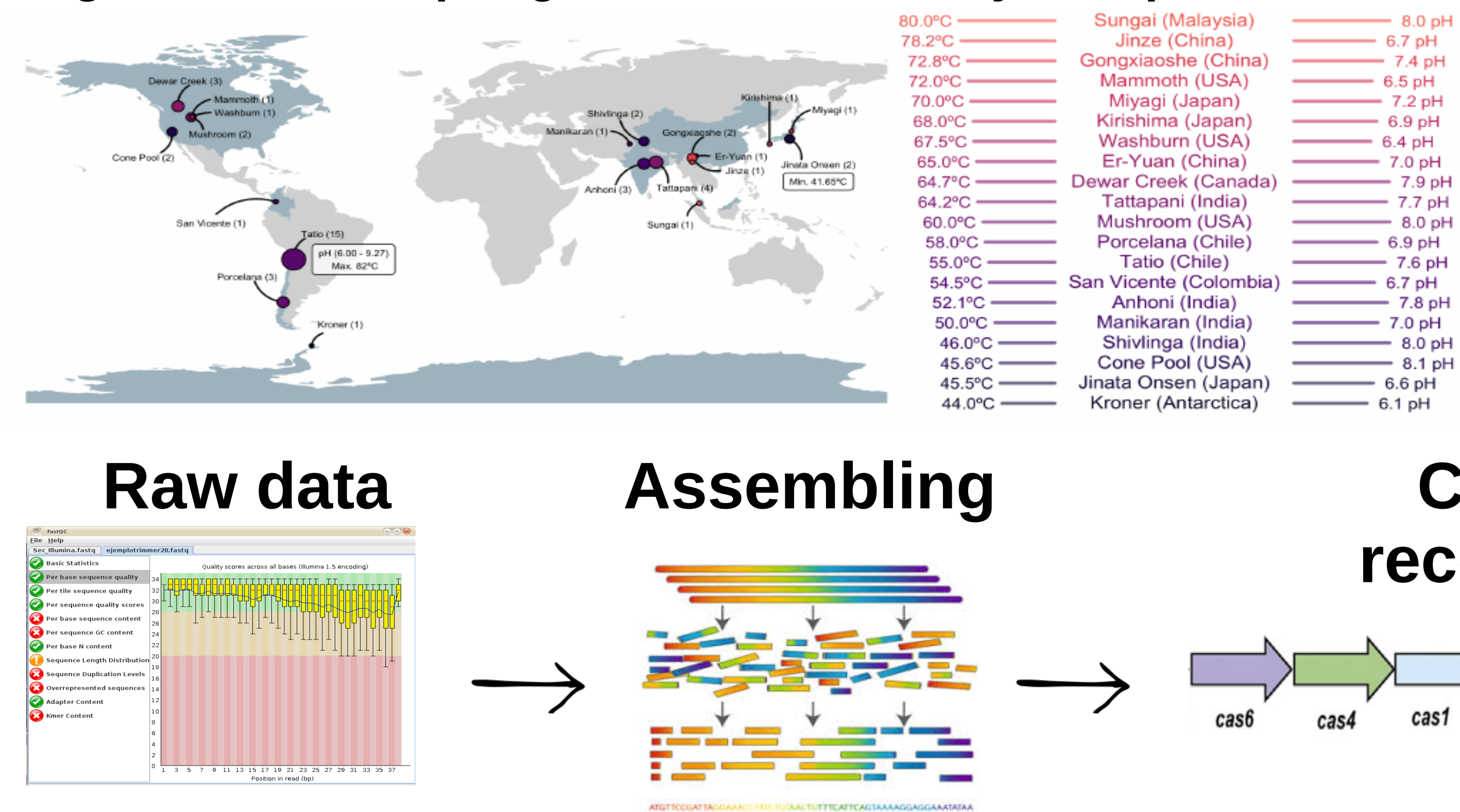


INTRODUCTION

Microorganisms threaten viral infections, which have caused the maintenance of molecular immunity mechanisms. The CRISPR-Cas systems (Clustered Regularly Interspaced Short Palindromic Repeats and CRISPR-associated proteins, Cas) perform the recognition and degradation of invading foreign nucleic acids as viruses. There are different types of CRISPR-Cas systems, and each one has a set of different Cas proteins. These proteins are highly diverse and mobile. Nevertheless, the Cas1 protein is recognized as the most conserved Cas and allows for the description of the represented groups of the CRISPR-Cas systems. On the other hand, some Cas1 sequences do not belong to CRISPR-Cas systems; instead, they are part of special transposons, called casposons due to Cas1 presence, an enzyme called casposase. CRISPR-Cas systems are more prevalent in microorganisms living in high-temperature environments; hence these environments are a perfect model for studying the diversity of CRISPR-Cas systems through Cas1 analyses due to the low complexity of hot springs communities compared with mesophilic environments. We used a phylogenetic and similarity approach to describe the ecological diversity of Cas1 proteins from hot springs. Our results show the novelty of Cas1 sequences from hot springs of the world.

