

# "CRISPR-associated protein 1 ( Cas1)"

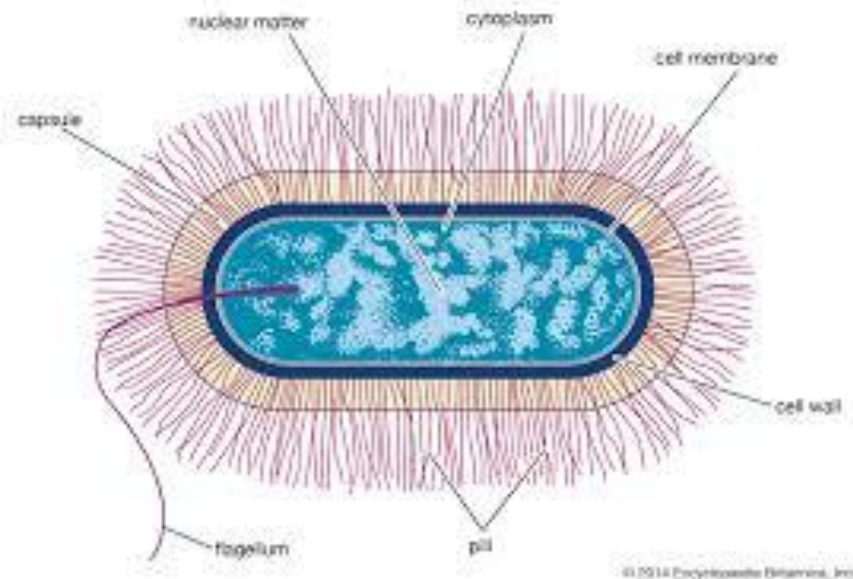
Molecular phylogenetics course project

---

Sedreh Nassirnia  
Fall 2019

# What is the story?

- Viruses catch bacterial cell in order to propagate
  - results in bacterial cell lysis and death
- Clustered regularly interspaced short palindromic repeat (CRISPR) and CRISPR associated (Cas) proteins immune defense system is the solution of bacteria to fight with viruses



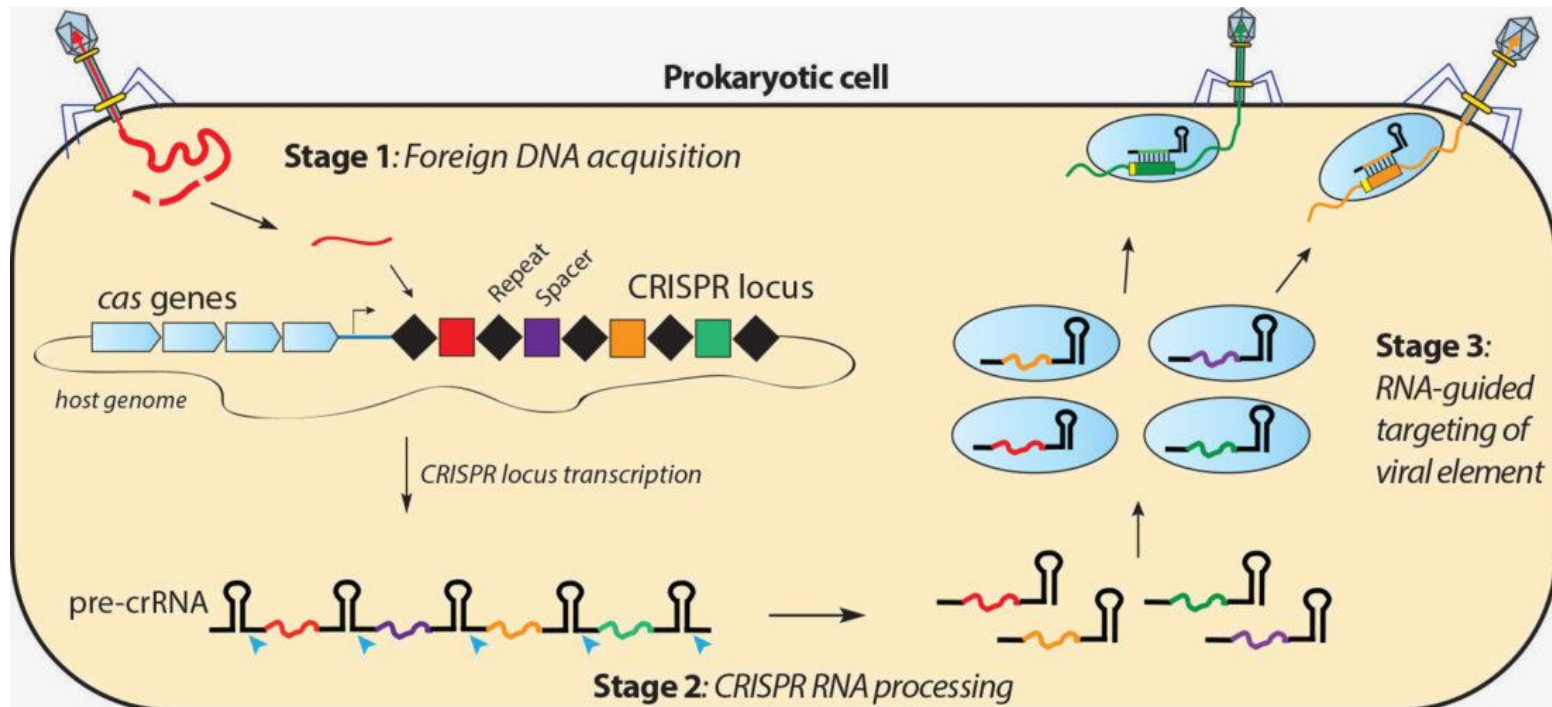
# CRISPR-Cas Systems



- **CRISPR (clustered regularly interspaced short palindromic repeat)**
  - **adaptive immune system that provides protection against mobile genetic elements such as genomic islands, plasmids, and transposon-like elements and viruses**
- **Important for**
  - **clinical microbiologists, ecologists and evolutionary biologists**
  - **CRISPR-Cas system potential uses**
    - **Detection and genotyping of microbial pathogens**
    - **Host identification in metagenomes**
    - **Analysis of viral genomes**
    - **Targeted genome engineering in both prokaryotic and eukaryotic cells**

