

# Report Group 5

## Ancient HGT to Firm5 (from unknown donor)

**Joaquim, Virginie & Kamil**

We are going to investigate which genes of our Firm5 group are potentially acquired by ancient horizontal gene transfer from unknown donor.

Our `GitHub` repository is:

[https://github.com/KamilSJaron/SAGE\\_Firm5\\_specific\\_HGT.git](https://github.com/KamilSJaron/SAGE_Firm5_specific_HGT.git)

On Vital-IT, you can access the same repo by typing:

```
cd /scratch/beegfs/monthly/mls_2016/claivaz_ricci/SAGE_Firm5_specific_HGT
```

In this folder, README.md comprises all informations about our **folders** and **files**. To make it short, **data** contains all inputs and outputs files whereas **scripts** contains all **scripts** for data analysis.

**Folders are in red, files (inputs and outputs) in green and scripts in blue.**

### **Summary:**

We are going to select all gene families which contains orthologous genes if at least 80% of all strains in a given group of bees are present. Then, we are going to extract each protein sequences of all orthologous genes of each group of bees and blast them against RefSeq database (ncbi database). Then, we are going to parse blast results according to parameters. Thanks to blast results, we are going to determine hierarchical taxonomy and hierarchical taxonomy distance for each blast hit. This will allow us to plot % of identity of blast hits against hierarchical taxonomy distance. In the end, all blast hits corresponding to 'Lactobacillus' (in addition to supposed contaminations) will be discarded as we are interested in HGT from species other than Lactobacillus. From blast results, we will perform 2 different analysis. The first one is constructing phylogenetic trees of gene families potentially acquired by HGT and the second one is to predict genomic islands of gene families potentially acquired by HGT.

### **Group of bees:**

**Bumble\_bees:** B\_pascuorum+B\_bohemicus

('F225','F230','F233','F234','F236','F237')+('F228','F245','F246','F247')

**Honey\_bees:** A\_mellifera

('JF72','JF73','JG30','JF74','JF75','JF76','F259','F260','F261','F262','F263','L183','L184','L185','L186')

**Bumble\_Honey\_bees:** Bumble\_bees+Honey\_bees

**Outgroup:** other strains which are not Bumble\_bees and Honey\_bees

('JG29','LA14','LA2','LDB','LGAS','LHV','LJP','WANG')

### **Steps:**

1. For each group of bees (Bumble\_bees, Honey\_bees, Bumble\_Honey\_bees, Outgroup), **subgroup\_sort.py** looks into each `Gene_family_*` (each line) of **GeneFamilies.txt**. If a

Gene\_family\_\* line contains at least 80% of all strains present in the group of bees, it stores the Gene\_family\_\* line into its corresponding file.

*Input: GeneFamilies.txt (in data)*

*Outputs: Bumble\_bees.txt, Honey\_bees.txt, Bumble\_Honey\_bees.txt and Outgroup.txt (in data/sort\_group)*

2. For each group of bees and for each Gene\_family\_\*, [proteinseqextract.py](#) extracts all protein sequences for further 'Blast' step. This function creates a folder for each group of bees which comprises fasta files for each Gene\_family\_\* containing protein sequences.

*Inputs: outputs of subgroup\_sort.py in data/sort\_group (Bumble\_bees.txt, Honey\_bees.txt, Bumble\_Honey\_bees.txt or Outgroup.txt), all\_proteins.fasta (in mls\_2016/blast folder)*

*Outputs: folders in data of Bumble\_bees\_proteins, Honey\_bees\_proteins, Bumble\_Honey\_bees\_proteins - each folder contains .fasta files of gene families (Gene\_family\_\*.fasta) - a file concerns one gene family and contains all reference\_genome|protein\_ID and their sequence in fasta format*

3. For every protein sequences present in each Gene\_family\_\*.fasta files in each folder of group of bees (data/Bumble\_bees\_proteins, data/Honey\_bees\_proteins, data/Bumble\_Honey\_bees\_proteins), [blast\\_gene\\_families.sh](#) and [function\\_blast.sh](#) perform blastp against RefSeq database (module Blast is previously added on Vital-IT, see [blast\\_gene\\_families.sh](#)). RefSeq database is suitable as it is big (and manually curated) and it will allow us to access taxonomic informations from NCBI for further 'Phylogeny' step. [blast\\_gene\\_families.sh](#) runs blastp command specifying queries and outputs. [function\\_blast.sh](#) creates output folders (data/blast/Bumble\_bees\_proteins, data/blast/Honey\_bees\_proteins, data/blast/Bumble\_Honey\_bees\_proteins) and assigns variables used in [blast\\_gene\\_families.sh](#) (variables correspond to every Gene\_family\_\*.fasta files present in each folder of group bees and stored in [list\\_files.txt](#)).

