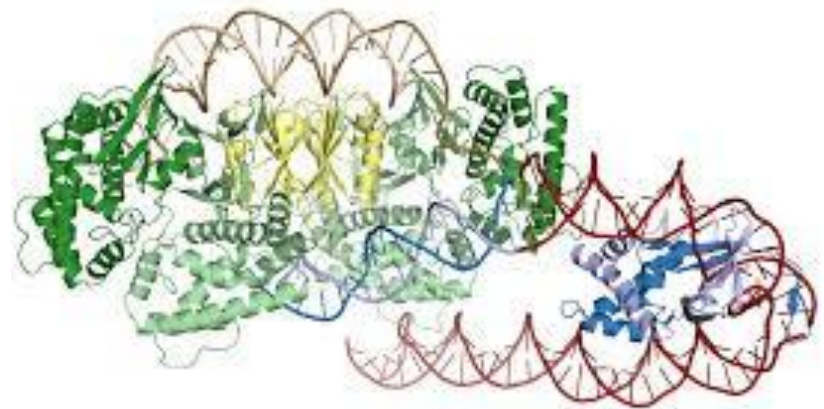# "CRISPR-associated protein 1 ( Cas1)"

## Molecular phylogenetics course project

Sedreh Nassirnia

Fall 2019

# Objective

- Study evolutionary relationships among bacteria family based on cas1

- What percentage of organisms evolved according to viral strains

# CRISPR-Cas Systems

- **CRISPR (clustered regularly interspaced short palindromic repeat)**
  - **adaptive immune system that provides protection against mobile genetic elements**
- **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**
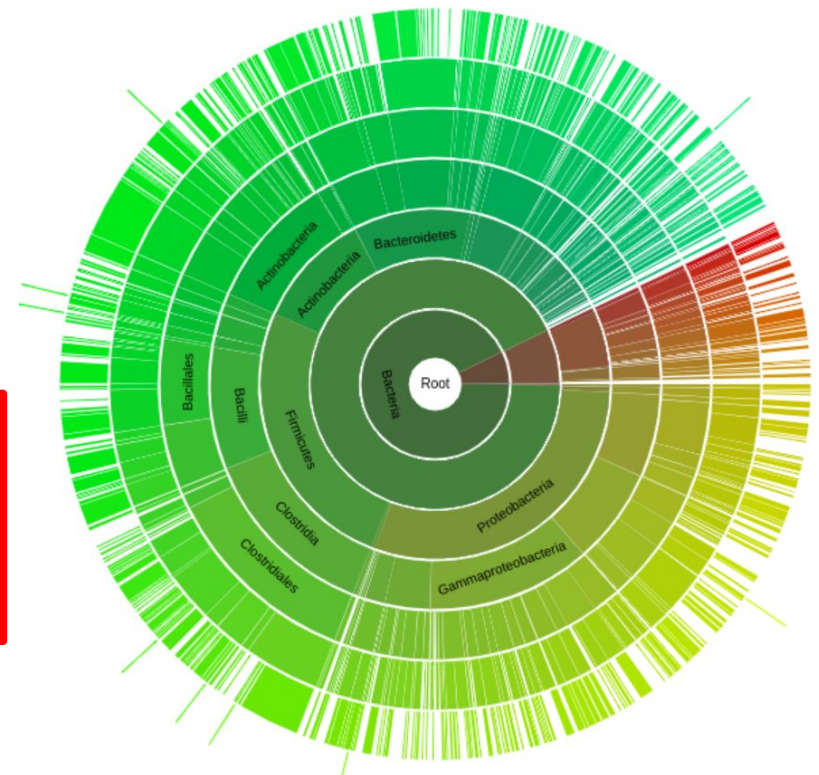
# Data



