## 1. Introduction

In this study, Prosthecochloris aestuarii and Rubrobacter xylanophilus were subjected to comparative genomic analysis. Orthologs were found between protein sequences from two different species using the Best Bidirectional Hits (BBH) method, and a phylogenetic tree was created by aligning the homologs of one example of orthologs. Speciation and Duplication events were marked on the tree, which was compared with a species tree based on rRNA resources. Multiple Sequences Alignment (MSA) was applied in the end to identify conserved regions in protein sequences and potential promotor regions in DNA sequences.

## 2. Methods

### 2.1. Compute orthologs using Best Bidirectional Hits

All protein sequences of Prosthecochloris aestuarii (RefSeq: GCF_000020625.1) and Rubrobacter xylanophilus (RefSeq: GCF_000014185.1) were downloaded from NCBI database (https://www.ncbi.nlm.nih.gov/data-hub/genome/) as FASTA files. Two FASTA files were used to build a database of Prosthecochloris aestuarii and a database of Rubrobacter xylanophilus, then protein sequences from one species were used to run protein blast with the database of another species' protein sequences and the database of its own protein sequences. Four output files were generated from protein blast (output_Pa.txt, output_Pa_self.txt, output_Rx.txt, ouput_Rx_self.txt), and a python script (BBH.py) was created to generate BBH for two species using these output files, and it also indicated all in-paralogs.

### 2.2. Check orthology with a phylogenetic tree

One protein sequence of orthologs from co-orthology was picked to run protein blast in NCBI with standard databse and retrieve homologs to that protein sequence, then these homologs and the BBH were used to do MSA and generate the phylogenetic tree using clastalw2. When building the tree in clastalw2, the function of excluding positions with gaps and correcting for multiple substitutions was turned on, and the bootstrap value was set as 1000. The tree file was visualized in iTol (http://itol.embl.de) including branch lengths and bootstrap values, and speciation events were marked as red squares and duplication events were marked as green squares. In addition, rRNA sequences of all homologs species were downloaded from silva rRNA database (https://www.arb-silva.de/), and a species tree was built using these rRNA sequences. The species tree was also visualized in iTol to discuss the differences from the phylogenetic tree of homologs.

### 2.3. Identify functional regions

The resulting file of MSA was used to run a python script (Conservation.py), which uses an algorithm to calculate the variability to measure conservation and results in a list of variabilities to every position in the sequence (Variability.txt). The list of variability was exploited to generate a variability plot in R, and the value of variability is equal to 1 means all sequences share the same amino acid in that position.

$$variability = \frac{N * k}{n}$$

N: The number of sequences in the alignment
k: The number of different amino acids at a given position

n: The frequency of the most common amino acid at that position

## 2.4. Identify promotor regions

In order to identify promotor regions, the DNA sequences to specific proteins (orthologs) were downloaded from the NCBI and moved to one FASTA file. Clastalw2 was used to do the MSA analysis of these DNA sequences, and PVS (http://imed.med.ucm.es/PVS/) was used to see the variability plot.

# 3. Results

## 3.1. Orthologs

In two FASTA files (R.x_protein.faa and P.a_protein.faa) for species, 3205 protein-coding genes in the genome of Rubrobacter xylanophilus, and 2302 protein-coding genes in the genome of Rubrobacter xylanophilu. After running protein blast and BBH.py, 804 pairs of BBH were found between the two species, and 184 pairs of in-paralogs were identified. Because both in-paralogs are orthologous to the BBH, a total of 1792 orthologs were found between the two species.

## 3.2. Phylogenetic trees

WP_012504624.1 was picked as the example to run protein blast, this is a S41 family peptidase of Prosthecochloris aestuarii, its in-paralog is WP_012505586.1 and the BBH of it is WP_011563862.1 that is a S41 family peptidase of Rubrobacter xylanophilus. After protein blast, 72 homologs sequences from 26 different species were downloaded in one FASTA file (homologs.txt), and rRNA sequences of these species were also downloaded in one FASTA file (rRNA_species.txt). The MSA analysis was done for each FASTA file (guidetree_homolos, guidtree_rRNA, align_homologs and align_rRNA), and A phylogenetic tree (Supplementary Figure 1 and Supplementary Figure 3) and a species tree (Supplementary Figure 2) were built to study the speciation events and duplication events.