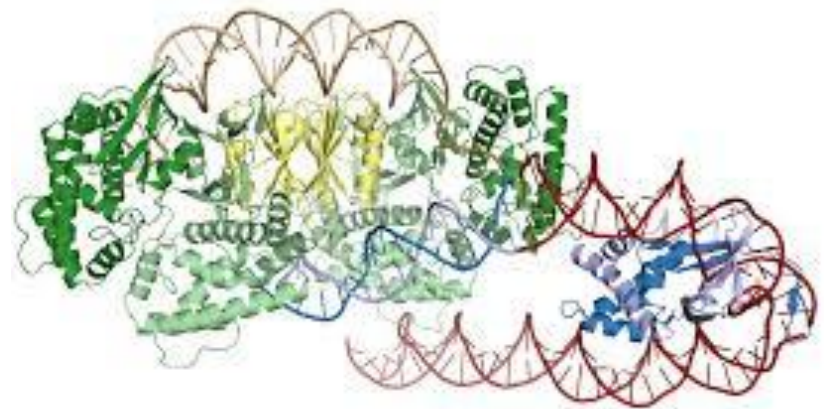
# CRISPR associated protein Cas1

- Cas1 responsible for the ability of the **CRISPR immune system** in **bacteria** to adapt to new **viral infections**
  - Identify the site in the genome where they insert **viral DNA**



# Objective

- Study evolutionary relationships among bacteria family based on cas1
- What percentage of organisms evolved according to viral strains






# Data

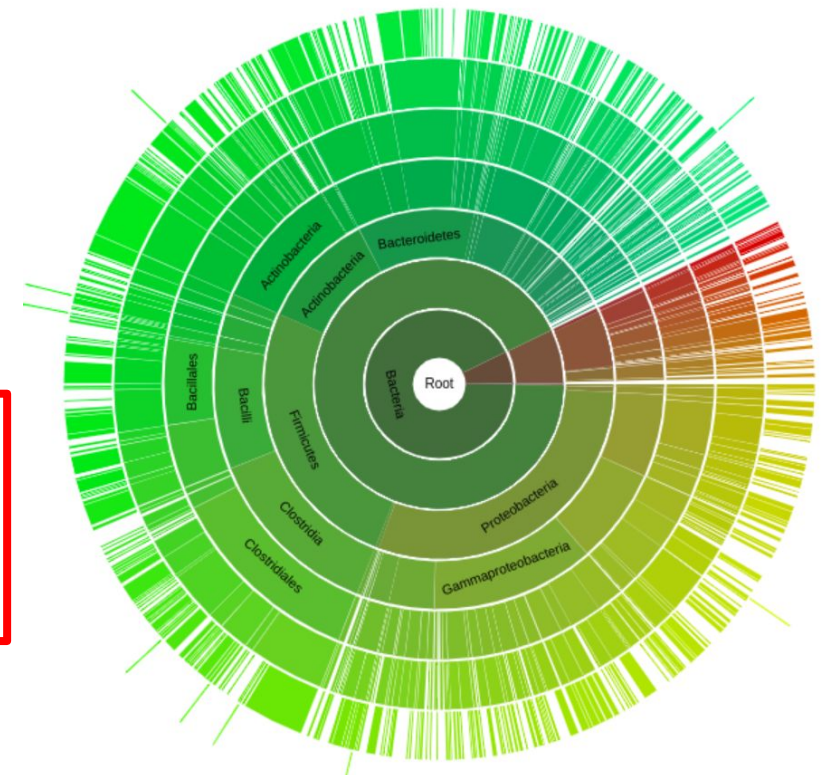
Downloading directly from Pfam doesn't give information about Genus and Species

Link database was used to fetch this information.

Example:

[https://www.kegg.jp/entry/azr:CJ010\\_02280](https://www.kegg.jp/entry/azr:CJ010_02280)

 <b>Azoarcus sp. DD4: CJ010_02280</b> <span>Help</span>			
<b>Entry</b>	CJ010_02280	CDS	T06048
<b>Definition</b>	(GenBank) subtype I-C CRISPR-associated endonuclease Cas1		
<b>KO</b>	K15342 CRISP-associated protein Cas1		
<b>Organism</b>	azr Azoarcus sp. DD4		
<b>Brite</b>	KEGG Orthology (KO) [BR:azr00001] 09180 Brite Hierarchies 03102 Protein families: genetic information processing 03400 DNA repair and recombination proteins [BR:azr03400] CJ010_02280 09183 Protein families: signaling and cellular processes 02048 Prokaryotic defense system [BR:azr02048] CJ010_02280 DNA repair and recombination proteins [BR:azr03400] Prokaryotic type DSB (double strand breaks repair) HR (homologous recombination) Other HR factor CJ010_02280 Prokaryotic defense system [BR:azr02048] CRISPR-Cas system Universal Cas proteins CJ010_02280 <a href="#">BRITE hierarchy</a>		
<b>SSDB</b>	<a href="#">Ortholog</a> <a href="#">Paralog</a> <a href="#">Gene cluster</a> <a href="#">GFIT</a>		
<b>Motif</b>	Pfam: Cas_Cas1 <a href="#">Motif</a>		
<b>Other DBs</b>	NCBI-ProteinID: QDF95462 UniProt: A0A4Y6KN73		
<b>LinkDB</b>	<a href="#">All DBs</a>		
<b>Position</b>	463572..464612 <a href="#">Genome map</a>		
<b>AA seq</b>	346 aa <a href="#">AA seq</a> <a href="#">DB search</a> MRRQLNTLYVTTEGAWLKKDGANIVMEVEGAERARLPVHMLSEVMCMGRVMVSPPLLGYC AEOGICVSFLSPNGKFLARMEGPVSGNVLLRREQYRRDDPARCGLVVRNLLIGKVHNR AVLGRALRDHGEPMPEADQALAHARERLRISARLLLLEKLDVLRGLGEAAQAYFGVF DHLIRVPEPALRFKGRSRRPLDAVNALLSFLYTLTTHDCRSALSVGLDPAVGFLHRDR PGRPSLALDLLLEFRPVMADRLALSLINRRQLGERDFVQLDNGAVSLKEESRKTVLTAQ ERKREEMRHAFLLEKFAVGLFPAVQAQLFARHLRGDLEAYPPFLWK		
<b>NT seq</b>	1041 nt <a href="#">NT seq</a> +upstream0 nt +downstream0 nt atgcgccgcagctcaacaccctgtatgtgaccacggaaggcgcttgctgaagaaggat ggcgcgaacatcgctcatggaggtggaaggcggaacgcgcgctttgcccgttcacatg ctggagagcatggtgtgtatgggcccgggtgatggtgctgcgcgcactgctcgctactgc gctgagcagggtattcgctcagctctctctcgcctcaatggcaagttctctggtatgcatg gaggggccggtttctggaacgtgctgctacggcgcgagcaataaccgcccgtaccgatg ccggcccgcgtgtggtctagtcgctcgttaactctgttgatcggaaggtacataaccaacgc gcggtgcttgcccgctgcgctgcgcatcacggtgagggcacgacggaagcgatcagacc gccctggcgacgcggggagcgctgcggcgcatctccgcgcgttgttgctggaagag aagctggatgtgtgctggcctggaaggagaggcggcaggcctatttcggcgctctc gatcacctcatcgctgcccggagcctgcttgcctttaaggggcgagcgccggcgccg ccactggacgcagtgatgcttctcctctctacacgttgcgtgacgcatgactgc		