*Inputs: outputs of proteinseqextract.py - Gene\_family\_\*.fasta in data/Bumble\_bees\_proteins, data/Honey\_bees\_proteins, data/Bumble\_Honey\_bees\_proteins (to select Gene\_family\_\*.fasta, function\_blast.sh will go through list\_files.txt - created in the terminal: ls Gene\_family\_\* > list\_files.txt)*

*Outputs: Gene\_family\_\*.out of its corresponding Gene\_family\_\*.fasta of its corresponding folder of group of bees (data/Bumble\_bees\_proteins, data/Honey\_bees\_proteins, data/Bumble\_Honey\_bees\_proteins) in its corresponding output folders (data/blast/Bumble\_bees\_proteins, data/blast/Honey\_bees\_proteins, data/blast/Bumble\_Honey\_bees\_proteins) - each Gene\_family\_\*.out file contains blast hits for each reference\_genome|protein\_ID's sequence - for each reference\_genome|protein\_ID's sequence (query): 5 first lines summarizing blastp information (version of BLASTP, query, database, fields and the total number of found hits) - fields correspond to the name of each resulting column (total of 14 columns) - subsequent lines correspond to every hits of the query*

4. For every blast results of each reference\_genome|protein\_ID's sequence (query) present in each Gene\_family\_\*.out files in each folder of group of bees (data/blast/Bumble\_bees\_proteins, data/blast/Honey\_bees\_proteins, data/blast/Bumble\_Honey\_bees\_proteins), [blast\\_hits\\_extract.py](#) parses hits according to different parameters (threshold\_alignment\_length, threshold\_ID, threshold\_eval, threshold\_bitscore):

- threshold\_alignment\_length : constant = 0.8
- threshold\_ID : starts at 50
- threshold\_eval : 0.00001
- threshold\_bitscore : 0

For a given *Gene\_family\_\*.out*, if there are less than 25 blast hits with parameters below, we relax the threshold\_ID by multiplying it by 0.8 until reaching at least 25 blast hits. *blast\_hits\_extract.sh* creates output folders (*data/parsed\_blast/Bumble\_bees\_proteins*, *data/parsed\_blast/Honey\_bees\_proteins*, *data/parsed\_blast/Bumble\_Honey\_bees\_proteins*) and assigns parameters as below.

*Inputs: outputs of blast\_gene\_families.sh and function\_blast.sh - Gene\_family\_\*.out in data/blast/Bumble\_bees\_proteins, data/blast/Honey\_bees\_proteins, data/blast/Bumble\_Honey\_bees\_proteins (to select Gene\_family\_\*.out, blast\_hits\_extract.py will go through list\_files.txt - created in the terminal: Is Gene\_family\_\* > list\_files.txt)*

*Outputs: Gene\_family\_\*\_parsed.out of its corresponding Gene\_family\_\*.out of its corresponding folder of group of bees (data/blast/Bumble\_bees\_proteins, data/blast/Honey\_bees\_proteins, data/blast/Bumble\_Honey\_bees\_proteins) in its corresponding output folders (data/parsed\_blast/Bumble\_bees\_proteins, data/parsed\_blast/Honey\_bees\_proteins, data/parsed\_blast/Bumble\_Honey\_bees\_proteins)*  
 - each *Gene\_family\_\*\_parsed.out* file contains the 'best' blast hits for each reference\_genome|protein\_ID's sequence - the header is '# Query\_ID Subject\_titles %\_Identity Alignment\_length evaluate bit\_score' - for each 'best' hits, informations are stored in its corresponding column