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

# Fetching IDs

```python
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/ITMO/semester3/Molecular_phylogenetic/COURSE_PROJECT/phyloproject_cas1/cas1_pfam.fasta'
        fasta.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.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/ITMO/semester3/Molecular_phylogenetic/COURSE_PROJECT/phyloproject_cas1/links.txt', 'r')
        for link in links:
            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:AREALGSMS7_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
MLYRSIYIGNPAYLKLKDQQMKIVCPETKAEKGSVPVEDLGLLMLDHFQITISHQLIQWLMGNNVVIISCDAHHLPHGQMLPLHGNAIYSQRIKDQIEASEPLKKQLWKQTIEC
>Corynebacterium striatum
MAYSEDAITFSTIPADHQVRLEDRVSFAYVEHAAIRQDRTGVVAYSVVDNSELEQRIQLPVGGLAVLMLGPGTSISAAAATSCTRSGTTIMFTGGGGVPAYTHAASLTSSARWA
>Candidatus Caldiarchaeum
MSELVIDKPGTYLGVRKGLFVVRTKGGGRSEFSPVELSHISIRCRGVGVSVDALRLACRFGIEVSVYSRGRPVGKVVGAFLGGGAVTRRAQLEAWGTERGLAVAREIVSAKLYN
>Corynebacterium singulare
MTTPHEVPLTRQALARVGDRISFLYAERCVINRDGNSLTIVDQRGTAHVPATQIAALLLGPGTKITYAAMALLGDAGVSAVWVGERGVRYYAHGRPPAKSSRMAEIQAEVVTHQ
>Clostridium tetani
MKRSYYIYNNGILKRKDNSMAFIDELGERRYIPIETANEIYVMSEMDFNTSLINYLSQYDVIIHFFNYYSFYTGSFQPRKKLVSGNLLVNQVNHYSDNSKRLEIAKKFVDGASY
>Chlorobaculum tepidum
MKKHLNTLFVTTQGSYLSKEGECVLISIDRVEKTRIPLHMLNGIVCFGQVSCSPFLLGHCAQLGVAVTFLTEHGRFLCQMQGPVKGNILLRRAQYRMADNYDQTATLARLFVIG
>Corynebacterium terpenotabidum
MNKIPFRSSVTTQGRNASGARSLWRATGMTDEDFEKPIIAVANSYTQFVPGHVHLKNVGDIVAEAVKEAGGVAREFNTIAVDDGIAMGHSGMLYSLPSREIISDSVEYMVNAHQ
>Hungateiclostridium thermocellum
MKKSAFIFSDGELKRKDSTVLFESEDSKNYLPIEDISDIYIFGEVTVTKKFLELATQKEILLHFYNYNEYYVGTYYPREHYNSGFMILKQAEHYLDEEKRMAIAKKFIHGSVKN
>Cronobacter turicensis
MSFVPLNPIPLNDRTSMIFLQYGHLDVLDGAFVLVDKTGVRTHVPVGAIACIMLEPGTRVSHAAIRLASQVGTLLVWVGEAGVRLYASGQPGGARADKLLYQAKLALDETLRLK
>Corynebacterium urealyticum
MRTPQQVPIERQSLSQMGDRISFLYVERAVVSRDGNALTVTDQRGVAHVPATQLAALLLGTGTRITNAAIALLGDSGVSTVWVGERGVRYYAHGRPPAKSSRLAELQARVVTNQ
>Corynebacterium ulcerans
MSYSNEALAFSTIPASEQIRLEDRVSFLYLEYCLIRQDRTGVIAVSRGDEKAPAELKDLPIKARIQLPVGGLAVLMLGPGTSISQPAATSCARAGVSVLFTGGGGVQAYSLSTP
>Corynebacterium ureicelerivorans
MGSRISFLYIERATVNRDGNALTITDQRSVAHVAATQLAVLLLGPGTRITYAAMALLGDAGVSIVWVGERGVRYYASGRPPAKSSRMAELQAEIVTNQRKRLACAKRMYSLRFP
```