# Fetching IDs

```
import os
import urllib
from bs4 import BeautifulSoup as bs
import sys
import requests
import re

def write_fasta(genus, species, ids, seq):
    with open('/home/sedreh/ITM0/semester3/Molecular_phylogenetic/COURSE_PROJECT/phyloproject_cas1/cas1_pfam.fasta'
              as fa, 'a') as fa:
        fa.write(">{} {} {}{}\n{}\n".format(genus, species, ids, seq))

def fetch_fasta(link):
    page = requests.get(link)
    soup = bs(page.content, 'html.parser')
    children = (list(soup.children)[2]).get_text()
    name = re.findall("(?<=\\xa0\\xa0\\xa0)(.*)?(?=\\:)", children)[0]
    genus, species = format(re.sub(r'\\d-', '', name)).split(" ", 1)
    species = species.strip(" ")
    ids = re.findall("(?s)(?<=UniProt:\\xa0)(.*)?(?=\\n)", children)
    seq = re.findall("(?s)(?<=aa \\n)(.*)?(?=\\nNT seq\\n)", children)
    seq = format(re.sub('\\n', '', seq[0]))

    return write_fasta(genus, species, ids, seq)

def main():
    list_of_links = list()
    with open('/home/sedreh/ITM0/semester3/Molecular_phylogenetic/COURSE_PROJECT/phyloproject_cas1/links.txt', 'r') as f:
        for link in f:
            list_of_links.append(link.strip('\\n'))

    for link in list_of_links:
        fetch_fasta(link)

aaa:Acav_0268      K15342 CRISP-associated protein Cas1 | (GenBank) CRISPR-associated protein Cas1
aaa:Acav_3874      K15342 CRISP-associated protein Cas1 | (GenBank) CRISPR-associated protein Cas1
aac:Aaci_2651      K15342 CRISP-associated protein Cas1 | (GenBank) CRISPR-associated protein Cas1
aacn:AANUM_1357    K15342 CRISP-associated protein Cas1 | (GenBank) CRISPR-associated Cas1 family protein
aacn:AANUM_1920    K15342 CRISP-associated protein Cas1 | (GenBank) cas1-2; CRISPR-associated protein cas1
aact:ACT75_00725   K15342 CRISP-associated protein Cas1 | (GenBank) type I-C CRISPR-associated endonuclease Cas1
aact:ACT75_02270   K15342 CRISP-associated protein Cas1 | (GenBank) type I-F CRISPR-associated endonuclease Cas1
aad:TC41_2954      K15342 CRISP-associated protein Cas1 | (GenBank) CRISPR-associated protein Cas1
aae:aq_369         K15342 CRISP-associated protein Cas1 | (RefSeq) hypothetical protein
aal:EP13_09450     K07486 transposase | (GenBank) transposase
aalg:AREALGMS7_00764 K03832 periplasmic protein TonB | (GenBank) transport protein TonB
aan:D7S_00182      K15342 CRISP-associated protein Cas1 | (GenBank) CRISPR-associated protein Cas1
aan:D7S_00548      K15342 CRISP-associated protein Cas1 | (GenBank) CRISPR-associated protein Cas1
aao:ANH9381_1474   K15342 CRISP-associated protein Cas1 | (GenBank) CRISPR-associated protein Cas1
aap:NT05HA_0335    K15342 CRISP-associated protein Cas1 | (GenBank) crispr-associated protein Cas1
aaqu:D3M96_05980   K15342 CRISP-associated protein Cas1 | (GenBank) cas1f; type I-F CRISPR-associated endonuclease Cas1
aar:Acear_0821     K15342 CRISP-associated protein Cas1 | (GenBank) CRISPR-associated protein Cas1
aat:D11S_1150      K15342 CRISP-associated protein Cas1 | (GenBank) CRISPR-associated protein Cas1
```



# Cleaning data in R

- Raw data
  - 1413 organisms
  - More than 4000 sequences
- Clean and filtered data
  - Each genera must contain at least 3 entries

