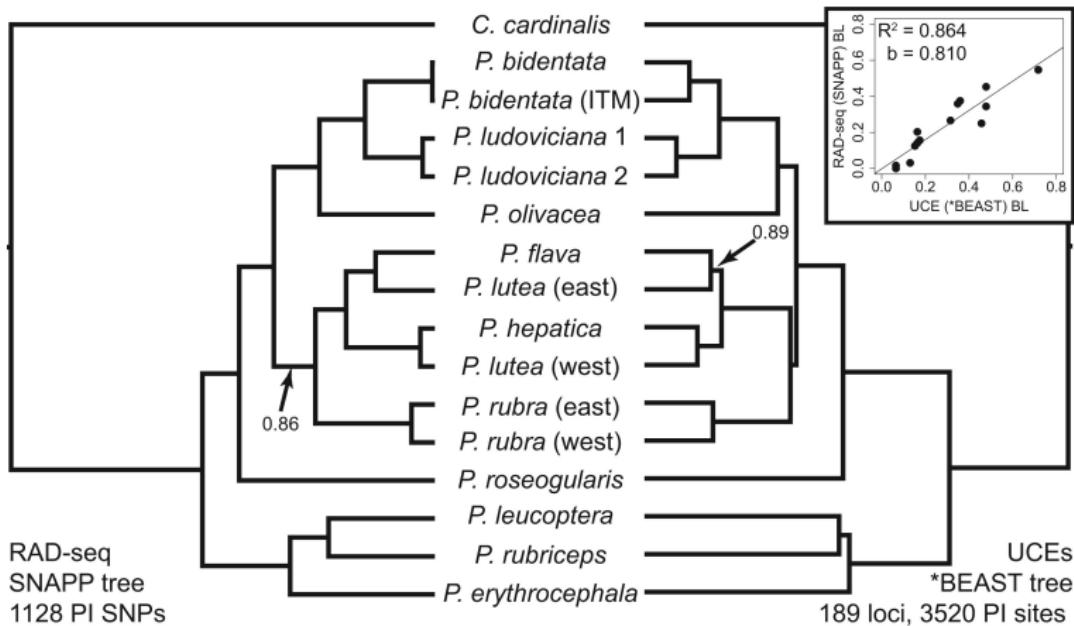
Single copy orthologs  
oooooooooooo

Reduced representation  
ooo●○

Alignment-free methods  
oooo

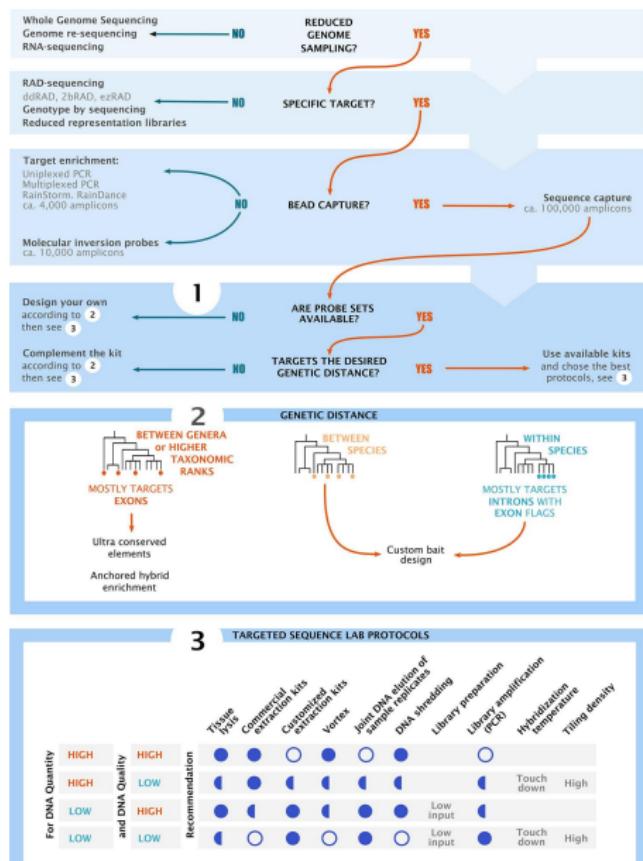
Outo  
o

# RAD-seq and UCE: example



<https://academic.oup.com/sysbio/article/65/4/640/1753369>

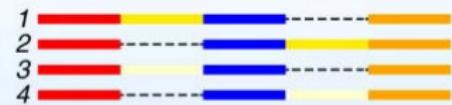
## Recommended!



