

Look for co-evolution

Elena Montenegro, Sara González, Sofía Cerdá, Fernando Freire

December 26, 2018

Contents

1	Look for co-evolution Streptococcus-Lactobacillus	2
1.1	Obtain data from servers	2
1.2	Compute substitution rates	9
1.3	Selection of first set of target proteins	11
1.4	Obtain protein pathways from servers.	11
1.5	Selection of final set of target proteins	11
2	Generate document outputs.	12

1 Look for co-evolution Streptococcus-Lactobacillus

Co-evolution is everywhere. Because no species can be considered totally isolated from the others. And not only on a biological level, we perceive it in the lines of change in our societies, in political, scientific, religious ideas and even in the evolution of software engineering.

The difficulty is to know how much co-evolution is due to the interaction between a pair of species, because the tandem can not be totally isolated from the rest of the universe either. No doubt the problem can not be solved only from the bioinformatic perspective, it is also necessary the contribution of other branches of knowledge such as microbiology, physics, mathematics and computer theory.

Our approach consists in comparing the rates of protein evolutionary change between strains of Streptococcus and Lactobacillus linked by symbiotic pathways, and the rates between the rest of species of each genus.

We postulate that the proteins that reflects a rate change with more significance, i.e., more speedy or more slowly than intra genus rates, are proteins that could be influenced by the symbiotic environment.

To do so we need to compute separately both species of symbiotic tandem. The detailed steps are:

- 1) Obtain all the proteome from this four groups of species:
 - Genus streptococcus
 - Strains of streptococcus termophilus.
 - Genus lactobacillus
 - Strains of lactobacillus bulgaricus.
- 2) For each protein of each group we obtain the phylogenetic tree and compute the mean branch length. We suppose that we have ultrametricity or almost ultrametricity. It can be not exact but it could serve as a reference. We use the clustalw multialignment method, but other methods as T-COFFEE or MUSCLE could be used too. The branch length is a measure of the substitution rate.
- 3) For each protein of the first and second group, we calculate the ratio between branch length of termophilus and branch length of protein in genus. And also we compute the overall mean of all ratios. Finally we select the proteins with the most extreme values (80% percentile, two tailed).
- 4) We do the same with lactobacillus (third and fourth groups). At this stage we have two sets of proteins.
- 5) We obtain the biologic pathways of each of the sets of proteins.
- 6) The proteins involved in the same pathways from lactobacillus and streptococcus are the target proteins influenced by the coevolution between the two species. It will be necessary a later microbiological study that confirms this expectancies and dive deeply in the details of this coevolution.

1.1 Obtain data from servers

Script 1.1.1 (python)

```
1 import requests, sys
```

Script 1.1.2 (python)

```
1 def load_taxa(scientific_prefix):
2     """
3     """
4     requestURL = "https://www.ebi.ac.uk/prot eins/api/taxonomy/name/" + scientific_prefix + \
5         "%20?pageNumber=1&pageSize=100&searchType=STARTSWITH&fieldName=SCIENTIFICNAM
6 E"
7
8     r = requests.get(requestURL, headers={ "Accept" : "application/json"})
9
10    if not r.ok:
11        r.raise_for_status()
12        sys.exit()
13
14    jsonBody = json.loads(r.text)
15    taxa = []
16    names = []
17    for taxonomy in jsonBody["taxonomies"]:
18        print(taxonomy['taxonomyId'])
19        print(taxonomy['scientificName'])
20        taxa.append(taxonomy['taxonomyId'])
21        names.append(taxonomy['scientificName'])
22    return taxa, names
23
24 termophilus_taxa, termophilus_names = load_taxa("Streptococcus thermophilus")
25 streptococcus_taxa, streptococcus_names = load_taxa("Streptococcus")
```

Output

```
264199
Streptococcus thermophilus (strain ATCC BAA-250 / LMG 18311)
299768
Streptococcus thermophilus (strain CNRZ 1066)
322159
Streptococcus thermophilus (strain ATCC BAA-491 / LMD-9)
767463
Streptococcus thermophilus (strain ND03)
1042404
Streptococcus thermophilus CNCM I-1630
1051074
Streptococcus thermophilus JIM 8232
1073569
Streptococcus thermophilus MTCC 5460
1073570
Streptococcus thermophilus MTCC 5461
1091038
Streptococcus thermophilus DSM 20617
1187956
Streptococcus thermophilus MN-ZLW-002
1263110
Streptococcus thermophilus CAG:236
```