## Final dataset:

- 47 genera
- 320 entries

```
>Capnocytophaga stomatis
MLYRSIYIGNPAYLKLKDQQMKIVCPETKAEGSVPVEDLGLMLDHFQITISHQLIQWLMGNNVVIISCD AHLPHGQMLPLHGNAIYSQRIKDQIEASEPLKKQLWKQTIEC
>Corynebacterium striatum
MAYSEDAITFSTIPADHQVRLIEDRVSFAYVEHAAIRQDRTGVVAYSVDNSELEQRIQLPVGGLAVLMLGPGTSSISAAAATSCTRSGTTIMFTGGGGVPAYTHAASLTSSARWA
>Candidatus Caldiarchaeum
MSELVIDKPGTYLGVKGLFVVRTKGGGRSEFSPVELSHSIRCRGVGVSDALRLACRFGIEVSVYSRGRPVGKVVGAFLLGGGAVTRRAQLEAWGTERGLAVAREIVSAKLYN
>Corynebacterium singulare
MTTPHEVPLTRQALARVGDRISFLYAERCVINRDGNSLTIVDQRGTAHV PATQIAALLLGPGTKITYAAMALLGDAGVSAVWVG ERGVRYAHGRPPAKSSRMAEIQAEVVTHQ
>Clostridium tetani
MKRSYIYNNGILKRKDN SMAFIDELGERRYIPIETANEIYVMSEMDFNTSLINYL SQYDVIHFFFNYSFYTG SFQPRKKLVSGNLLVNQVNHYS DNSKRLEIAKKFVDGASY
>Chlorobaculum tepidum
MKKHLNTLFVTTQGSYLSKEGECVLISIDRVEKTRIPLHMLNGIVCFGQVSCSPFLLGHCAQLGVAVTFLTEHGRFLCQM QGPVKGNILLRRAQYRMADNYDQTATLARLFVIG
>Corynebacterium terpenotabidum
MNKIPFRSSVTTQGRNASGARSLWRATGMTDEDFEKPIIAVANSYTQFVPGHVHLKNVGDIVAEAVKEAGGVAREFNTIAVDDGIAMGHSGMLYSLPSREIISDSVEYVMVNAHQ
>Hungateiclostridium thermocellum
MKKSAFIFSDGELKRKDS TVLFESEDSKNYLP IEDISDIYIFGEVTVTKKFLELATQKEILLHFYNYNEYVVG TYPREHYNSGF MILKQAEHYLDEEKRM AIAKKFIHGSVK N
>Cronobacter turicensis
MSFVPLNPIPLNDRTSMIFLQYGHLDVLDGAFVLVDKTGVRTHVPVGAIACIMLEPGTRVSHAAIRLASQVGTLLVWVGEAGVRLYASGQP GGARADKLLYQAKLALDETLRLK
>Corynebacterium urealyticum
MRTPQQVPIERQSLSQMGDRISFLYVERAVVSRDGNALTVTDQRGVAHV PATQLAALLLGTGTRITNAAIALLGDSGVSTVWVG ERGVRYAHGRPPAKSSRLAELQARVVTNQ
>Corynebacterium ulcerans
MSYSNEALAFSTIPASEQIRLEDRVSFLYLEYCLIRQDRTGVIAVSRGDEKAPAEKLDLP IKARIQLPVGGLAVLMLGPGTSSISQPAATSCARAGVSVLFTGGGGVQAYSLS TP
>Corynebacterium ureicelerivorans
MGSRIISFLYIERATVNRDGNALTITDQRSVAHVAATQLAVLLLGPGRITRYAAMALLGDAGVSIWVG ERGVRYASGRPPAKSSRMAELQAEIVTNQRKRLACAKRMYSLRFP
```



# Method



**Step 1: All sequences aligned using Muscle multiple sequence alignment**

**Step 2: Model selection for aligned data using IQ-tree**

**Step 3: Bootstrapping for 1000 trees in IQ-tree**

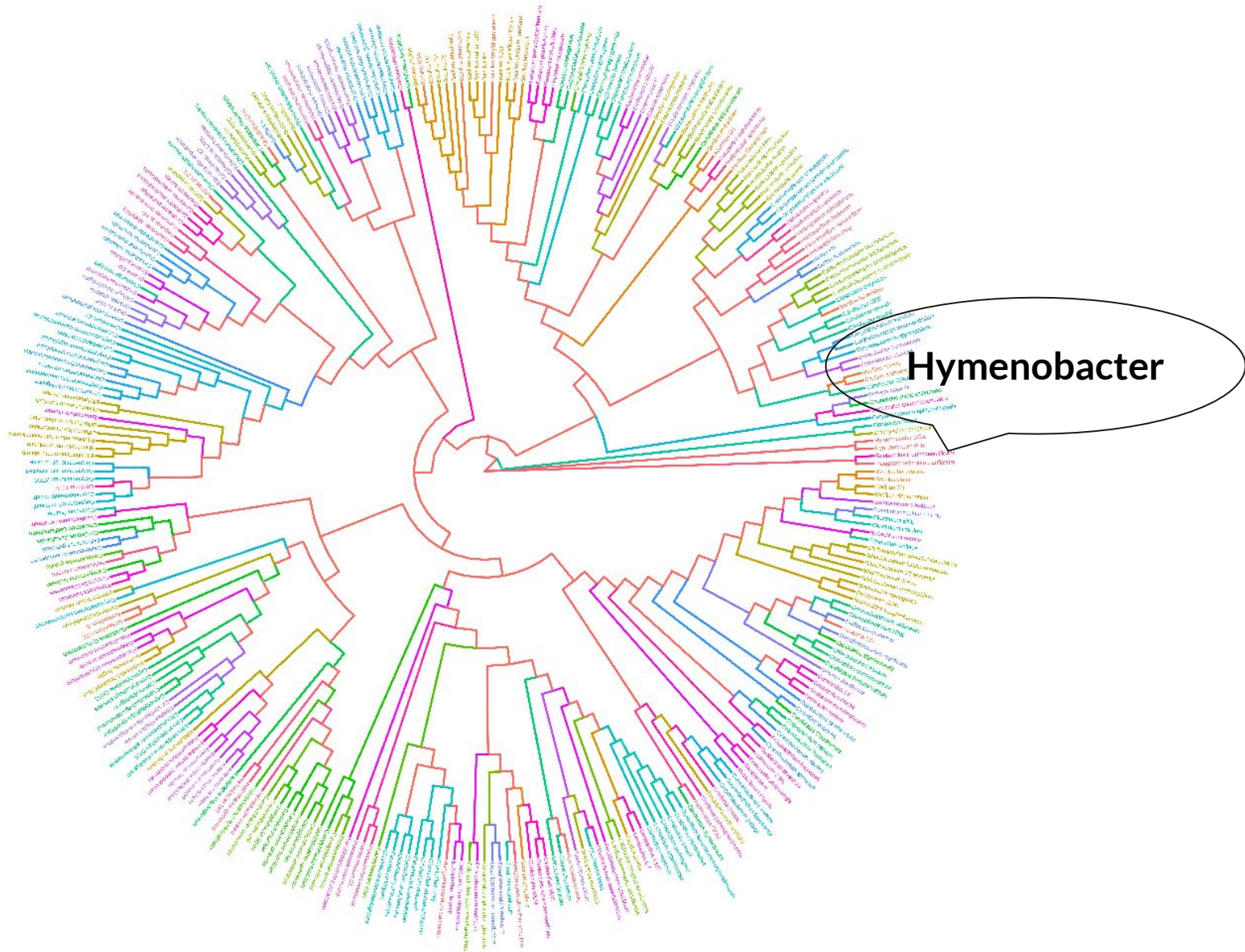
**Step 4: Root the tree and collapse clades with bootstrap support < 95%**