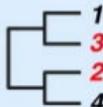
# Alignment-free methods



## A classical approach

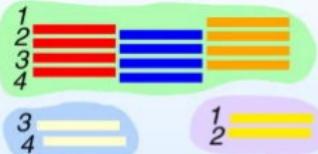


phylogenetic inference

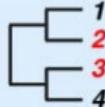


phylogenetic tree

## B alternative approach

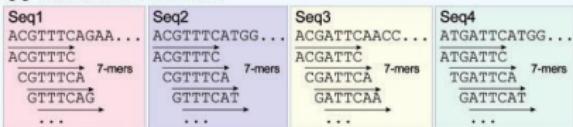


phylogenetic inference

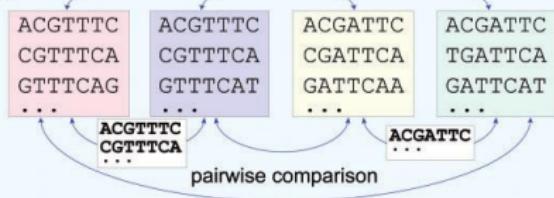


## Alignment-free methods

## A extraction of $k$ -mers



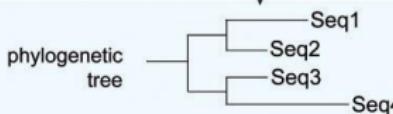
B



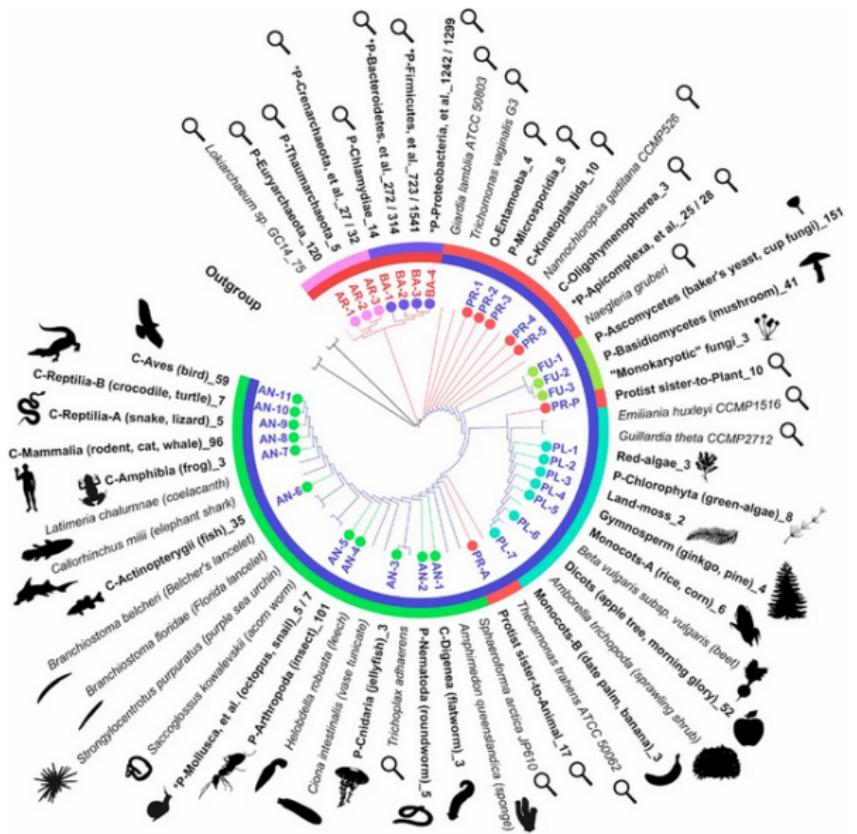
C

	Seq1	Seq2	Seq3	Seq4
pairwise distance matrix	Seq1	0.0		
	Seq2	0.3	0.0	
	Seq3	0.5	0.4	0.0
	Seq4	0.7	0.6	0.4
				0.0

D



## Alignment-free methods



Markers and phylogenetics  
oooooooooooooooooooo

Phylogenomics intro  
oooooooooooo

Single copy orthologs  
oooooooooooo

Reduced representation  
oooooo

Alignment-free methods  
ooo●

Outro  
o

## Упражнение №5: восстановление филогении 12 видов млекопитающих по аминокислотному составу протеомов

Markers and phylogenetics



Phylogenomics intro



Single copy orthologs



Reduced representation



Alignment-free methods



Outro



# Интегративная таксономия