5. From every **parsed** blast results of each reference\_genome|protein\_ID's sequence (query) present in each *Gene\_family\_\*\_parsed.out* files in each folder of group of bees (*data/parsed\_blast/Bumble\_bees\_proteins*, *data/parsed\_blast/Honey\_bees\_proteins*, *data/parsed\_blast/Bumble\_Honey\_bees\_proteins*), *extract\_taxonomy\_hierarchy.py* first creates a list of subject IDs (hit IDs, for example 'WP\_052726720'). Subsequently, thanks to Entrez (from Bio package), the function extracts the hierarchical taxonomy of each subject IDs (from protein database of NCBI, for example: 'Bacteria;Firmicutes;Bacilli;Lactobacillales;Lactobacillaceae;Lactobacillus'). *The example ('WP\_052726720') is used as reference hierarchical taxonomy.* Finally, the function looks for the most recent taxa which is shared between each hierarchical taxonomy of each subject IDs and the reference hierarchical taxonomy. Hierarchical taxonomy distance are subjective:

- Lactobacillus = 1
- Lactobacillaceae = 2
- Lactobacillales = 3
- Bacilli = 4
- Firmicutes = 5
- Bacteria = 6
- None = 7 (which either corresponds to Archae or Eukaryota - contaminations)

The function creates a summary file (*hierarchical\_taxonomy.txt*) which is located in *data/parsed\_blast*.

*Inputs: Gene\_family\_\*\_parsed.out of its corresponding folder of group of bees (data/parsed\_blast/Bumble\_bees\_proteins, data/parsed\_blast/Honey\_bees\_proteins, data/parsed\_blast/Bumble\_Honey\_bees\_proteins)*

*Outputs: hierarchical\_taxonomy.txt in data/parsed\_blast - the first column corresponds to the strain IDs, the second one corresponds to the hierarchical taxonomy and the third one corresponds to the subjective hierarchical taxonomy distance*

6. For every **parsed** blast results of each reference\_genome|protein\_ID's sequence (query) present in each *Gene\_family\*\_parsed.out* files in each folder of group of bees (*data/parsed\_blast/Bumble\_bees\_proteins*, *data/parsed\_blast/Honey\_bees\_proteins*, *data/parsed\_blast/Bumble\_Honey\_bees\_proteins*), *plot\_identity\_taxonomy.R* extracts its percentage of identity and its hierarchical taxonomy distance (in *hierarchical\_taxonomy.txt*) and plot them against each other. Moreover, for each *Gene\_family\*\_parsed.out*, on one hand, we will perform linear regression model in addition to polynomial model and investigate if models differ. If they differ, it would mean that polynomial model is the best one and suggest potential horizontal gene transfer.

*Inputs: Gene\_family\*\_parsed.out of its corresponding folder of group of bees (data/parsed\_blast/Bumble\_bees\_proteins, data/parsed\_blast/Honey\_bees\_proteins, data/parsed\_blast/Bumble\_Honey\_bees\_proteins) and hierarchical\_taxonomy.txt in data/parsed\_blast*

*Outputs: Bumble\_taxo\_plot.pdf, Honey\_taxo\_plot.pdf, Bumble\_Honey\_taxo\_plot.pdf and potential\_HGT.txt in data/parsed\_blast - each \*.pdf file contains all plots of a given bee group - potential\_HGT.txt contains a list of Gene\_family\*\_parsed.out that are orthologous genes potentially acquired by HGT*