## 3.3. Conserved regions

A variability plot (Figure 1) was generated using the resulting file from MSA analysis. The variability is lower, the sequence region is more likely to be a conserved region. In the variability list (Variability.txt), positions 141T, 170I, 195E, 223I, 252K, 284G, 286Q, 289S, 290V, 291K, 293A, 303V, 308T, 337T, 347T, 348I, 351S, 352E, 355A, 362D, 367V, 372Q, 374L, 376R, 397I,398Q,400Q,403G have the lowest variability 1 (start with 0), and from the plot, sequence region from 330 to 410 has relatively low variability.
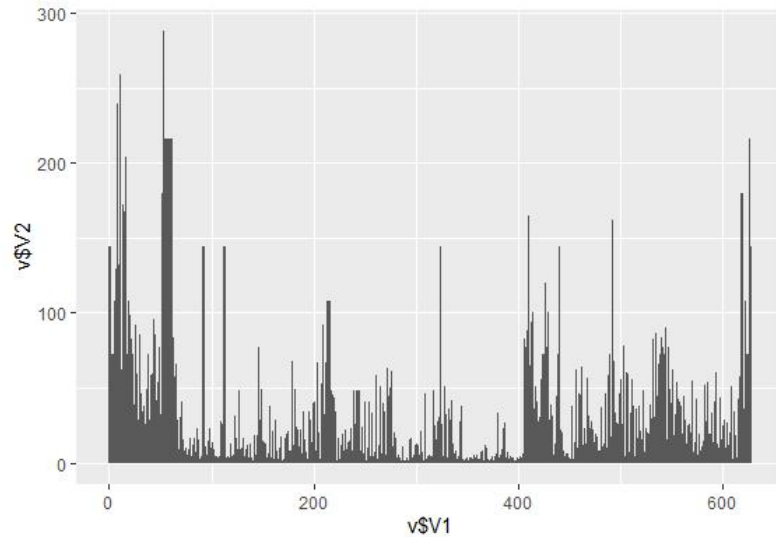
Figure 1: Variability plot of 72 homologs. X-axis represents each position in the sequence, Y-axis represent the variability value of each position.

### 3.4. Promotor regions

Three protein sequences were used to study the orthology, they are WP_012504624.1, its in-paralog WP_012505586.1 and the BBH WP_011563862.1 Both WP_012504624.1 and WP_012505586.1 are orthologous to WP_011563862.1, so a FASTA file containing three DNA sequences related to these three proteins were generated (orthologs.txt). Then the FASTA file was used to do the MSA analysis and generate an output file (align_orthologs). The alignment file was uploaded to PVS to create the variability plot (Figure 2). In the plot, positions from 0 to 70 have a region with relatively low variability.
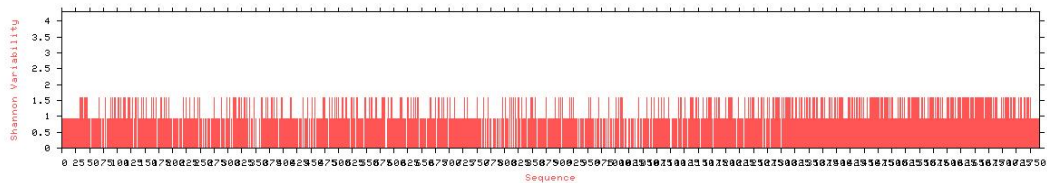


Figure 2: Variability plot of 3 orthologs. X-axis represents each position in the sequence, Y-axis represent the variability value of each position.

### 4. Discussion

804 pairs of BBH, 184 pairs of in-paralogs and a total of 1792 orthologs were found between proteomes of Prosthecochloris aestuarii (3205 protein-coding genes) and Rubrobacter xylanophilus (2302 protein-coding genes), which means a relatively large percent of protein-coding genes are orthologous between two species and some of them are even co-orthologous, this discovery indicates that these protein-coding genes are evolved from a common ancestral gene and retained similar functions in two species.