# 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



Root tree using
Hymenobacter
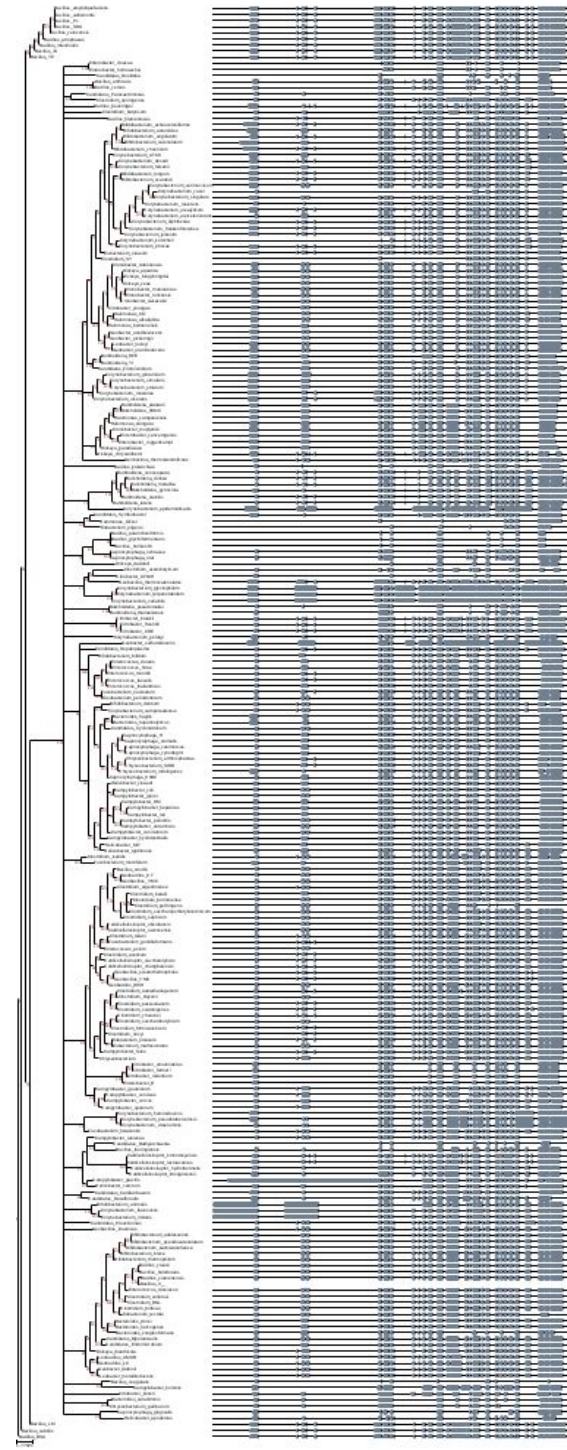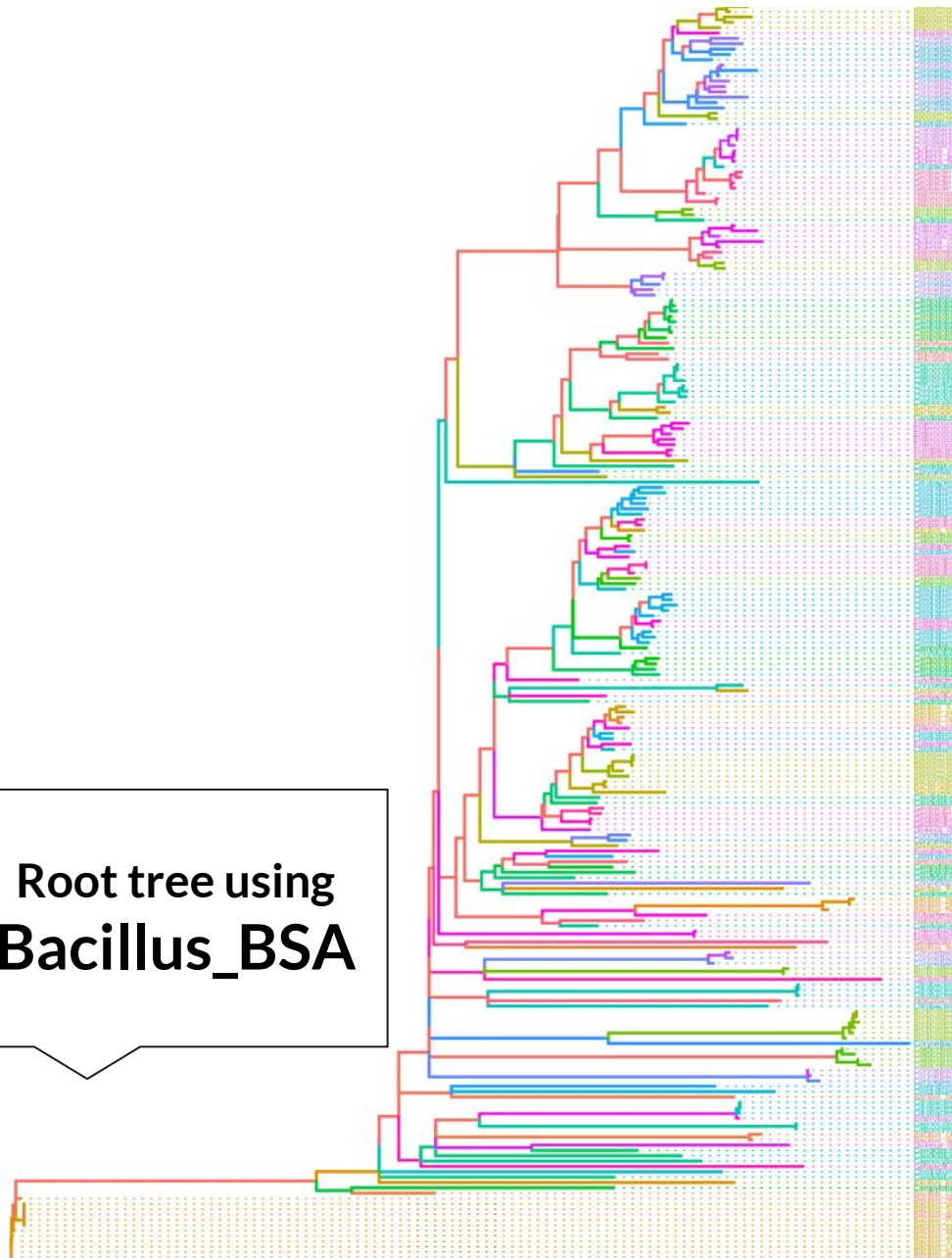
# Results-Unrooted IQtree