**Step 5: Analysing the results**

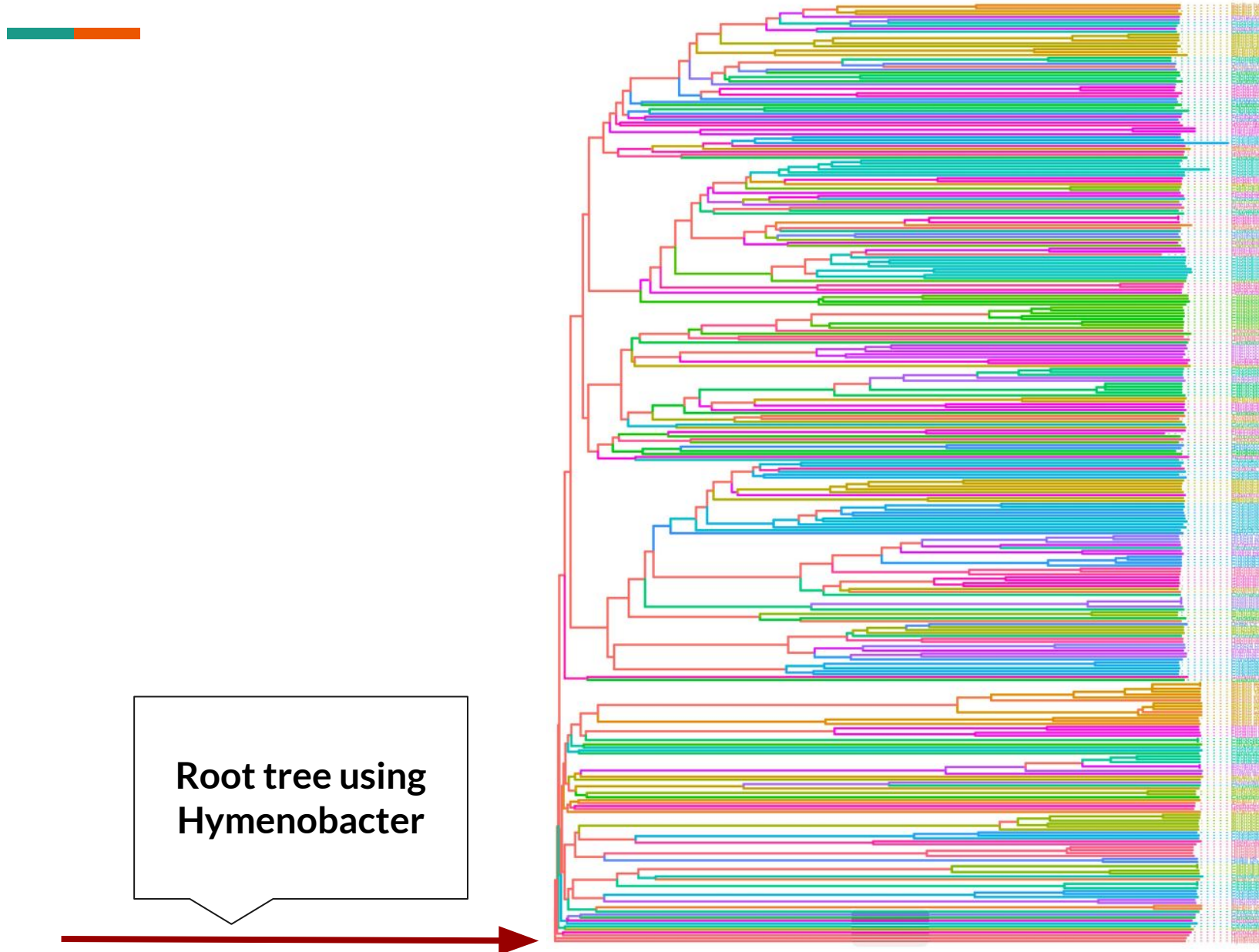
# Results-Alignment and GBLOCKS



# Results-NJ tree construction



# Results- Rooted