

PHYLOGENETIC RECONSTRUCTION /Roadmap/

BLAST

Potential homologs candidates detection will depend on:

- Sample size
- E-value
- identity

Clean fasta headers¹

`cut -d " " -f 1 myBlast.fasta > myBlast_short.fasta`

Obtain Taxonomy

Extract ACCESSION numbers:

`grep ">" myBlast_short.fasta | sed -E 's/>/' | tr "\n" " "`

ACCESSION to Gis²

Run the script:

`perl change_fasta_headers.pl myBlast_short.fasta > myBlast_GIs.fasta`

Redundancy³

READ instructions. Use the program CD-HIT:

`cd-hit -i myBlast_GIs.fasta -o myBlast_GIs.nr98.fasta -c 0.98`

Alignment⁴

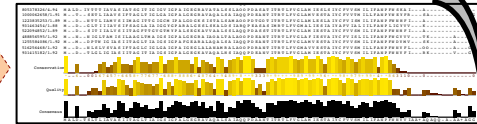
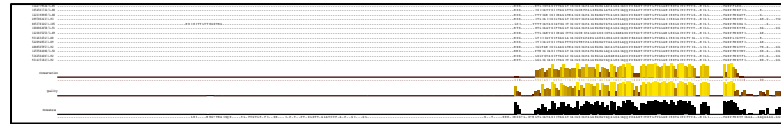
Use the program MAFFT:

`mafft myBlast_GIs.nr98.fasta > myBlast_GIs.nr98.aln.fasta`

Manual edition of the alignment?

"Before of manual edition"

"After of a manual edition"



Fasta to Phylip format

Transform FASTA to PHYLIP format (use online resources):

`myBlast_GIs.nr98.aln.fasta → sample.phylip → myBlast_GIs.nr98.aln.phy`

Phylogenetic tree⁵

Use RAXML to reconstruct a phylogenetic tree:

`raxmlHPC -m PROTGAMMAWAG -s myBlast_GIs.nr98.aln.phy -n myBlast_GIs.nr98.raxml -p 12345`

Bootstrapping⁶

RENAME the tree:

`nw_rename RAXML_result.myBlast_GIs.nr98.raxml myLABELS.txt > myTREE_labels.tree`

PLOT the tree, identify the HGT, and make the tree pretty in FigTree.

Final and crucial step to support the tree topology. **FIRST**, use RAXML to make a FAST-BOOTSTRAPPING and ML search:

`raxmlHPC-PTHREADS -f a -m PROTGAMMAWAG -p 12345 -x 12345 -# 100 -s myBlast_GIs.nr98.aln.phy -n myBlast_GIs.nr98.100_BOOTSTRAP_TREE.raxml`

SECOND, make a CONSENSUS tree using PHYLIP: `consense`

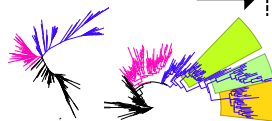
→

`myBlast_GIs.nr98.100_BOOTSTRAP_TREE.raxml`

→

Y

THIRD, rename the tree branch labels.



13

PARAMETRIC APPROACH /Roadmap/

GENOME

NCBI GENOME database:
<ftp://ftp.ncbi.nlm.nih.gov/genomes/>

LATEST version (CDS)

NCBI GENOME database. Use CDS
(coding sequence) file.

HGT annotated QUERY GENE

QUERY GENE with HGT and Genes
from the host/recipient genome.

Genomic Signatures

Estimate multiple genomic signatures:

```
codonw myCDS_GENOME.fasta myCDS_GENOME.outfile myCDS_GENOME.blk -all_indices -coa_cu  
-coa_rscu -cu -cutab -nomenu -nowarn -silent
```

Gene annotation

Make your own gene annotation: HOST and HGT.

Add a new column to the outfile created by codonw named as "myCDS_GENOME.outfile", and saved.

Use the R-script to make a simple plot GC% vs CAI.

```
Rscript make_plot.R myCDS_GENOME.outfile myPLOT_GC_vs_CAI.pdf mySAVED_Rsession.RData
```

BEST Genomic
signature for HGT

The best strategy (but not the unique) is show both signatures, one from the donor and other from the recipient.

Statistical Test

codonw offers multiple estimations calculated from the genome in study. Others analysis can show support a statistical or significant support. See the outfiles named as *.coa and read the manual of codonw for more info.

