

Gene family expansion and contraction analysis

Overview

1. Preparing the pep files
2. Gene family cluster using Orthofinder
3. Reconstruct the ultrametric phylogeny tree with time using MCMCTree
4. Gene family expansion and contraction analysis using CAFE

Overview

This document is the pipeline to do gene family expansion and contraction analysis. Before you follow the pipeline, the **MOST IMPORTANT** thing is to choose the species you want to include, because the classification of orthogroups is always changing if you use different data set. Here are some tips to guide you to choose species:

https://davidemms.github.io/orthofinder_tutorials/orthofinder-best-practices.html

If you want to see the function of interested gene family, it's better to include *Arabidopsis thaliana* in your analysis, so that you could check the function of the gene/gene family easily in TAIR:

<https://www.arabidopsis.org>.

And also here are some blogs about how to do gene family expansion and contraction analysis:

<https://www.jianshu.com/p/8c6ef557cc71>

<https://www.jianshu.com/p/dc75116b0099>

<http://www.chenlianfu.com/?p=2974>

1. Preparing the pep files

Delete invalid transcripts, which include stop codon, lack start codon;

Delete transcripts less than 50 AA (150bp);

Only keep the primary transcripts;

run_pipeline_orthofinder.sh

Bash 复制代码

```

1 #1. extract valid transcripts,filter following transcripts:
2 ##Genes with internal stop codons(-g -V),lack start or end codon(-J);
3 ##Genes which length less than 50 amino acids(-l 150);
4 ##discard redundant transcripts(-M -K -Q)
5 ##remain coding only(-C)
6
7 :<<!
8 for i in `cat list_13sp`
9 do
10 gffread -g $i.fa -y $i.ext.merge.pep -V -M -K -J -Q -C -l 150 $i.gff3 &
11 gffread -g $i.fa -x $i.ext.merge.cds -V -M -K -J -Q -C -l 150 $i.gff3 &
12 done
13 !
14
15 #:<<!
16 #2. remain the primary (longest) transcripts as representation of the gene
17 ##the script was downloaded from the website
18 #make sure the geneID in fasta is uniq for each transcript
19 #if something wrong when running the script, FIRST check your data: fast
a, gff
20
21 python3 ../extract_primaryTranscript.py Juglans_mandshurica_NFU.ext.merge.
cds Juglans_mandshurica_NFU.gff3 Juglans_mandshurica_NFU.ext.cdsL &
22 python3 ../extract_primaryTranscript.py Juglans_mandshurica_NFU.ext.merge.
pep Juglans_mandshurica_NFU.gff3 Juglans_mandshurica_NFU.ext.pepL &
23 !

```

2. Gene family cluster using Orthofinder

put all protein files from last step into a directory, for example: 0703_PPJF

run_pipeline_orthofinder.sh

Bash 复制代码

```

1 nohup orthofinder -f 0703_PPJF -t 10 >log 2>err &

```

If everything goes well, you may see file **Orthogroups_SingleCopyOrthologues.txt** in the directory named Orthogroups. These single copy genes in file Orthogroups_SingleCopyOrthologues.txt will be used to reconstruct the phylogeny tree in the following analysis.

3. Reconstruct the ultrametric phylogeny tree with time using MCMCTree

The pipeline here is similar to the preparation the input of partitionFinder, just concentrate codon1, codon2 and codon3 of all the single copy genes into three big matrix. It's not definitely proper way to reconstruct the phylogeny, but for MCMCtree, it's no need to analyse each gene seperately, please search for detail in this blog if you want: <http://www.chenlianfu.com/?p=2974>

```

1 #0. modify the format of cds's name
2 #Mru:
3 #sed -r "/>/s/>(\S+) \[.*\] \[.*\] \[.*\] \[.*\]/>\1/g" Mru.cds >Mru.cds.1
4 #sed -r "/>/s/(\S+)_cds_(\S+)_(\S+)/>\2/" Mru.cds.1 >Mru.cds.2
5 :<<!
6 ln -s /data/data/Juglandaceae/Platycarya/13_geneExCon/genome/*.ext.merge.c
ds ./
7 sed -r 's/rna-gnl\|WGS:JAEDWW\|//g' Cil.ext.merge.cds >Cil.ext.merge.cds.2
8 sed -r 's/rna\-gnl\|WGS\:RXIC\|mrna\./g' Mru.ext.merge.cds >Mru.ext.merg
e.cds.2
9 !
10
11 #step 0 1 2
12 :<<!
13 for i in `cat list_11sp.index`
14 do
15     l=`echo ${i%_*}`
16     sp=`echo ${i#*_}`
17     colnu=${l+1}
18     #0. get the orthlist of each sp
19     # cut -f $colnu ./0_orthlist/Orthogroups_SingleCopyOrthologues.txt.2 |sed
'1d' >./0_orthlist/orthlist_$sp
20     #1. get singlecopy orthologues of each sp
21     # awk '/^>/&&NR>1{print "";}{ printf "%s",/^>/ ? $0"\n":$0 }' ./1_genefa/
$sp".ext.merge.cds" >./1_genefa/$sp".cds1"
22     # grep -A 1 -f ./0_orthlist/orthlist_$sp ./1_genefa/$sp".cds1" -w | sed "/"
^--$/d" > ./1_genefa/$sp".orth.cds"
23     #2. modify name of cds to orthfroup's name
24     # perl 2_substitute_orthname.pl ./0_orthlist/Orthogroups_SingleCopyOrtholo
gues.txt.2 $l ./1_genefa/$sp".orth.cds" &
25 done
26 !
27
28 #3. genefa to indfa
29 #perl 3_genefa2indfa.pl 0_orthlist/orthlist_0G ../list_11sp 1_genefa 2_ind
fa &
30
31 :<<!
32 #step 45 can be substitute by translatorx!!!
33 for i in `cat ../0_orthlist/orthlist_0G`
34 do
35     translatorx_vLocal.pl -i ${i}.ind.fa" -o ${i} -p maFft
36 done
37 !
38