The species tree indicates that there exist at least five or six speciation events between these two species from the common ancestral species and the branch length is relatively long between the two species, which means the orthologous protein-coding genes that are still retained in the two species should have essential functions. In the phylogenetic tree with 72 different homologous

proteins, the orthologous proteins (BBH: WP 011563862.1 and WP 012505586.1) stay very close and the in-paralogous protein (WP 012504624.1) is nearby. There is only one speciation event before BBH, which indicates these orthologous proteins share a similar function, and this function should be essential for these two species, so this function has never changed a lot during evolution.

If proteins have similar functions, they should have similar sequences. For conserved regions among these homologs, they are likely to have some functions like the Post-translational modification (PTM) in the interaction, and these sites are more likely to be close to each other.

In the end of this study, I was trying to find the promotor regions in DNA sequences of these orthologs. One study (Jin et al., 2006) showed that promotor regions are more likely to be conserved in orthologs, and most sites of the region at the beginning of these DNA sequences from 0 to 70 have variability of 1, so this region is the potential promotor region for these orthologs.

## 5. References

Jin, V. X., Singer, G. A. C., Agosto-Pérez, F. J., Liyanarachchi, S., & Davuluri, R. V. (2006). Genome-wide analysis of core promoter elements from conserved human and mouse orthologous pairs. *BMC Bioinformatics*, *7*(1), 114–114. https://doi.org/10.1186/1471-2105-7-114

## 6. Appendix

**BBH.py:**
```
import pandas as pd
import numpy as np

BBH = {}
inparalogs = {}

# Create arrays for four different protein BLAST results
# output_Pa.txt: protein BLAST of protein sequences of Prosthecochloris aestuarii, database is protein
# sequences of Rubrobacter xylanophilus
# output_Rx.txt: protein BLAST of protein sequences of Rubrobacter xylanophilus, database is protein
# sequences of Prosthecochloris aestuarii
# output_Pa_self.txt: protein BLAST of protein sequences of Prosthecochloris aestuarii, database is protein
#sequences of Prosthecochloris aestuarii
# output_Rx_self.txt: protein BLAST of protein sequences of Rubrobacter xylanophilus, database is protein
# sequences of Rubrobacter xylanophilus
# In these four arrays, 1st column is the query name, 2nd column is the target name, 3rd column is the
# coverage rate, 4th column is the evalue
pa_df = pd.read_csv('D:/blast-BLAST_VERSION+/bin/output_Pa.txt', sep="\t", header=None)
pa_protein_hits = np.array(pa_df.iloc[:, [0, 1, 2, 10]])
print(len(pa_protein_hits))


rx_df = pd.read_csv('D:/blast-BLAST_VERSION+/bin/output_Rx.txt', sep="\t", header=None)
rx_protein_hits = np.array(rx_df.iloc[:, [0, 1, 2, 10]])
print(len(rx_protein_hits))


pa_self_df = pd.read_csv('D:/blast-BLAST_VERSION+/bin/output_Pa_self.txt', sep="\t", header=None)
pa_self_protein_hits = np.array(pa_self_df.iloc[:, [0, 1, 2, 10]])
print(len(pa_self_protein_hits))


rx_self_df = pd.read_csv('D:/blast-BLAST_VERSION+/bin/output_Rx_self.txt', sep="\t", header=None)
rx_self_protein_hits = np.array(rx_self_df.iloc[:, [0, 1, 2, 10]])
print(len(rx_self_protein_hits))


# For each protein sequence that finds the best hit in output_Pa.txt, check if each sequence has the
# in-paralogs
for i in range(len(pa_protein_hits)):
    protein = pa_protein_hits[i, 0]
    BH = pa_protein_hits[i, 1]
    evalue = pa_protein_hits[i, 3]

    # For each protein sequence that finds the best hit in output_Rx.txt, check if that sequence is the BH and
```

```python
        # if the BH of this sequence is the protein
        for m in range(len(rx_protein_hits)):
            if rx_protein_hits[m, 0] == BH and rx_protein_hits[m, 1] == protein:
                BBH[protein] = BH
                # Now we have find one pair of BBH, and the next is to check if there exists inparalogs
                # First check if rx_protein has an inparalogs
                rxprotein = rx_protein_hits[m, 0]
                rxevalue = rx_protein_hits[m, 3]
                for n in range(len(rx_self_protein_hits)):
                    # If two sequences from one species are more similar to each other than any sequence in
                    # the other specie, this sequence has inparalogs
                    # , and both sequences are orthologous to the BBH
                    if rx_self_protein_hits[n, 0] == rxprotein and rx_self_protein_hits[n, 2] != 100 and
rx_self_protein_hits[n, 3] < rxevalue:
                        rxinparalogs = rx_self_protein_hits[n, 1]
                        inparalogs[rxprotein] = rxinparalogs


                # Then check if pa_protein has inparalogs
                for j in range(len(pa_self_protein_hits)):
                    # If two sequences from one species are more similar to each other than any sequence in
                    # the other specie, this sequence has an inparalogs
                    # , and both sequences are orthologous to the BBH
                    if pa_self_protein_hits[j, 0] == protein and pa_self_protein_hits[j, 2] != 100 and
pa_self_protein_hits[j, 3] < evalue:
                        painparalogs = pa_self_protein_hits[j, 1]
                        inparalogs[protein] = painparalogs

print("The number of BBH pairs:", len(BBH))
print("The number of inparalogs:", len(inparalogs))
print("The number of orthologs:", len(BBH) * 2 + len(inparalogs))
```