MATERIALS AND METHODS

Figure 1. 48 hot springs used in this study with pH and



Cas1-solo analyses

IQ-TREE
Efficient software for phylogenomic inference

CRISPR-Cas classification

Russel88/
CRISPRCasTyper

Phylogeny & Network

IQ-TREE
Efficient software for phylogenomic inference

Cytoscape

Taxonomic assignment

GTDB
TK

NCBI

MAGs recovery

metaWRAP

Healing of sequences

NCBI

RESULTS AND DISCUSSION

We obtained 2150 Cas1 sequences over 0.1% of abundance from the 48 metagenomes analyzed.

Taxonomic distribution is according to predominant *phyla* described in hot springs (Figure 2A), where Class I CRISPR-Cas systems were observed (Figure 2B).

Figure 2. Taxonomic distribution of Cas1 sequences (A) and CRISPR-Cas classification (B).

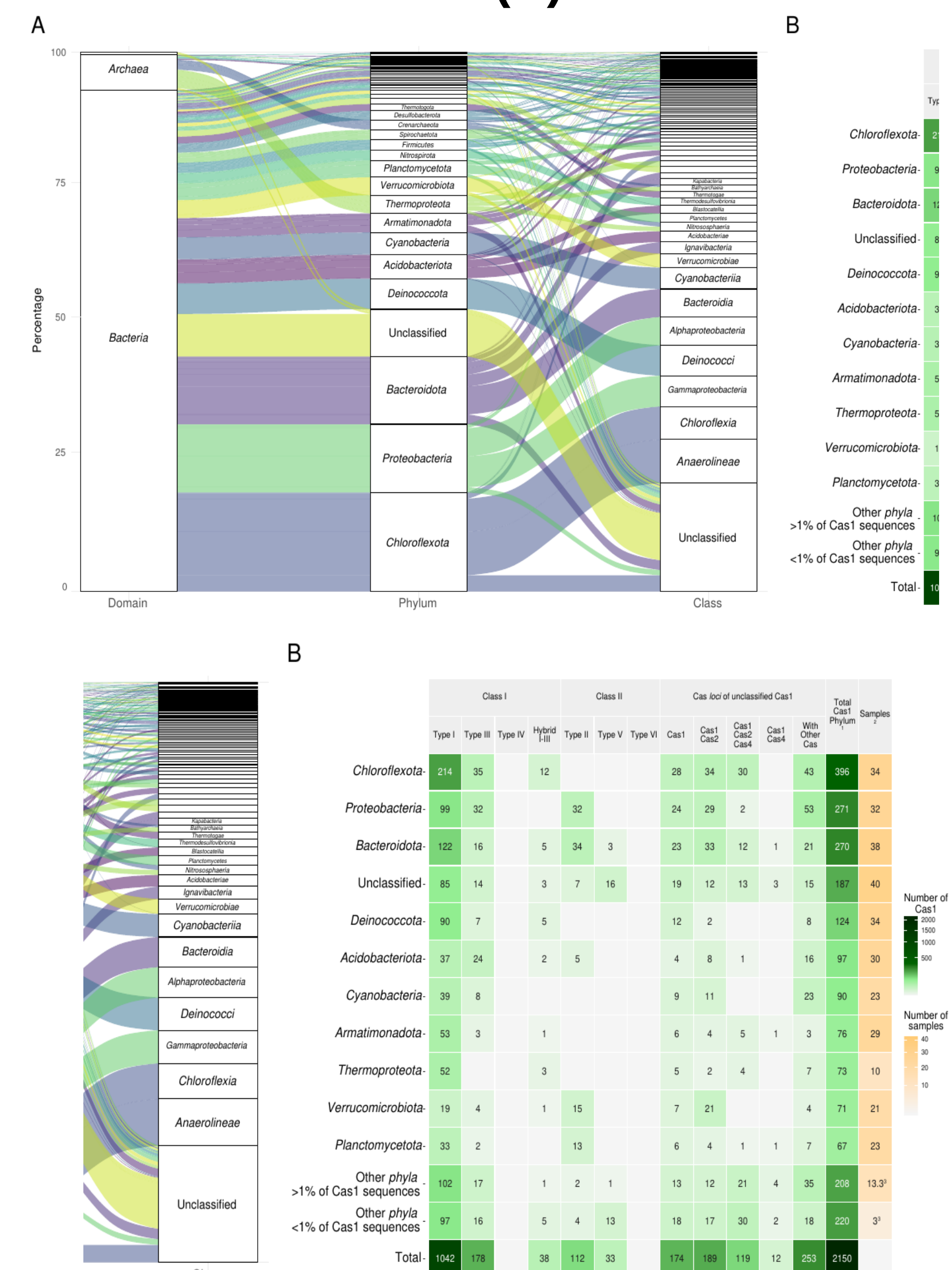


Figure 3. nMDS of 16S rRNA gene (A) and Cas1 gene (B).

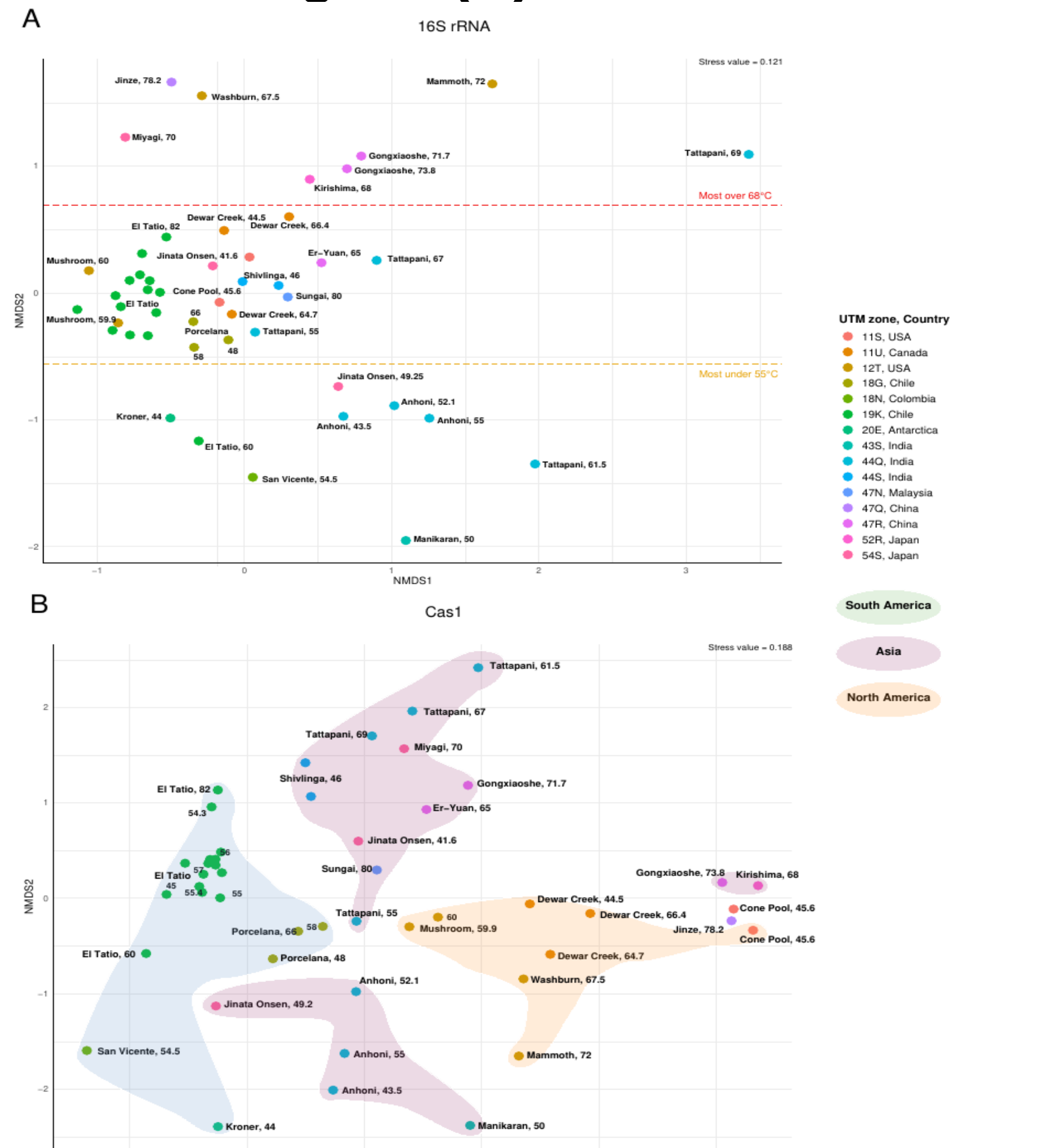


Figure 4. Phylogenetic ML tree of 2150 Cas1 from hot spring with 93 reference Cas1.

