# CAFE:
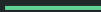## aCcelerated Alignment-FrEe sequence analysis

Name: Elahe Badali
Student No: 98209072

# INTRODUCTION

A general look

- Problem Definition

- Alignment-Free Methods

- CAFE, A Good Solution

# Problem Definition

- The dominant tools for sequence comparison are alignment-based methods, including global and local sequence alignments.

- But these tools have challenges in some cases:
  - Gene regulatory regions ⟶ not highly conserved
  - large amounts of short reads from NGS ⟶ assembling challenges
  - virus-host infectious associations

- Solution: Alignment-free sequence comparison

# Alignment-free methods

- Several types of alignment-free approaches based on:

  - counts of k mers

  - longest common subsequences

  - shortest absent patterns

- This paper is based on k-mer counts.

- Example: d2∗ and d2S

- Advantages:  perform well theoretically, solve 3 challenges of alignment-based methods.

- Disadvantages: relatively slow due to the requirement of calculating the expected k-mer counts.

# CAFE, A Good Solution

- speeds up the calculation of recently developed measures.
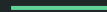- reduced memory requirement.

**Features:**

- Can work with:
  - assembled genomic sequences .
  - unassembled shotgun sequence reads from NGS technologies.
- Counts k-mers by JELLYFISH.
- The resulting pairwise dissimilarities: a symmetric matrix in PHYLIP format.
- four types of built-in downstream visualized analyses:
  - dendrograms using the UPGMA algorithm
  - heatmap visualization of the matrix
  - projecting the matrix to a 2D space using principal coordinate analysis (PCoA)
  - network display

# MATERIALS

CAFE in detail

- 3 Type Of Methods

- Measures

- Workflows

# 3 Type Of Methods

- Conventional k-mer counts: normalize k-mer counts into the k-mer frequencies.

- Adjusted k-mer counts: the Markov models for the sequences are assumed as the

  underlying generative models.

- presence/absence of k-mers: binarize k-mer counts into presence/absence indicators.

# Measures

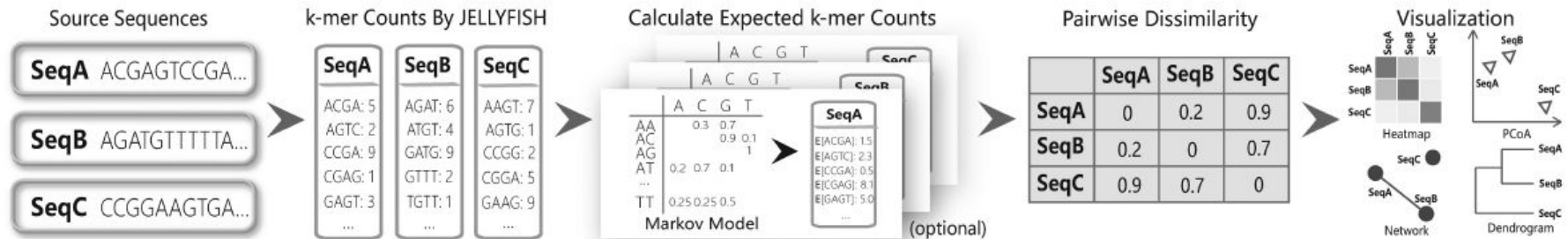| Conventional k-mer counts | Adjusted k-mer counts | presence/absence of k-mers, |
|---|---|---|
| 1. Canberra<br>2. Ch<br>3. Cosine<br>4. Co-phylog<br>5. D2<br>6. Eu,<br>7. FFP,<br>8. JS,<br>9. Ma<br>10. Pearson | 1. CVTree<br>2. d∗2<br>3. d2S | 1. Anderberg<br>2. Antidice<br>3. Dice<br>4. Gower<br>5. Hamman<br>6. Hamming<br>7. Jaccard<br>8. Kulczynski<br>9. Matching<br>10. Ochiai<br>11. Phi<br>12. Russel<br>13. Sneath<br>14. Tanimoto<br>15. Yule. |