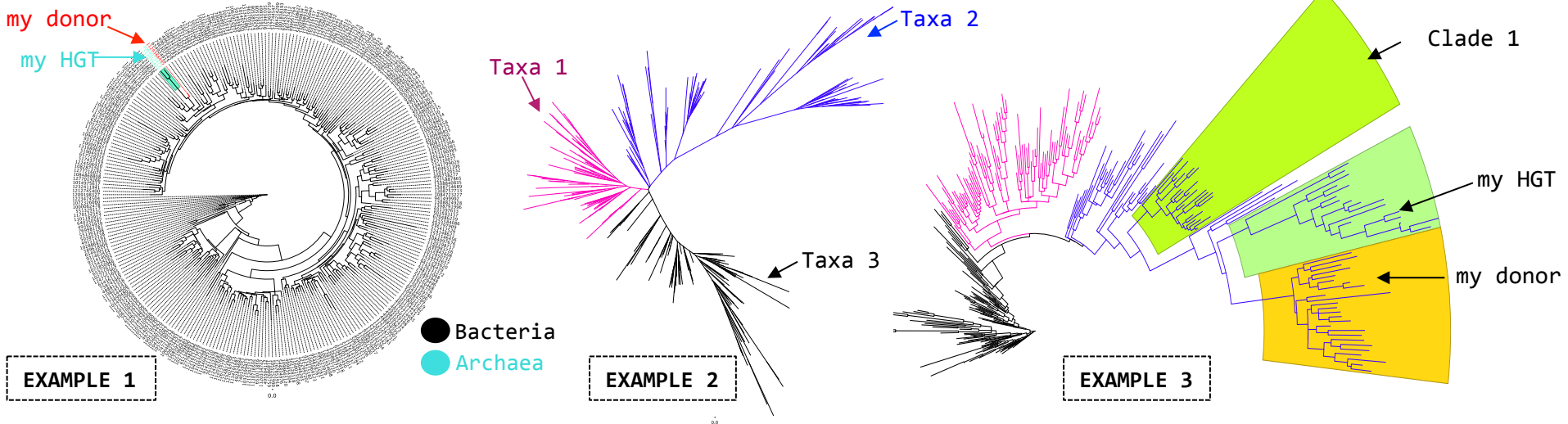
NOTES on "Phylogenetic Reconstruction"

Each step is critical, here are some advices:

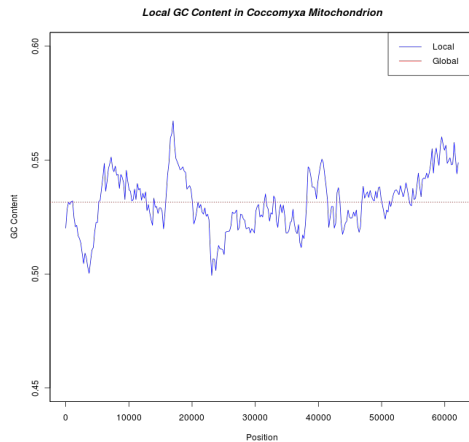
1. Clear the fasta header helps to manipulate the data, making it less confuse.
2. GI number contains max 10 characters and they are unique for each gene.
3. Why make a redundancy test? Gene copies or highly similar genes (e.g. from strains) might be remove.
4. Delete regions manually in the alignment has to be done in many cases. Why? For instance, it is used as criterion to improve the alignment or the homologs searching.
5. Previous to this step, is mandatory—if you don't know how your gene(s) 'evolve'—to identify their best substitution model (aa or nts). Use jModelTest or Prottest to achieve this purpose. Read our manual (step 11). Also, phylogenetic trees based on Maximum Likelihood (ML) or distances are good options, you should try a third approach: the Bayes methods. One good strategy/exercise show both tree topologies, ML and Bayes.

PHYLOGENETIC RECONSTRUCTION

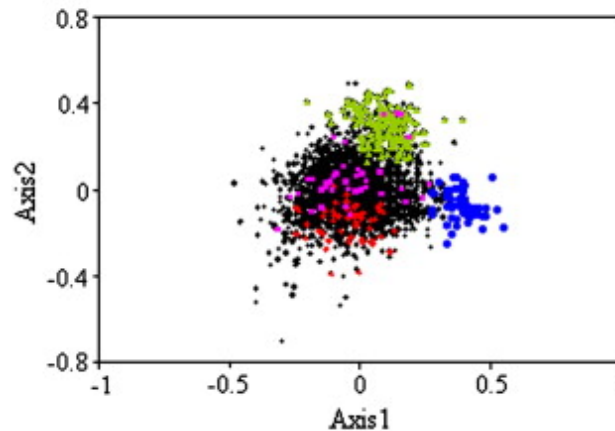
Make your tree pretty using [FigTree](#)



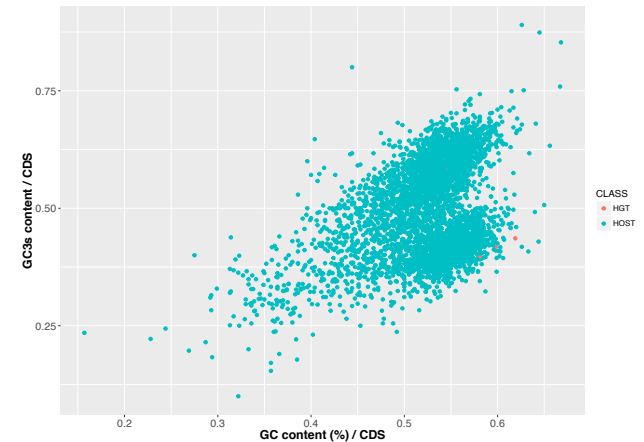
PARAMETRIC METHODS



EXAMPLE 1



EXAMPLE 2



EXAMPLE 3