1268061
Streptococcus thermophilus DGCC 7710
1408178
Streptococcus thermophilus ASCC 1275
1415776
Streptococcus thermophilus TH1435
1423145
Streptococcus thermophilus TH1436
1433288
Streptococcus thermophilus MTH17CL396
1433289
Streptococcus thermophilus M17PTZA496
1435972
Streptococcus thermophilus TH985
1435974
Streptococcus thermophilus TH982
1435981
Streptococcus thermophilus 1F8CT
1436725
Streptococcus thermophilus TH1477
1302
Streptococcus gordonii
1303
Streptococcus oralis
1304
Streptococcus salivarius
1305
Streptococcus sanguinis
1306
Streptococcus sp.
1307
Streptococcus suis
1308
Streptococcus thermophilus
1309
Streptococcus mutans
1310
Streptococcus sobrinus
1311
Streptococcus agalactiae
1313
Streptococcus pneumoniae
1314
Streptococcus pyogenes
1317
Streptococcus downei
1318
Streptococcus parasanguinis
1319
Streptococcus sp. 'group B'
1320
Streptococcus sp. group G

1324
Streptococcus sp. G148
1325
Streptococcus sp. GX7805
1326
Streptococcus acidominimus
1328
Streptococcus anginosus
1329
Streptococcus canis
1332
Streptococcus criae
1333
Streptococcus criceti
1334
Streptococcus dysgalactiae
1335
Streptococcus equinus
1336
Streptococcus equi
1337
Streptococcus hyointestinalis
1338
Streptococcus intermedius
1339
Streptococcus macacae
1340
Streptococcus porcinus
1341
Streptococcus rattii
1343
Streptococcus vestibularis
1345
Streptococcus ferus
1346
Streptococcus iniae
1348
Streptococcus parauberis
1349
Streptococcus uberis
10728
Streptococcus pneumoniae phage HB-3
10747
Streptococcus phage Cp-1
10748
Streptococcus phage Cp-7
10749
Streptococcus phage Cp-9
12344
Streptococcus phage Cp-5
12366
Streptococcus pyogenes phage H4489A

12402
Streptococcus phage EJ-1
28037
Streptococcus mitis
29389
Streptococcus alactolyticus
29390
Streptococcus gordonii str. Challis
33040
Streptococcus milleri
33972
Streptococcus sp. 'group C'
35344
Streptococcus pyogenes phage T12
36470
Streptococcus sp. 'group A'
39425
Streptococcus phage T270
40041
Streptococcus equi subsp. zooepidemicus
40287
Streptococcus mutans serotype C
45634
Streptococcus cristatus
53354
Streptococcus gallolyticus
55085
Streptococcus thoraltensis
59310
Streptococcus macedonicus
64186
Streptococcus virus Sfi21
68891
Streptococcus peroris
68892
Streptococcus infantis
69017
Streptococcus sp. (strain 19909)
72638
Streptococcus virus Sfi19
73422
Streptococcus phage TP-J34
73492
Streptococcus pyogenes phage
74382
Streptococcus phage SFi18
76860
Streptococcus constellatus
78535
Streptococcus viridans
78541
Streptococcus virus Sfi11

82269
Streptococcus sp. 28D
82348
Streptococcus pluranimalium
82806
Streptococcus ovis
83545
Streptococcus sp. KN1
83546
Streptococcus sp. KN2
83547
Streptococcus sp. KN3
83549
Streptococcus sp. TW1
85154
Streptococcus phage 01205
86065
Streptococcus pyogenes phage H10403
90410
Streptococcus virus DT1
99822
Streptococcus dysgalactiae subsp. dysgalactiae
102143
Streptococcus phage S3b
102144
Streptococcus phage S92
102145
Streptococcus phage ST3
102146
Streptococcus phage ST64
102147
Streptococcus phage Sfi16A
102150
Streptococcus phage J1
102684
Streptococcus infantarius
102886
Streptococcus didelphis
104215
Streptococcus sp. 1400-98
112023
Streptococcus virus 7201
113107
Streptococcus australis
114652
Streptococcus orisratti
116154
Streptococcus sp. Z1227
116155
Streptococcus sp. Z12
116156
Streptococcus sp. Z89

```

118670
Streptococcus sp. PSH2
118671
Streptococcus sp. PSH1a
118672
Streptococcus sp. PSH1b
119224
Streptococcus phocae
119602
Streptococcus dysgalactiae subsp. equisimilis
119603
Streptococcus dysgalactiae group