# Workflows



**Figure 1.** The workflow of CAFE. The JELLYFISH software parses the input sequence files (in Fasta format), counts *k*-mers and saves compressed information into separate databases. CAFE subsequently loads the databases and constructs a symmetric dissimilarity matrix among the inputs. CAFE also integrates four types of visualized downstream analysis, including dendrograms, heatmap, principal coordinate analysis (PCoA) and network display.
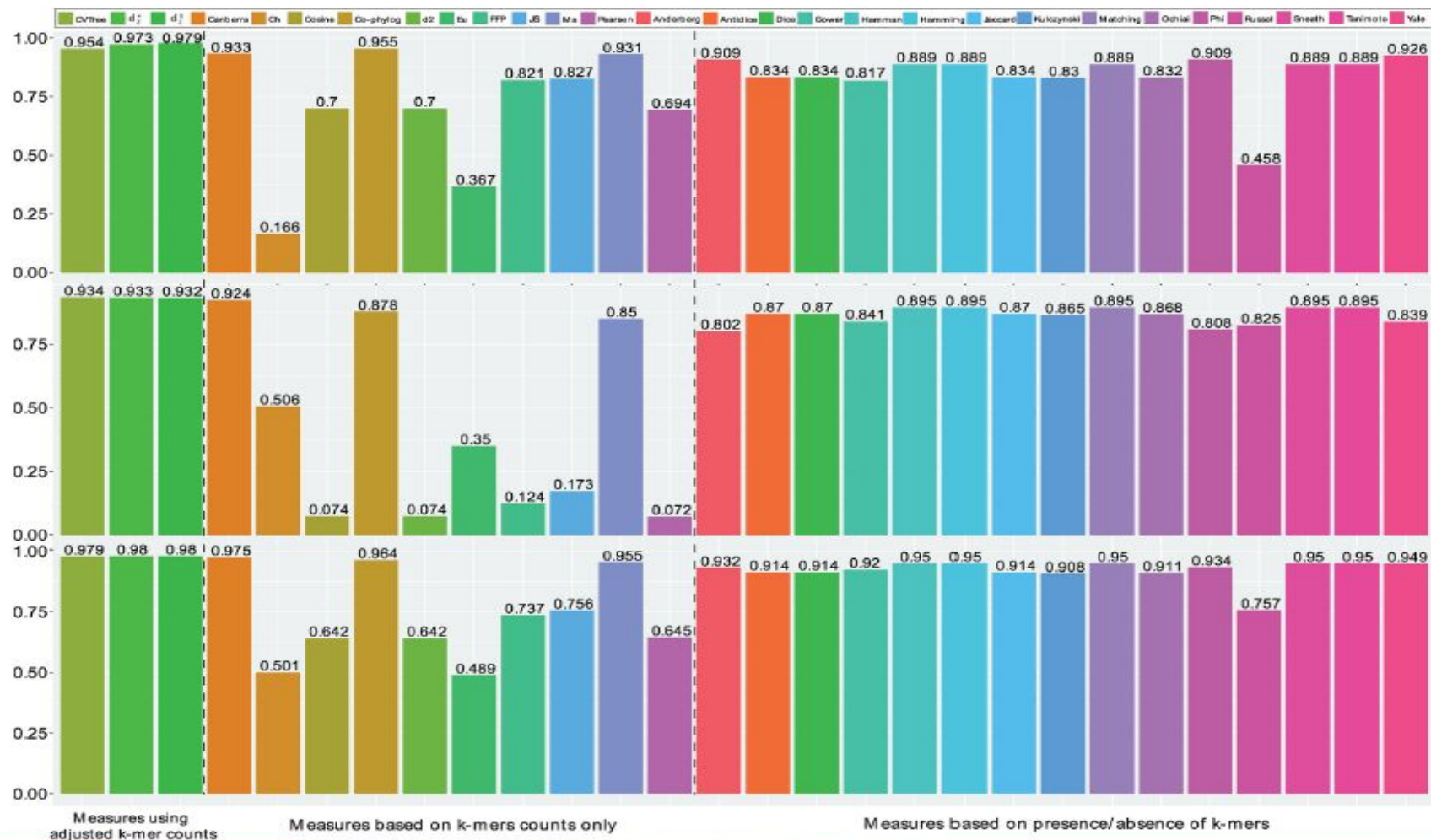
# RESULTS

Evaluation in 3 Cases

- Application to primate and vertebrate genomic sequences
- Application to microbial genomic sequences
- Application to metagenomic samples