Bacillus_BSA ⟵ Outgroup
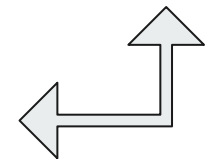
Best-fit model:
VT+F+G4

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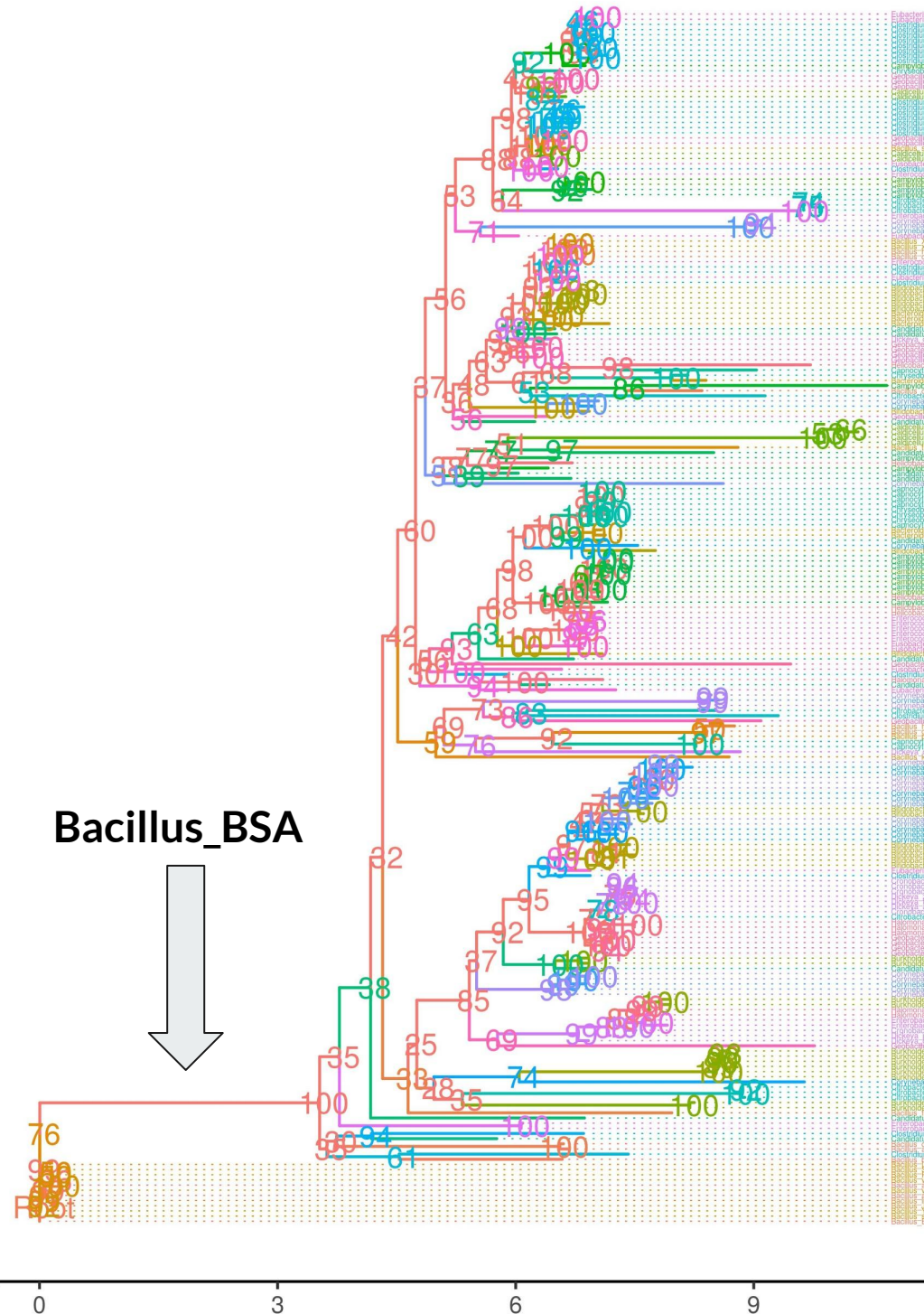Phylogenetic tree (newick) viewer

