# - JSMC practical course script -
## Inferring phylogeny based on sequence information

## Collection of homologous sequences using synteny

Search for the homologous sequences of the following genes from *Arabidopsis thaliana*. The sequence search will be conducted on the Brassica Database Webseite (http://brassicadb.org/brad/searchSyntenytPCK.php).

Use the following Locus IDs for your search to find the homologous genes in the species *Brassica rapa*, *Schrenkiella parvula*, *Leavenworthia alabamica* and *Capsella rubella*.

**Locus-IDs:**

| | |
|---|---|
| APETALA1 | AT1G69120 |
| APETALA3 | AT3G54340 |
| PISTILLATA | AT5G20240 |
| AGAMOUS | AT4G18960 |
| SEPALLATA3 | AT1G24260 |
| AGAMOUS-LIKE6 | AT2G45650 |

The result will look as in Figure 1. The row highlighted in green marks the sought-after locus. If there is a homologous gene in the corresponding species, this will be indicated by a green dot. Click on the green dot to display and save the "gene sequence" (which is actually the coding sequence of the gene). Copy the coding sequences of all five species for all six genes to a fasta-file. To simplify recognition in the phylogeny that will be reconstructed, make sure to assign informative names to all of the sequences (including species and gene names and using underscores instead of blanks, e.g. Arabidopsis_thaliana_SEP3).

## Multiple sequence alignment (see above)

Translate the collected coding sequences into the corresponding amino acid sequences using 'Transeq' (http://www.ebi.ac.uk/Tools/st/emboss_transeq/).
Ensure that all coding sequences have been translated correctly by loading the amino acid sequences into Jalview (http://www.jalview.org/).
If all sequences have been translated correctly create a multiple sequence alignment using MAFFT (implemented in Jalview, web service > alignment > run MAFFT with preset).
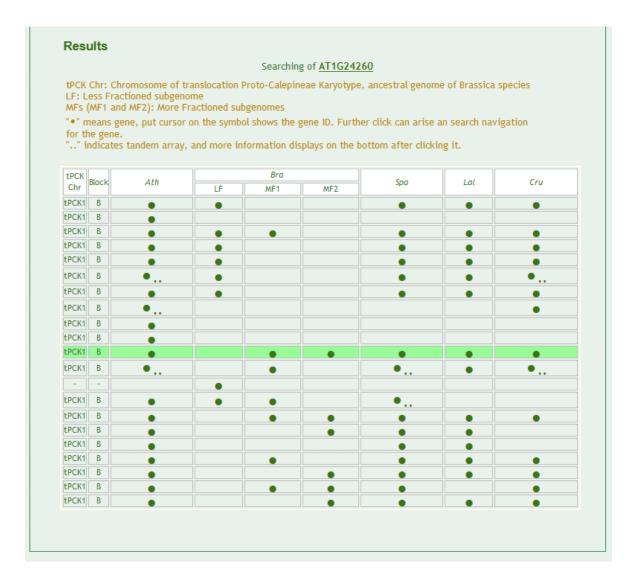
Figure 1: Result of synteny search

## Improving gene predictions using FGENESH+

Check the MAFFT alignment in Jalview and identify sequences that do not fit nicely into the alignment. Figure 2 shows two examples.
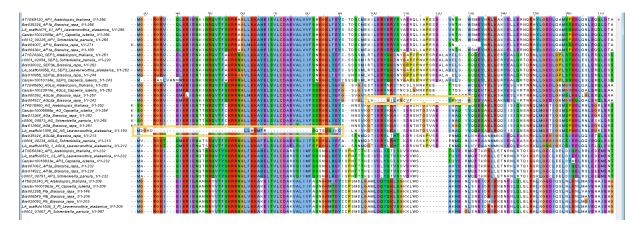


Figure 2: MAFFT alignment where sequences that do not fit nicely into the alignment are boxed.