```

Script 1.1.3 (python)

```

1 def load_proteome(taxids, size=10, protein=["LDH"]):
2     """
3     """
4     taxids_str = ",".join(str(x) for x in taxids)
5     protein_str = ",".join(x for x in protein)
6     print(taxids_str)
7     requestURL = "https://www.ebi.ac.uk/prot eins/api/prot eins?offset=0&size=" + str(size) +
8         ↳ "&taxid=" + \
9         taxids_str + "&reviewed=false"
10    if protein != []:
11        requestURL += "&gene=" + protein_str
12    print(requestURL)
13    r = requests.get(requestURL, headers={ "Accept" : "text/x-fasta"})
14
15    if not r.ok:
16        r.raise_for_status()
17        sys.exit()
18
19    proteome = r.text
20    return proteome
21
22 termophilus_taxids = termophilus_taxa[0:19]
23 streptococcus_taxids = streptococcus_taxa[0:19]
24 print(streptococcus_taxids)
25 print(termophilus_taxids)
26
27 streptococcus_proteome = load_proteome(streptococcus_taxids, -1, protein = ["LDH", "CAS2",
28     ↳ "CAS3"])
29 termophilus_proteome = load_proteome(termophilus_taxids, -1, protein = ["LDH", "CAS2",
30     ↳ "CAS3"])

```

Output

```

[1302, 1303, 1304, 1305, 1306, 1307, 1308, 1309, 1310, 1311, 1313, 1314, 1317, 1318, 1319,
↳ 1320, 1324, 1325, 1326]

```



```
[264199, 299768, 322159, 767463, 1042404, 1051074, 1073569, 1073570, 1091038, 1187956,
↪ 1263110, 1268061, 1408178, 1415776, 1423145, 1433288, 1433289, 1435972, 1435974]
1302,1303,1304,1305,1306,1307,1308,1309,1310,1311,1313,1314,1317,1318,1319,1320,1324,1325,1326
https://www.ebi.ac.uk/prot eins/api/prot eins?offset=0&size=-1&taxid=1302,1303,1304,1305,1306,1
↪ 307,1308,1309,1310,1311,1313,1314,1317,1318,1319,1320,1324,1325,1326&reviewed=false&gene=
↪ LDH,CAS2,CAS3
264199,299768,322159,767463,1042404,1051074,1073569,1073570,1091038,1187956,1263110,1268061,1
↪ 408178,1415776,1423145,1433288,1433289,1435972,1435974
https://www.ebi.ac.uk/prot eins/api/prot eins?offset=0&size=-1&taxid=264199,299768,322159,76746
↪ 3,1042404,1051074,1073569,1073570,1091038,1187956,1263110,1268061,1408178,1415776,1423145
↪ ,1433288,1433289,1435972,1435974&reviewed=false&gene=LDH,CAS2,CAS3
```

1.2 Compute substitution rates

Script 1.2.1 (python)

```
1 import re
2
3 def proteome2dict(proteome_fasta):
4     """
5     Returns a dict with keys protein accession and values the list of fasta format for all
↪ taxids
6     This is the basis for clustalw alignments and tree generation
7     """
8     proteome = {}
9     key_found = False
10    for line in proteome_fasta.splitlines():
11        if len(line) > 0:
12            if line[0] == ">":
13                if key_found:
14                    if key in proteome:
15                        proteome[key].append(seq)
16                    else:
17                        proteome[key] = [seq]
18                key_found = True
19                search_gene_name = re.search('GN=(\w*)', line)
20                key = search_gene_name.group(1).upper()
21                #print(key)
22                seq = line + '\n'
23            elif key_found:
24                seq += line + '\n'
25    if key_found:
26        if key in proteome:
27            proteome[key].append(seq)
28        else:
29            proteome[key] = [seq]
30    return proteome
```