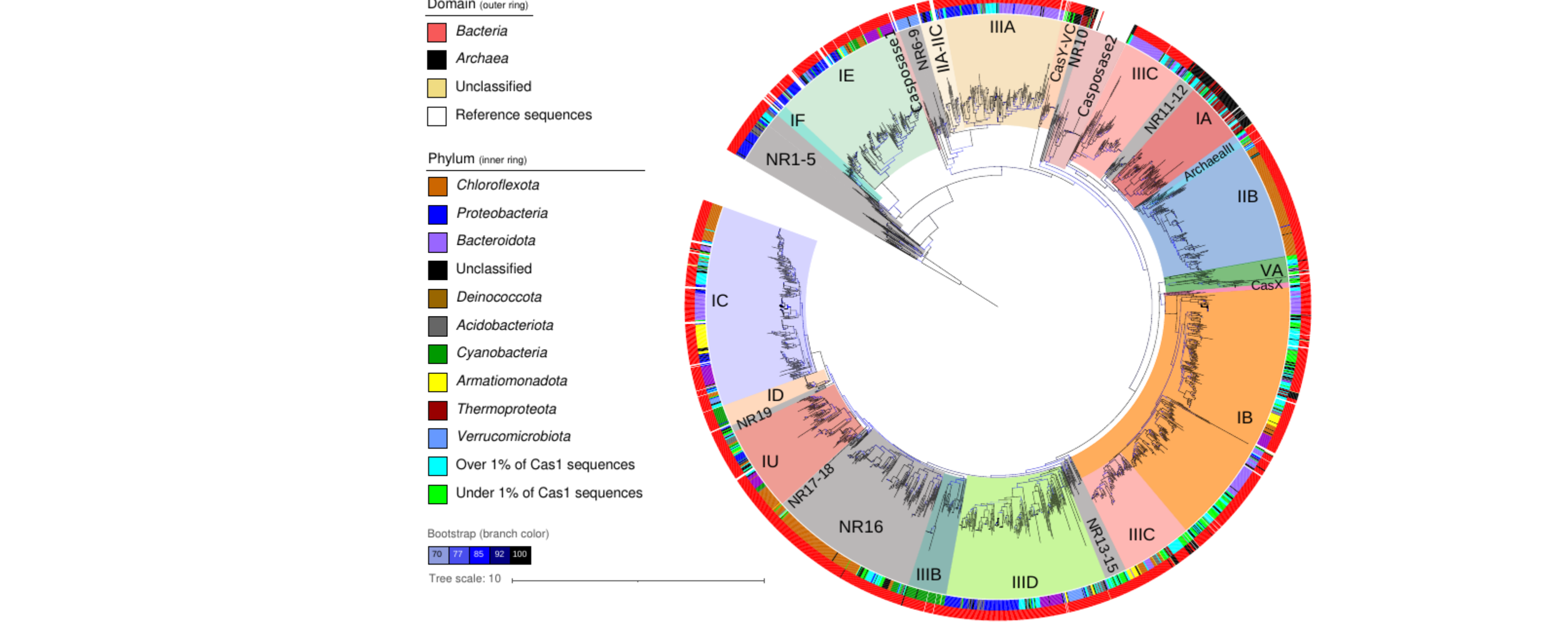
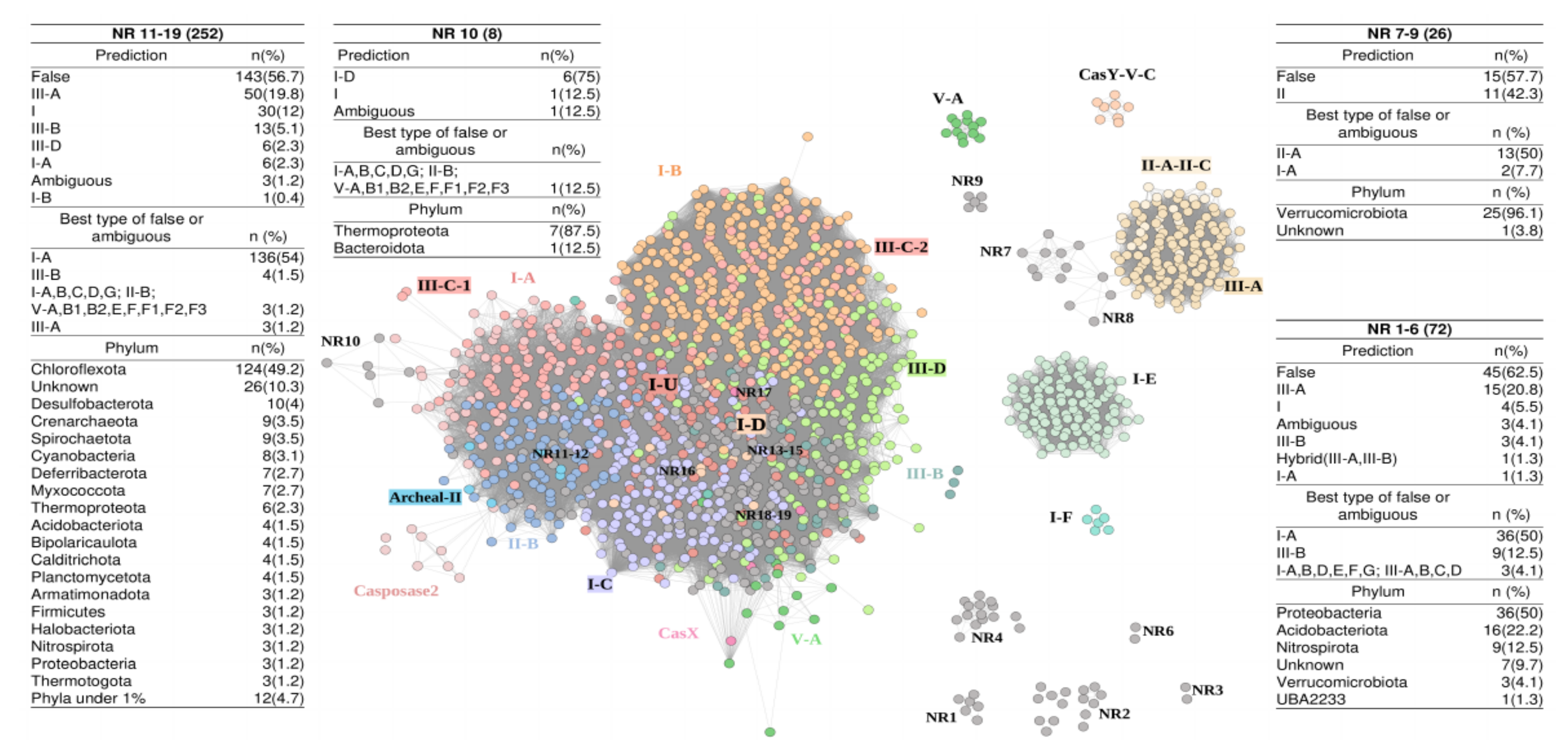