Get the genomic locus of the suspiciously looking sequences using CoGe BLAST (https://genomevolution.org/coge/CoGeBlast.pl) as follows. Choose the corresponding species and click "+ Add" (Figure 3). Enter the coding sequence in the field "Query Sequence(s)" and click "Run CoGe BLAST.



Figure 3: BLAST search on CoGe.

On the table with the BLAST results, click on the HSP# 1 and on the then appearing pop-up window, click on the display of the subject locus (Figure 4). A new window will open with the genome browser of the corresponding species at the desired locus. Zoom out of the locus such that at least 20 Kb are

shown and view the genomic sequence by choosing "Sequence" > "Save track data" > "View".
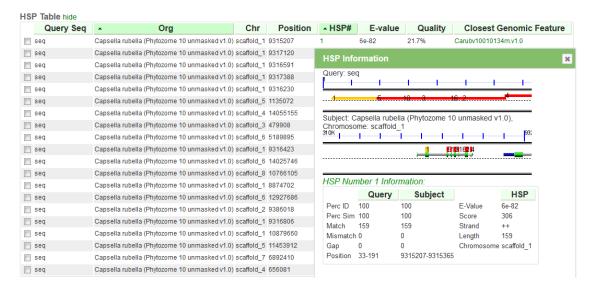


Figure 4: Result table of BLAST search on CoGe.

Copy the genomic sequence to the gene prediction tool FGENESH+ (http://www.softberry.com/cgi-bin/programs/gfs/fgenes_plus.pl). Paste the genomic sequence into the text area "Paste nucleotide sequence here". Into the other text area "Paste protein sequence here" enter the protein sequence of the most closely related gene from *Arabidopsis thaliana*. Under "Select organism specific gene-finding parameters" choose „Dicot plants, Arabidopsis (generic)" and press "search".

With the resulting gene prediction, first conduct a BLAST search with the protein sequence on NCBI (https://blast.ncbi.nlm.nih.gov/Blast.cgi?PROGRAM=blastp&PAGE_TYPE=BlastSearch&LINK_LOC=blasthome). If the BLAST search returns the expected results, replace the corresponding sequence in your fasta-file with the nucleotide sequence of the gene prediction.

Repeat this step until you predicted new genes for all sequences that did not fit nicely into the alignment.

## New phylogeny reconstruction

Add the sequences in your fasta-file to the sequence collection of 24 coding sequences from yesterday. Delete duplicate sequences from your new fasta-file (i.e. delete AP1, AP3, PI, AG and SEP3 from *Arabidopsis thaliana* **once** from your dataset [not both copies!]).

Translate the nucleotide sequences into amino acid sequences (see above). Align the sequences using MAFFT (see above). Remove alignment parts with low conservation using TrimAl (http://phylemon2.bioinfo.cipf.es/utilities.html > choose "start as anonymous user" > Utilities > Alignment Utilities > TrimAl (v. 1.3)). Upload your protein alignment and choose method "strict". Save the trimmed alignment and reconstruct a RAxML phylogeny via the online tool CIPRES (https://www.phylo.org/portal2/login!input.action). Upload the trimmed alignment, select the input file and the calculation tool (RAxML) and adjust the following parameter:

**Maximum Hours to Run (click here for help setting this correctly): 2**

**Please select the Data Type: Protein**

**Outgroup (one or more comma-separated outgroups, see comment for syntax): CgMADS1**

**Advanced Parameters**

**Conduct a rapid Bootstrap analysis and search for the best-scoring ML tree in one single program run. (-f a): Ticked**

**Bootstrap iterations (-#|-N): 1000**

Save the parameters and start the tree calculation.

After the tree calculation is completed display the output data via the 'output' button.

First of all open the error report file and check for any error messages.

The maximum likelihood tree is located in the 'RAxML_bipartitions.result' file. Download the file and visualize the tree with the FigTree software (http://tree.bio.ed.ac.uk/software/figtree/).