Script 1.2.2 (python)

```

1  # Phylo tree with clustalw. We need to measure the substitution rate.
2  import matplotlib
3  import matplotlib.pyplot as plt
4  %matplotlib inline
5  from Bio import Phylo
6  from io import StringIO
7  import os
8  from Bio.Align.Applications import ClustalwCommandline
9  CLUSTALW = r"./clustalw2"
10 assert os.path.isfile(CLUSTALW), "Clustal W executable missing"
11 plt.rcParams["figure.figsize"] = (20,30)
12 matplotlib.rc('font', size=12)
13
14 def compute_mean_subst_rate(proteome, verbose=False, show_tree=False):
15     """
16     """
17     clustalw_cline = ClustalwCommandline(CLUSTALW, infile=proteome + ".fasta")
18     stdout, stderr = clustalw_cline()
19     f = open(proteome + ".dnd", "r")
20     s_tree = f.read()
21     f.close()
22     #print(s_tree)
23     branch_len = 0
24     num_branches = 0
25     search_branch_length = re.findall(':([-.0123456789]*)', s_tree)
26     for branch_length in search_branch_length:
27         #print(branch_length)
28         if branch_length != "0.00000":
29             branch_len += float(branch_length)
30             num_branches += 1
31     if verbose: print(branch_len, num_branches, branch_len/num_branches)
32     if show_tree:
33         tree = Phylo.read(proteome + ".dnd", "newick")
34         Phylo.draw(tree)
35     return branch_len/num_branches
36
37 def compute_subst_rates(proteome, proteome_name, verbose=False):
38     """
39     """
40     subst_rates = {}
41     for protein in proteome.keys():
42         if verbose: print(protein)
43         protein_sequence = ""
44         # Only for proteins with enough sequences to make a tree
45         if len(proteome[protein]) >= 3:
46             for sequence in proteome[protein]:
47                 protein_sequence += sequence
48             fasta_file_name = proteome_name + "_" + protein
49             f = open(fasta_file_name + ".fasta", "w")
50             if verbose: print(protein_sequence)
51             f.write(protein_sequence)

```

```

52         f.close()
53         mean_subst_rate = compute_mean_subst_rate(fasta_file_name)
54         subst_rates[protein] = mean_subst_rate
55     return subst_rates
56
57
58     # Proteomes in fasta to dictionaries
59     proteome_termophilus = proteome2dict(termophilus_proteome)
60     proteome_streptococcus = proteome2dict(streptococcus_proteome)
61
62     # Compute branch lengths
63     subst_rates_groups = {}
64     subst_rates_groups["termophilus"] = compute_subst_rates(proteome_termophilus, "termophilus",
65         ↪ False)
66     subst_rates_groups["streptococcus"] = compute_subst_rates(proteome_streptococcus,
67         ↪ "streptococcus", False)
68     print(subst_rates_groups)
69
70     # Compute branch ratios
71     # Compute mean of branch ratios and standard deviation
72     # Obtain the most extreme values. These are the proteins that could have been a slowdown or
73     ↪ from his initial state

```

Output

```

{'termophilus': {'CAS2': 0.10117692307692308, 'LDH': 0.007116666666666667}, 'streptococcus':
↪ {'CAS2': 0.05127011363636366, 'LDHA': 0.09308272727272726, 'LDH_2': 0.10562888888888889,
↪ 'LDH_1': 0.12399615384615387, 'LDHA_1': 0.14062904761904763, 'LDHA_3':
↪ 0.1900181818181818, 'LDH': 0.038460561797752806, 'LDHA_2': 0.13067238095238096, 'LDHD':
↪ 0.04888}}

```

1.3 Selection of first set of target proteins

At this step is necessary to calculate the ratios, the means of the ratios and select the more extreme between them (two tailed percentile 80%).

1.4 Obtain protein pathways from servers.

At this stage we obtain the pathways for proteins.

1.5 Selection of final set of target proteins

At this stage we select the proteins that belong tho the same pathways on Lactobacillus and Streptococcus. These are the proteins that could be affected by co-evolucion.

2 Generate document outputs.

Script 2.0.1 (text)

```
1 %%bash
2 #cd /Users/nandoide/Desktop/uni/STRBI.practical
3 jupyter nbconvert --to=latex --template=~/.report.tplx coevolution.ipynb 1> /dev/null
4 pdflatex -shell-escape coevolution 1> /dev/null
```

Output

```
[NbConvertApp] Converting notebook coevolution.ipynb to latex
[NbConvertApp] Writing 41101 bytes to coevolution.tex
```