7. Manually, we decide which gene families are good candidates to investigate if orthologous genes are acquired by HGT. These gene family candidates are listed in *list\_files\_HGT\_to\_trees.txt* (in *data/parsed\_blast*). *Gene\_family\_1058\_parsed.out* and *Gene\_family\_991\_parsed.out* (*data/parsed\_blast/Bumble\_Honey\_bees\_proteins*) correspond to 'confirmed' HGT candidates. *Gene\_family\_1099\_parsed.out* (*data/parsed\_blast/Bumble\_Honey\_bees\_proteins*) and *Gene\_family\_1674\_parsed.out* (*data/parsed\_blast/Bumble\_bees\_proteins*) correspond to 'supposed' HGT candidates. *Gene\_family\_1048\_parsed.out* (*data/parsed\_blast/Bumble\_Honey\_bees\_proteins*) corresponds to non-HGT candidate. First, *extract\_seq\_blast\_hits.py* extracts the protein sequence of every blast hits of selected *Gene\_family\*\_parsed.out* (listed in *list\_files.txt*). From this, *align\_potential\_HGT.sh* aligns the sequence of every blast hits of selected *Gene\_family\*\_parsed.out* (listed in *list\_files.txt*) in its folder of group of bees (*data/phylogeny\_potential\_HGT/Bumble\_bees\_proteins*, *data/phylogeny\_potential\_HGT/Honey\_bees\_proteins*, *data/phylogeny\_potential\_HGT/Bumble\_Honey\_bees\_proteins*). Then, *tree\_wo\_bootstrap.sh* runs RAXML tool to infer phylogenetic trees (without bootstrapping) of all aligned protein sequences of every blast hits of selected *Gene\_family\*\_parsed.out* (in its bee folder in *data/phylogeny\_potential\_HGT*).

*extract\_seq\_blast\_hits.py:*

*Inputs: list\_files.txt in data/parsed\_blast, Gene\_family\*\_parsed.out of its corresponding folder of group of bees (data/parsed\_blast/Bumble\_bees\_proteins, data/parsed\_blast/Honey\_bees\_proteins, data/parsed\_blast/Bumble\_Honey\_bees\_proteins) listed in list\_files.txt*

*Outputs: amino\_acid\_seq\_Gene\_family\*\_parsed.fasta in data/parsed\_blast - each output file contains protein sequences of every blast hits of a selected gene family. The ID name of each sequence is defined as follows: StrainName\_SubjectID\_HierarchicalTaxonomyDistance*

`align_potential_HGT.sh`:

*Inputs: amino\_acid\_seq\_Gene\_family\_\*\_parsed.fasta in data/parsed\_blast*

*Outputs: Gene\_family\_\*\_parsed.multifasta in its corresponding folder of group of bees (data/phylogeny\_potential\_HGT/Bumble\_bees\_proteins, data/phylogeny\_potential\_HGT/Honey\_bees\_proteins, data/phylogeny\_potential\_HGT/Bumble\_Honey\_bees\_proteins) - each output file contains aligned protein sequences of every blast hits of a selected gene family*

`tree_wo_bootstrap.sh`:

*Inputs: Gene\_family\_\*\_parsed.multifasta in its corresponding folder of group of bees (data/phylogeny\_potential\_HGT/Bumble\_bees\_proteins, data/phylogeny\_potential\_HGT/Honey\_bees\_proteins, data/phylogeny\_potential\_HGT/Bumble\_Honey\_bees\_proteins)*

*Outputs: \*ML\_wo\_bootstrap\_Gene\_family\_\*\_parsed.out in data/phylogeny\_potential\_HGT/RAXML\_results*

>> RAXML\_bestTree.ML\_wo\_bootstrap\_Gene\_family\_\*\_parsed.out can then be opened using FigTree tool (available at <http://tree.bio.ed.ac.uk/software/figtree/>)

8. The inference of genomic islands is performed using Islandviewer tool (available at <http://www.pathogenomics.sfu.ca/islandviewer/>). For this purpose, two reference genomes from Bumble bees group (**F5\_237** and **F5\_245**) and one from Honey bees group (**Lb\_183**) are analysed. For each reference genome, `extract_coordinates_function_from_gbkfiles.py` extracts the protein name of every genes present in each Gene\_family\_\* from `GeneFamilies.txt` (in `data`). Thanks to GeneBank files (`mls_2016/genome_files`), the coordinates and the function of each protein are recovered. Then, we map manually the specific position of the protein onto its reference genome. Results are present in `data/IslandViewerResults`, `*.png` shows the circular representation of each reference genome and the position of the genomic islands, `*.tsv` contains the coordinates of each genomic island and which method allows this inference. Moreover, `MappingIntoGenome.txt` sums up the different results for this part.

*Input: GeneFamilies.txt in data, GeneBank files in mls\_2016/genome\_files*

*Outputs: \*.png files , \*.tsv files and MappingIntoGenome.txt in data/IslandViewerResults*