# Results-Rooted IQtree
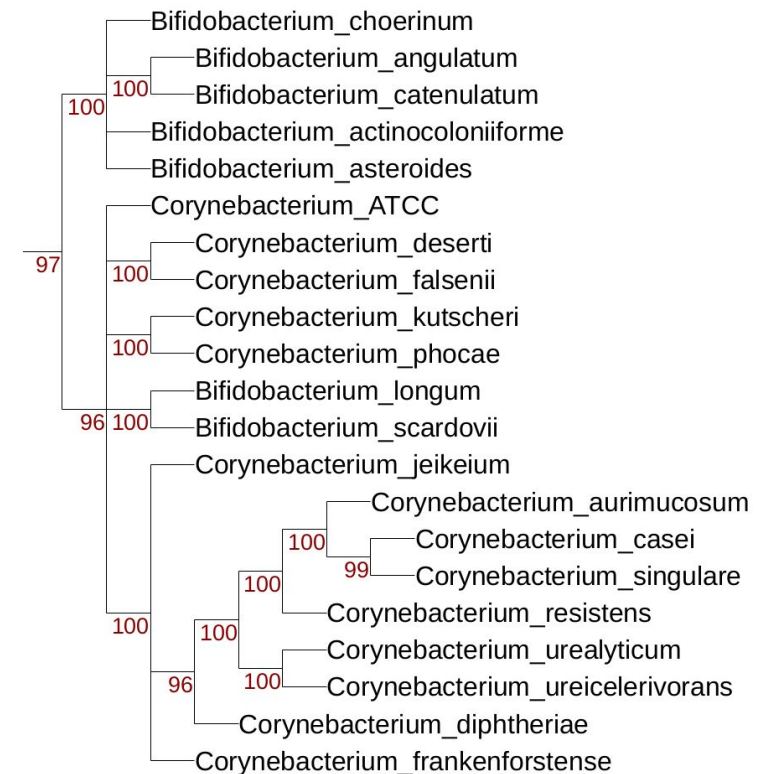


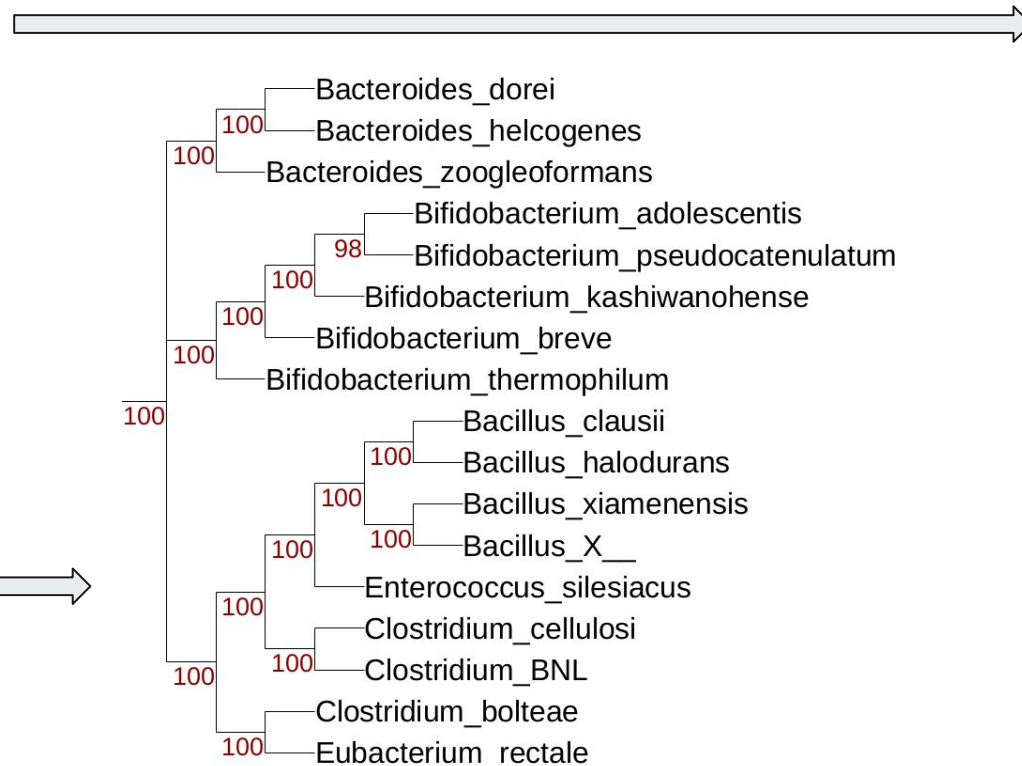Root tree using
**Bacillus_BSA**

Rooted tree
with alignment
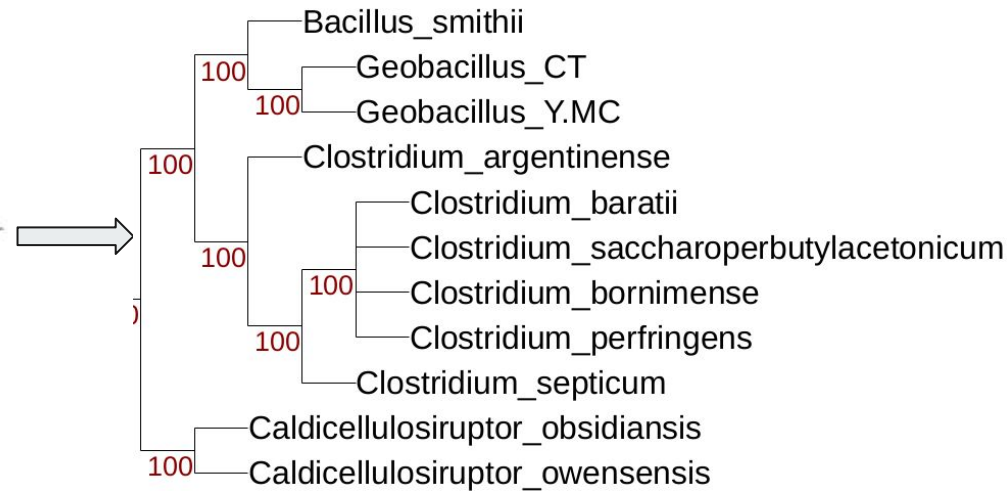
13

# Results- IQtree with bootstrapping



**Bacillus_BSA**

# Results- IQtree with bootstrapping



Bacillus_BSA

- **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%

16

# 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)