Figure 5. Network analyses of Cas1 from hot spring with 93 reference Cas1.



The phylogenetic tree reveals that Cas1 tends to group according to CRISPR-Cas systems type. However, several clades do not group with reference Cas1 sequences from databases (NR, Figure 4). These clades are composed of sequences of several taxa and tend to be at the root of significant clades, which is also observed for Cas1 of recently discovered CRISPR-Cas systems or casposons Cas1.

The similarity of Cas1 through network analyses shows that most Class I Cas1 shapes a big module separated from Class II sequences and NR sequences (Figure 5). This suggests that Cas1 proteins from hot springs are new for databases and could represent new CRISPR-Cas systems.

In-depth phylogenetic and syntenic analyses of Cas1 that group close to casposase (Cas1-solo) indicate that hot spring harbor new casposons, which are located in a branch distant from described casposase families and vestigial Cas1-solo (Figure 6).

Figure 6. Phylogenetic ML tree of Cas1-solo from hot springs.



CONCLUSIONS

- Cas1 genes from hot springs are represented in predominant *phyla* of these environments.
- Cas1 from hot springs tends to be specific to each geographical location.
- Novel phylogenetic clades are observed using Cas1 from hot springs and reference sequences, which suggest new CRISPR-Cas systems.
- Casposases from hot springs reveal a new family of casposons.