NJ tree





# Results-Unrooted IQtree

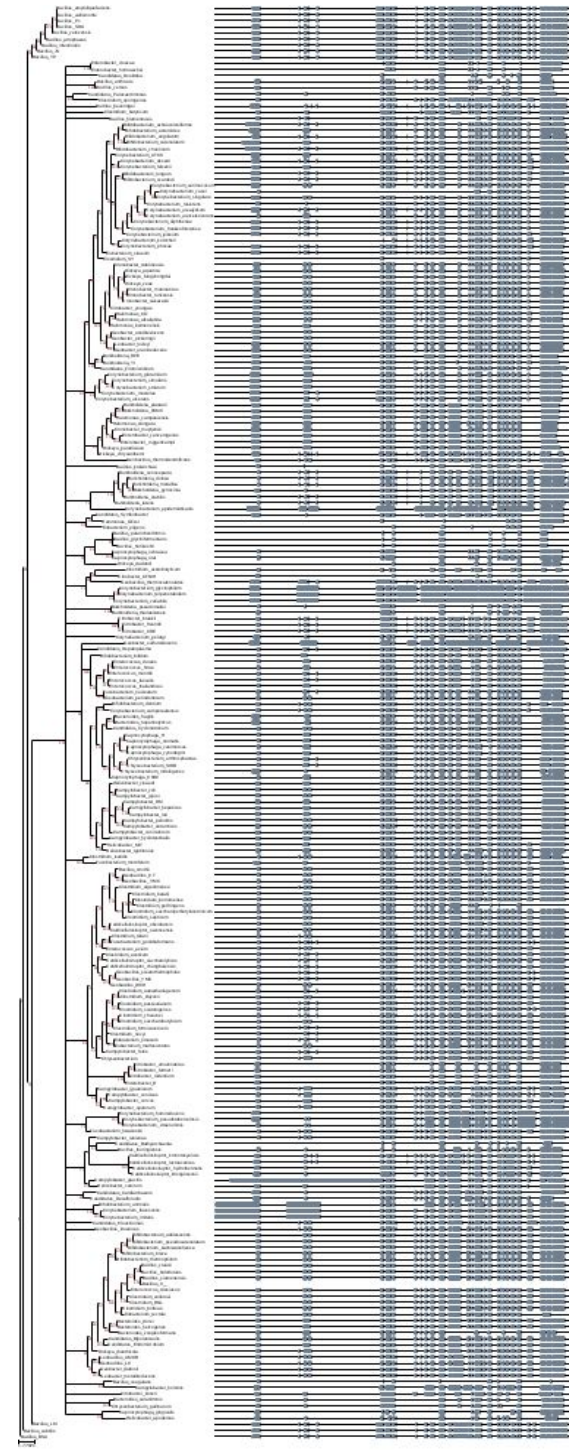
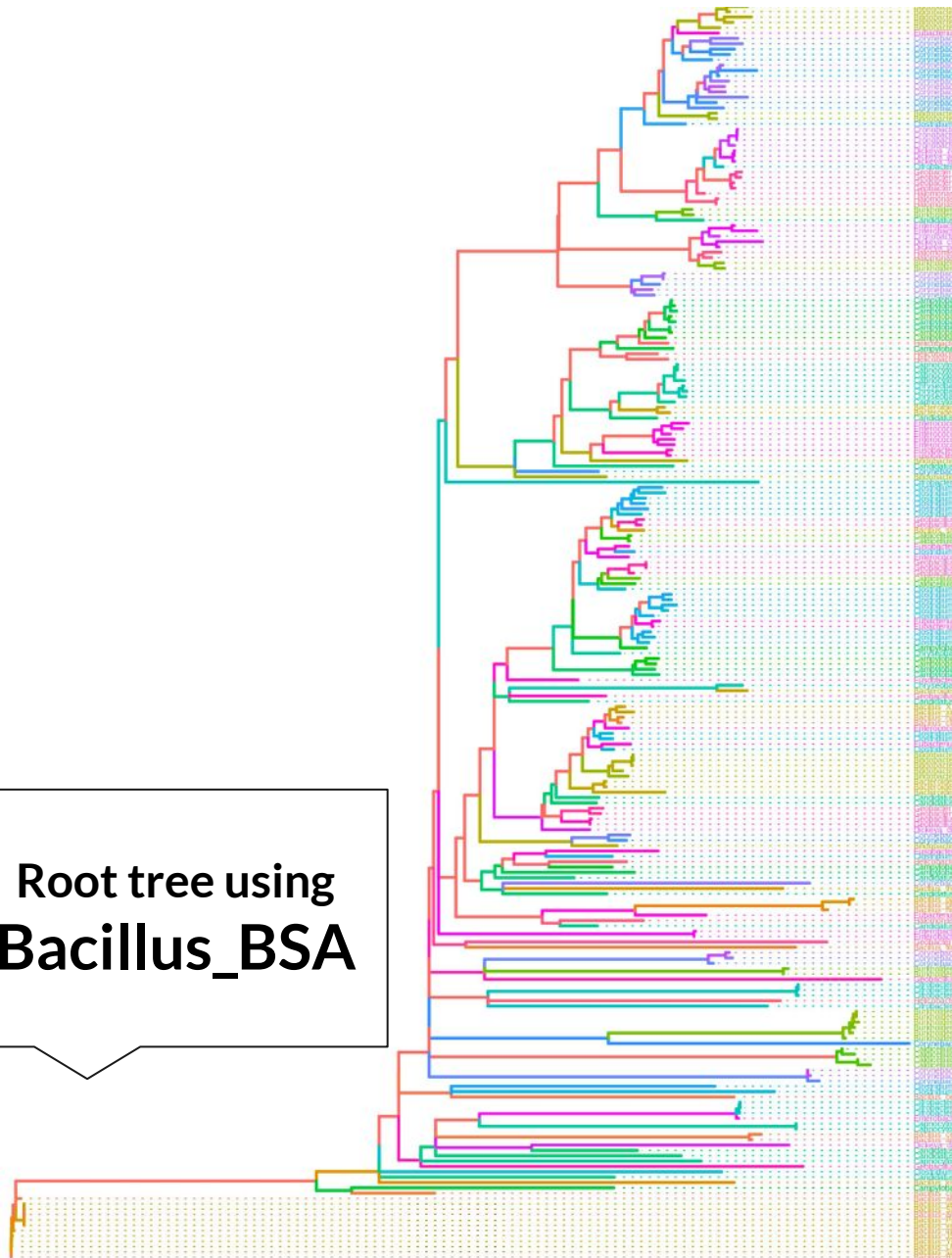
Bacillus\_BSA ← Outgroup