# Application to microbial genomic sequences

1. **Dataset 1:** 21 primates species
2. **Dataset 2:** 28 vertebrate species
3. **Dataset 3:** Combined  the two datasets

How to Evaluate:

The Spearman correlation of various dissimilarity measures with the evolutionary distances using maximum likelihood approach across many genomic regions based on above datasets.

| Sequence Model | Original Implementation | | CAFE | | |
|---|---|---|---|---|---|
| | Wall time | Peak memory | Wall time | Speedup | Peak memory |
| order=0 | 0:42'32" | 64.0G | 0:6'09" | 6.9x | 31.1G |
| order=1 | 1:44'18" | 64.0G | 0:6'13" | 16.8x | 31.1G |
| order=2 | 2:11'32" | 64.0G | 0:6'12" | 21.2x | 31.1G |
| order=3 | 2:34'28" | 62.4G | 0:5'05" | 30.4x | 24.8G |
| order=4 | 2:34'11" | 62.3G | 0:6'10" | 25.0x | 31.1G |
| order=5 | 3:24'43" | 64.0G | 0:5'08" | 39.9x | 24.8G |
| order=6 | 2:53'08" | 63.9G | 0:5'14" | 33.1x | 24.8G |
| order=7 | 2:40'04" | 64.0G | 0:6'29" | 24.7x | 31.1G |
| order=8 | 2:33'19" | 64.0G | 0:6'08" | 25.0x | 48.1G |
| order=9 | 2:37'50" | 64.2G | 0:6'19" | 25.0x | 48.2G |
| order=10 | 2:22'18" | 64.7G | 0:5'15" | 27.1x | 48.5G |
| order=11 | 2:05'55" | 60.4G | 0:6'29" | 19.4x | 49.6G |
| order=12 | 1:53'40" | 74.6G | 0:6'39" | 17.1x | 37.0G |

**Figure 4.** Wall time, peak memory usage and speedup ratio comparison between CAFE and the original implementation to calculate $d_2^*$ dissimilarity between a pair of genomes for $k = 14$.

# Application to microbial genomic sequences

- **Dataset:** 27 E. coli and Shigella genomes dataset

- 6 E. coli reference (ECOR) **groups:** A, B1, B2,D, E and S.

- Used **UPGMA** to cluster the samples based on the calculated pairwise dissimilarity matrix

- **The Markov order 1**

- k = 14

- Result here based on CVTree, d2* and d2S
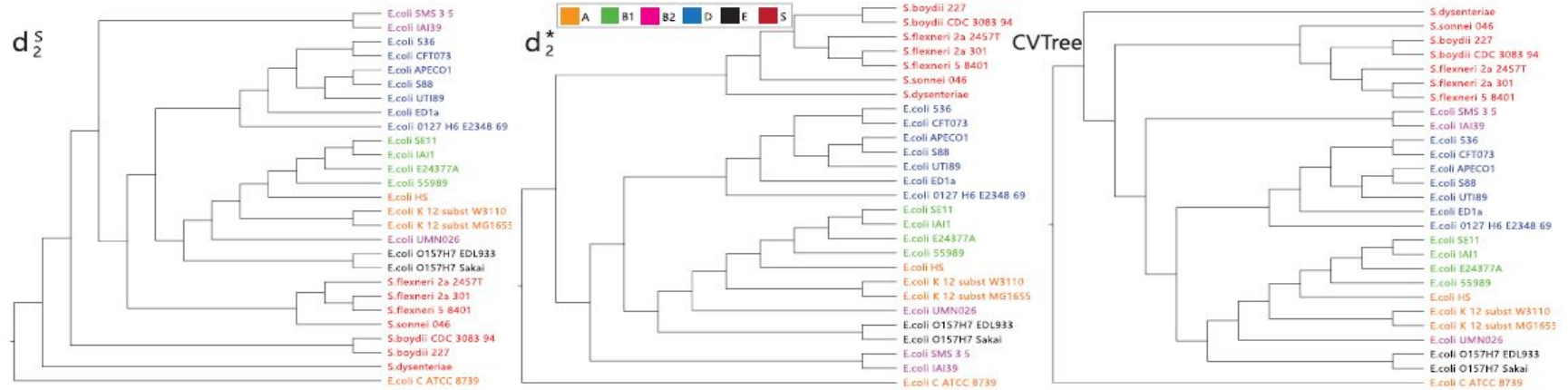
# Result



**Figure 5.** The clustering results of 27 *Escherichia coli* and *Shigella* genomes using measures based on background adjusted 14-mer counts: $d_2^S$, $d_2^*$ and *CVTree*. The Markov order of the sequences were set at 1. The colors indicate the six different *E. Coli* reference groups.