**Conservation.py:**

```python
import pandas as pd
import numpy as np

# Obtain the multiple sequence alignment, and turn to array
MSA_df = pd.read_csv('C:/Users/LNH/Desktop/Study/Comparative/Comparative
Genomics/Assignment/align_homologs', sep="\t", header=None, skiprows=1)
MSA = np.array(MSA_df)

# First of all, we need to delete all useless information before sequence, which means to
# delete 36 positions in each line
# then, we need to remove all lines that do not contain sequence
remove = []
for i in range(len(MSA)):
    MSA[i, 0] = MSA[i, 0][36:]
    if '.' in MSA[i, 0] or ':' in MSA[i, 0] or '*' in MSA[i, 0]:
        remove.append(i)
i = 0
for j in range(len(remove)):
    MSA = np.delete(MSA, remove[j] - i, axis=0)
    i = i + 1

# In this research, we have 72 different protein sequences, MSA divides them into 13 regions,
# the last region contains 29 sites, and each of the rest contians 50 sites
proteins = 72
region_sites = len(MSA[1, 0])
last_regoin_sites = len(MSA[-1, 0])
regions = len(MSA)//proteins
columns = (regions - 1) * region_sites + last_regoin_sites
# We create a empty list with 936 rows and 50 columns
seq = np.empty((len(MSA), region_sites), dtype=str)

# For convenience, we add "1" to fill in the last region, so each region has 50 sites
for i in range(len(MSA) - proteins, len(MSA)):
    for j in range(region_sites - last_regoin_sites):
        MSA[i, 0] = MSA[i, 0] + "1"

# We move every site from MSA to the list 'seq', so each element in the 'seq' represents one
# site on the sequence
for i in range(len(MSA)):
    for j in range(region_sites):
        seq[i, j] = list(MSA[i, 0])[j]

# We reshape the seq, now we have 13 regions like MSA result, each region has 72 rows
```