Best-fit model:  
VT+F+G4

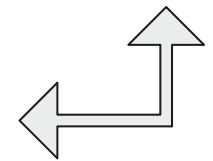
I got clear ml tree from  
Phylogenetic tree (newick) viewer

# Results-Rooted IQtree

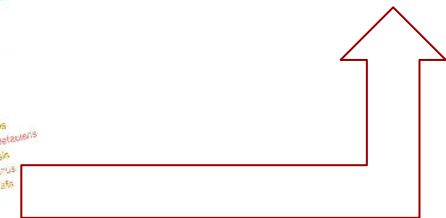
Root tree using  
*Bacillus\_BSA*



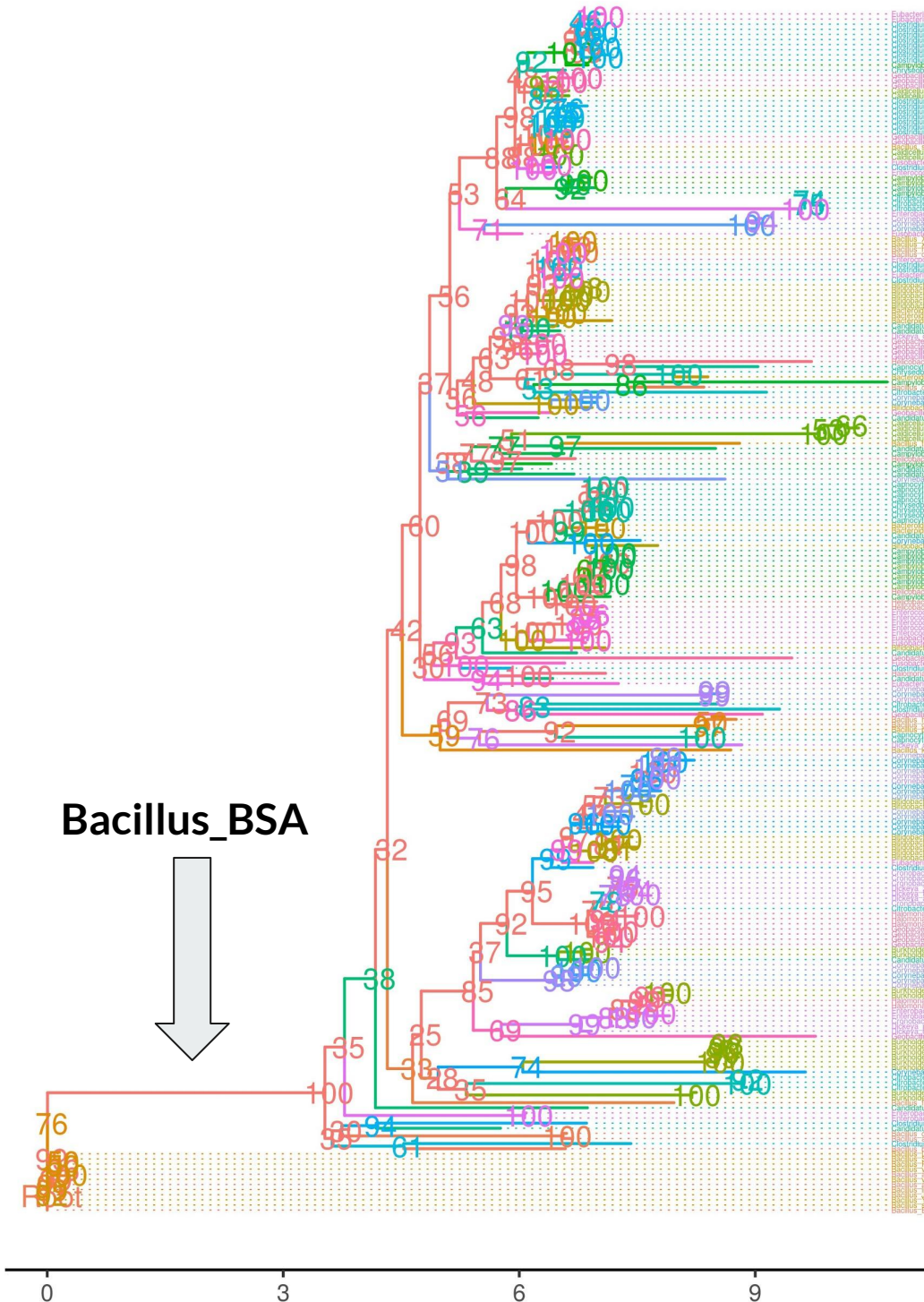
Rooted tree  
with alignment







# Results- IQtree with bootstrapping

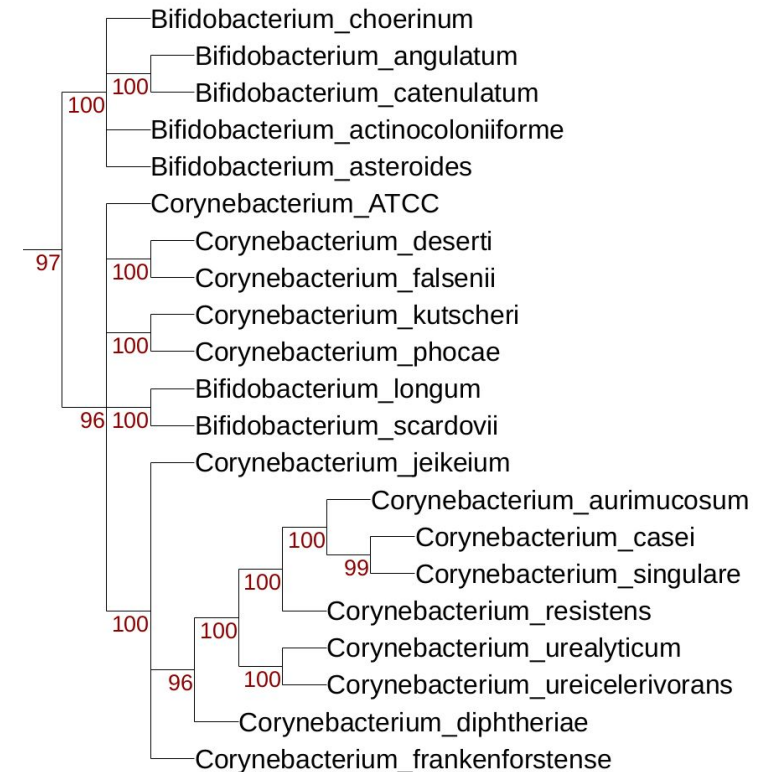
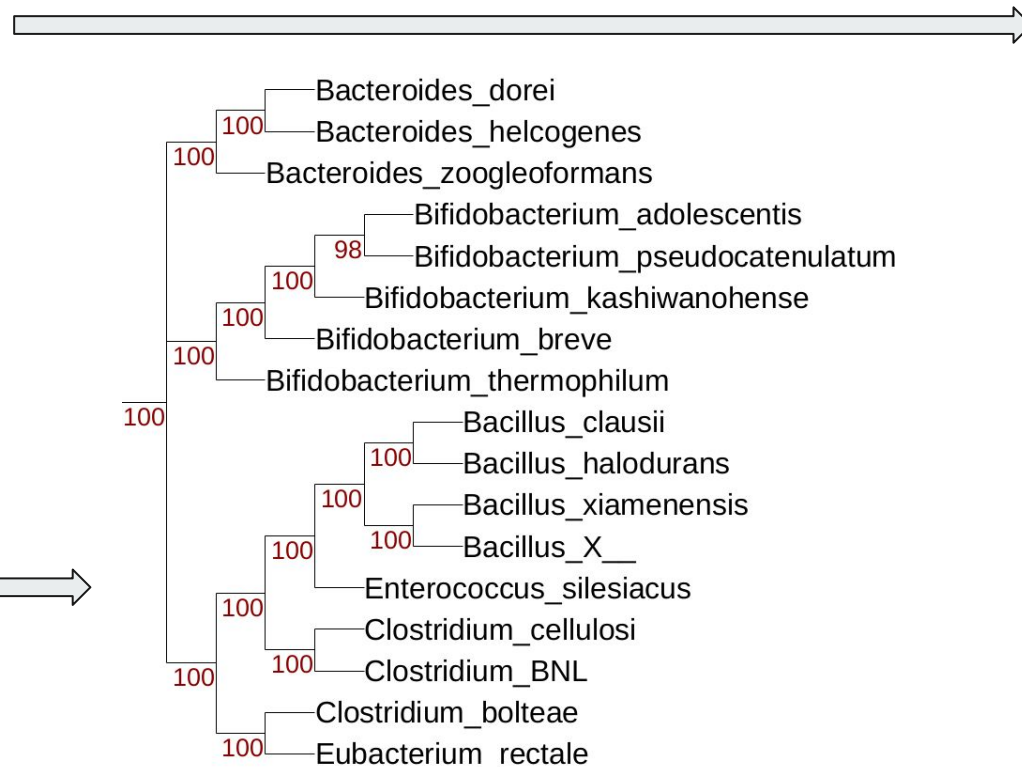
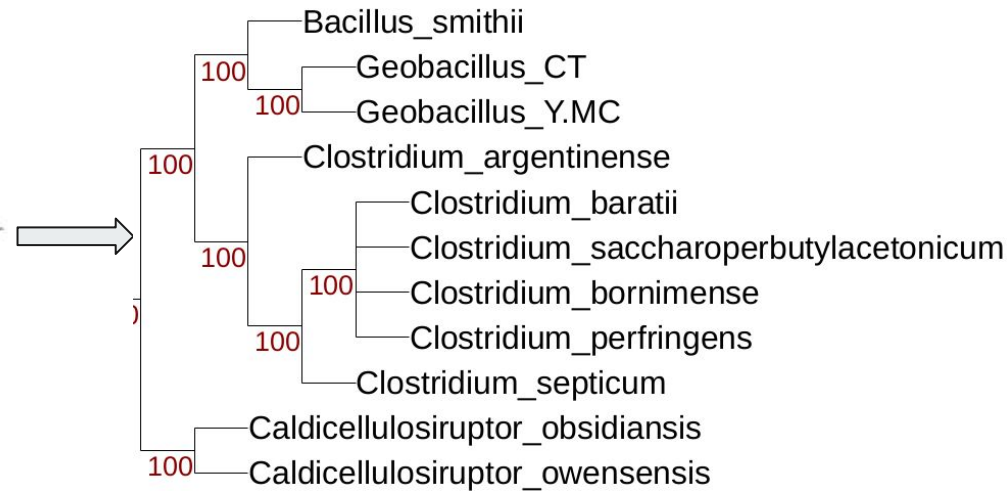


- Tree is huge and I have to compress it as much as possible!
- All clear photos are available on Github



## After collapsing (Final tree)

- “collapse” clades with bootstrap support < 95%



# Discussion



- 1.) After bootstrapping “Bacillus\_BSA” confirmed as an outgroup!
- 2.) It can be seen that cas 1 gene have been found scattered in most of the species from Bacilli (Bifidum group) and Clostridia (clostridium group) class!
  - a.) both are gram positive species where bacillus are aerobic and clostridium are anaerobi
  - b.) Bacillus mycoides is source organism of cas1
- 3.) most of the organisms evolved according to viral strains
- 4.) Also significant percentage of species from Actinobacteria class grouped with species from Bacilli and Clostridia class!

# Source



All codes will be available on GitHub:

[https://github.com/Sedreh/phylogenetics semester3](https://github.com/Sedreh/phylogenetics_semester3)