# Application to metagenomic samples

- Dataset: mammalian gut metagenomic dataset comprised of NGS short reads from 28 samples.

- 3 groups: 8 hindgutfermenting, herbivores, 13 foregut-fermenting herbivores 7 simple-gut carnivores.

- UPGMA to cluster

- TheMarkov order 0

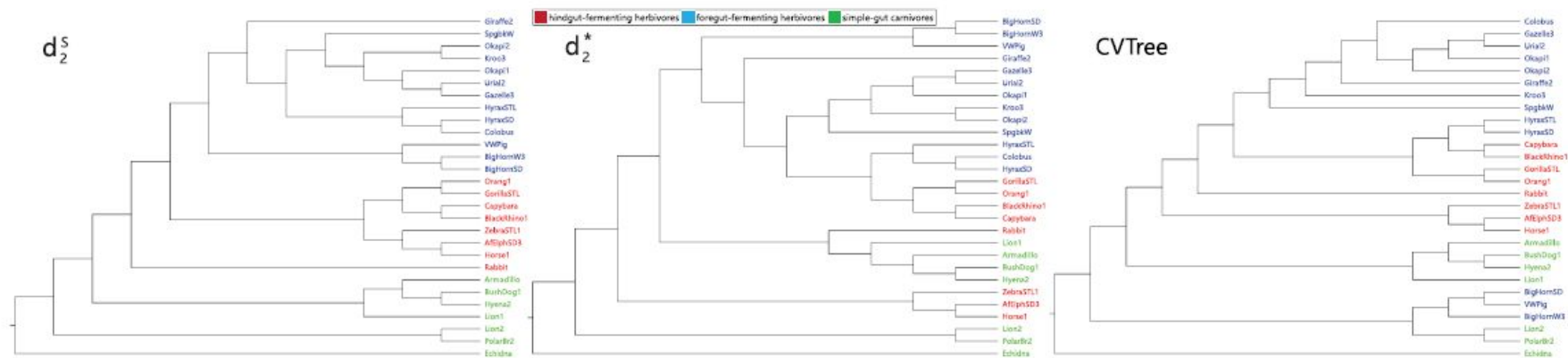- k = 5

- Result here based on CVTree, d2* and d2S

# Result



**Figure 6.** The clustering results of the mammalian gut samples using measures based on background adjusted $k$-mer counts: $d_2^S$, $d_2^*$ and *CVTree*.

# Other Articles

- This article Cited by 14 articles

- For Example:

  - **Alignment-Free Sequence Analysis and Applications**

    - an updated review

  - **Reads Binning Improves Alignment-Free Metagenome Comparison**

    - imperfection in d2S and d2∗ ⟶ neglect the heterogeneity

    - different reads bins

  - **Afann: Bias Adjustment for Alignment-Free Sequence Comparison Based on Sequencing Data Using Neural Network Regression**

    - Alignment-free methods ⟶ more time and memory efficient

    - However, dissimilarity can be overestimated

    - Afann ⟶ adjusts this bias

# Other Works of Author

- **Alignment-free oligonucleotide frequency dissimilarity measure improves prediction of hosts from metagenomically-derived viral sequences**
  - Viruses and their host genomes often share similar oligonucleotide frequency (ONF) pattern
  - k=6 and $d_2^*$
- **Improving contig binning of metagenomic data using *dS*2 oligonucleotide frequency dissimilarity**
  - model contigs using relative *k*-tuple composition, followed by measuring dissimilarity between contigs using d2s
- **.VirFinder: a novel *k*-mer based tool for identifying viral sequences from assembled metagenomic data**
  - the first *k*-mer frequency based, machine learning method for virus contig identification that entirely avoids gene-based similarity searches.

# Thank You For Your Attention