```python
# (protein) and 50 columns(sites)
seq.shape = (regions, proteins, region_sites)
print(seq)


# Now we calculate varibality for each position
# Be careful! The position contain "1" should be ignored!
# We put all varibalities into one list 'Vs'
Vs = []
for i in range(regions):
    for j in range(region_sites):
        if seq[i][0][j] != "1":
            eachSite = []
            for m in range(proteins):
                eachSite.append(seq[i][m][j])
            n = 1
            for e in eachSite:
                if e != '-' and eachSite.count(e) > n:
                    max = e
                    n = eachSite.count(e)
            k = len(set(eachSite))
            N = proteins
            Variability = k * N / n
            Vs.append(Variability)


# Print all sites that has variability equal to 1, which means that amino acid conserved in all
# protwins
for i in range(len(Vs)):
    if Vs[i] == 1:
        print(i)
        print(seq[i//region_sites][0][i%regions])



# Create csv file for variablity, and we can plot the values in R
Vs_df = pd.DataFrame(data=Vs)
Vs_df.to_csv('C:/Users/LNH/Desktop/Study/Comparative/Comparative
Genomics/Assignment/Variability.txt', encoding = 'gbk', header=None)
```
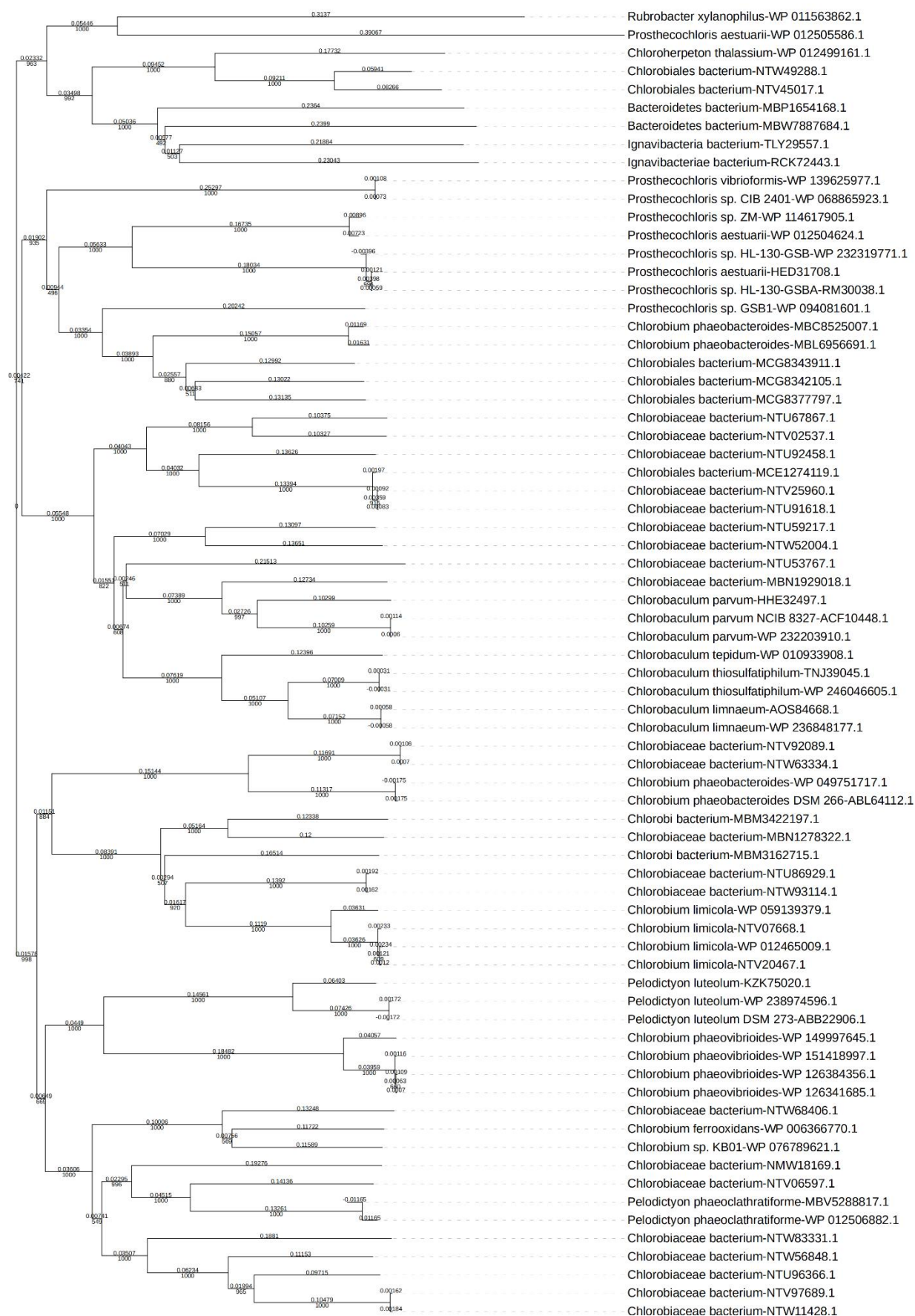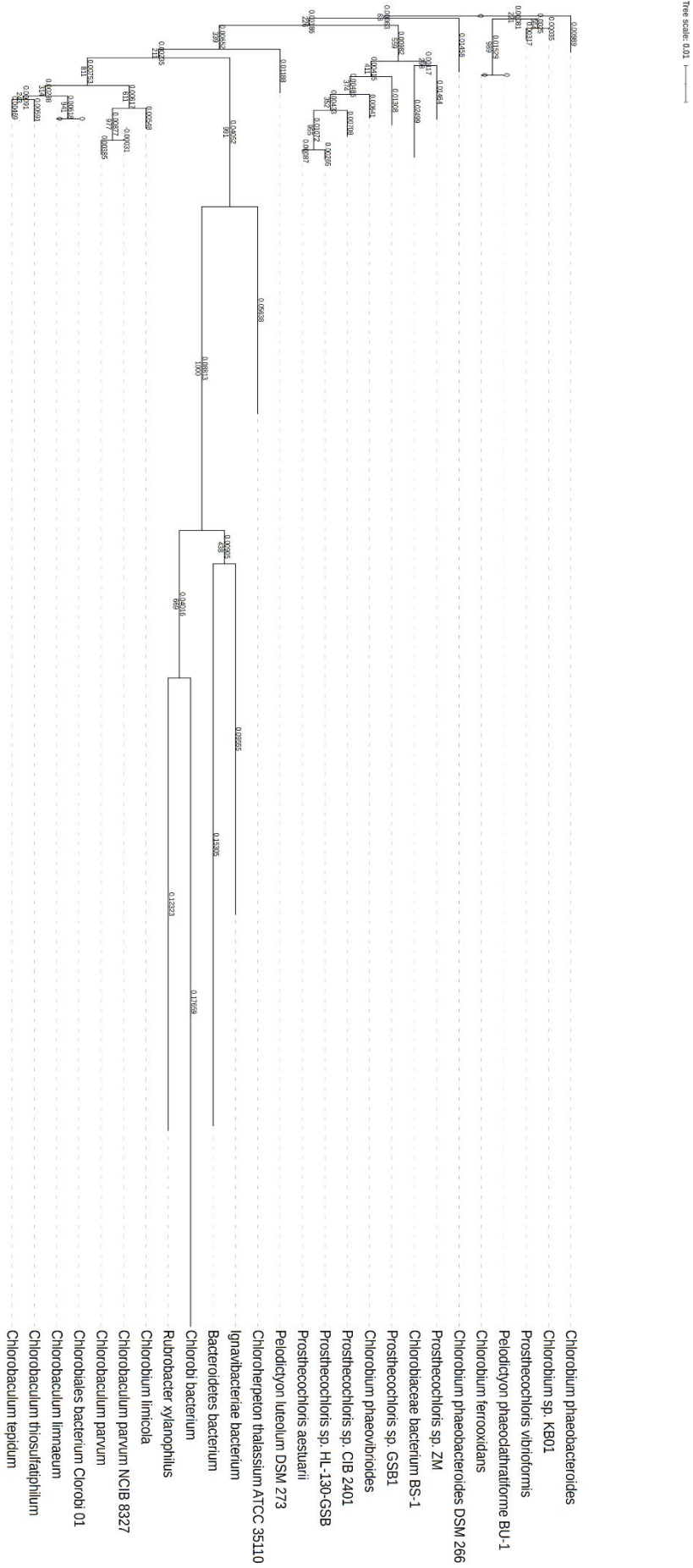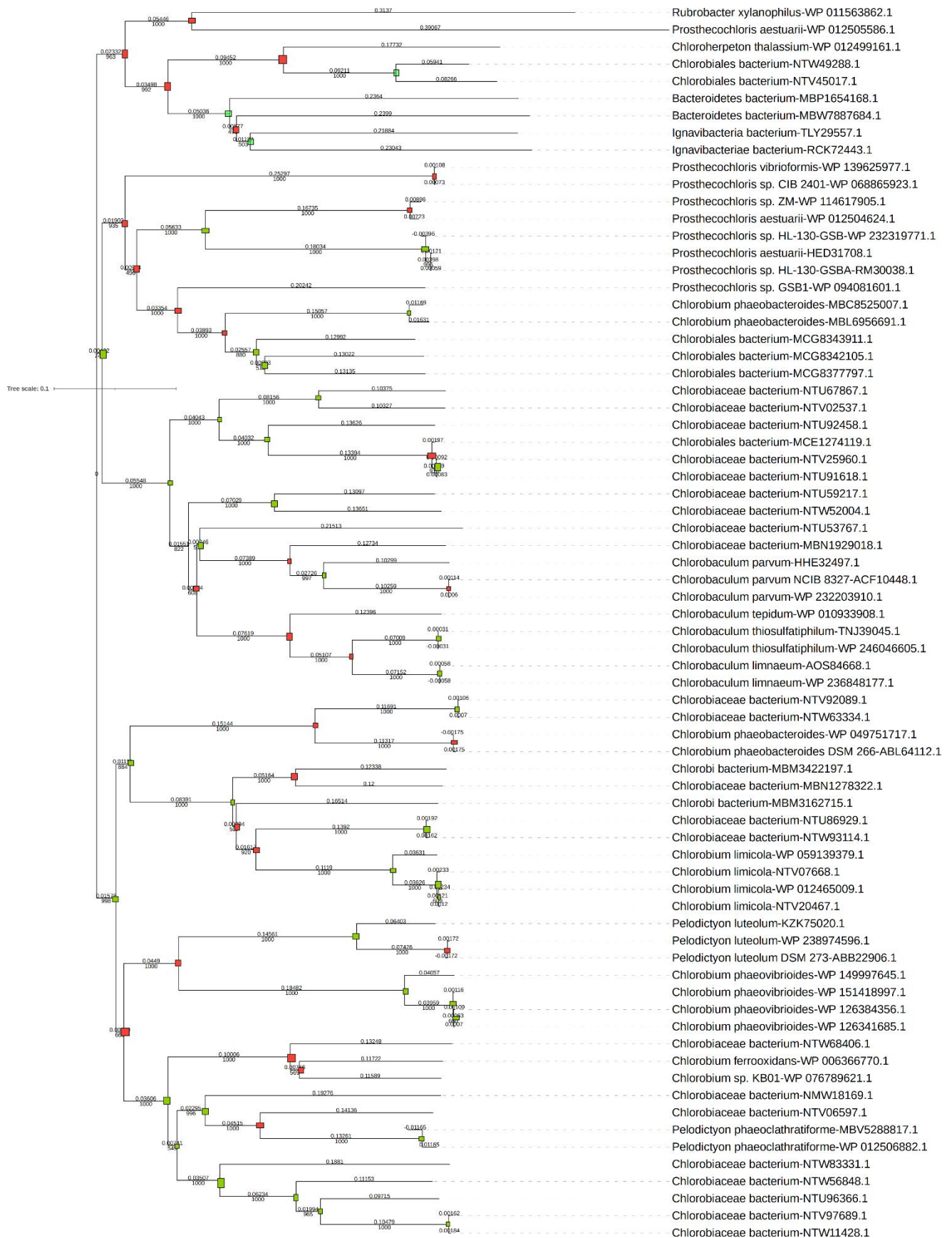
**Supplementary Figure 1**: Phylogenetic tree of 72 homologs

**Supplementary**
**Figure 2**: Species tree of 26 species

**Supplementary Figure 3**: Phylogenetic tree of 72 homologs with speciation events and duplications events. Red squares represent speciation events, and green squares represent duplication events.