```

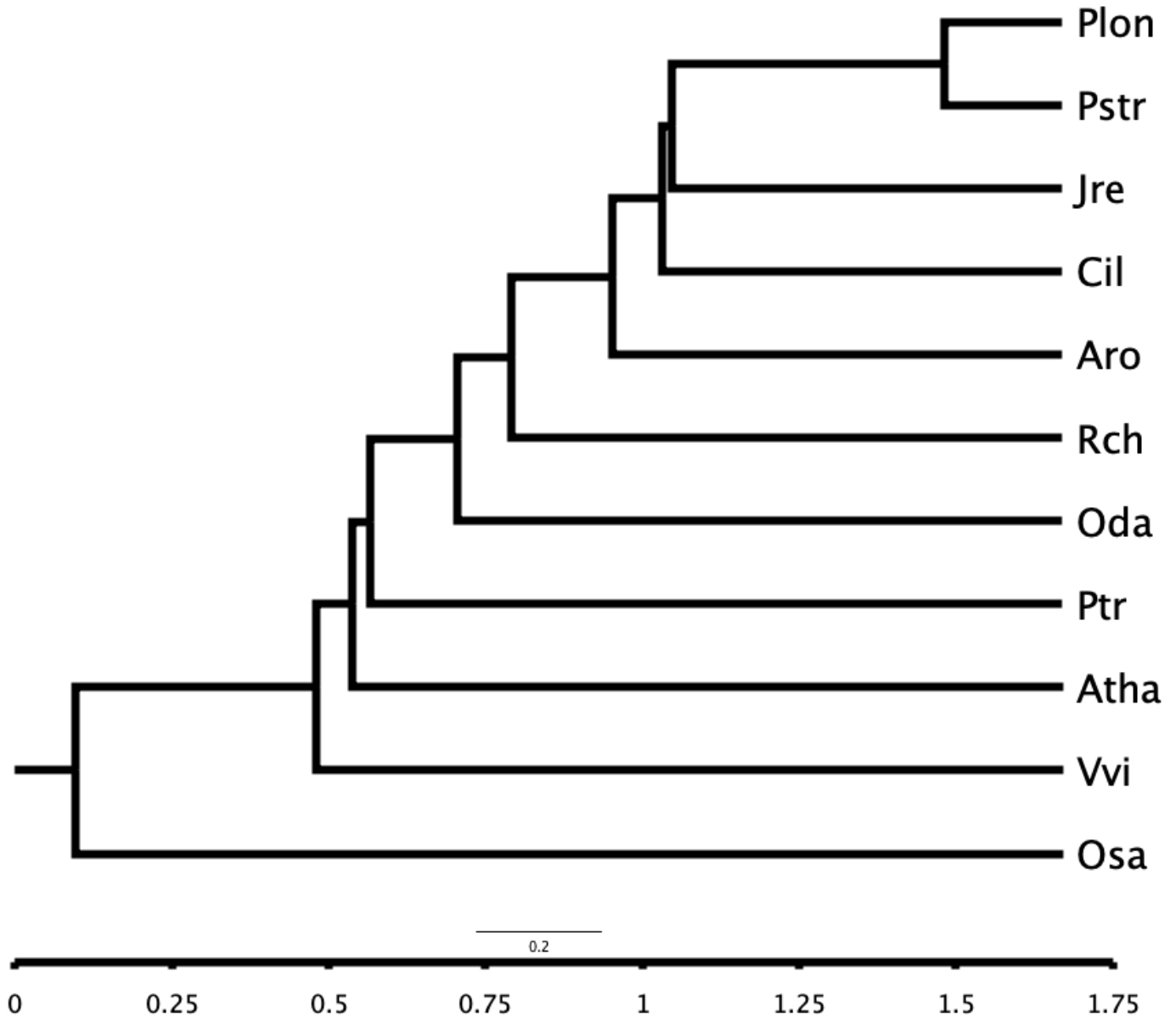
```

39  :<<!
40  #4. pal2nal
41  #nohup sh 4.1_mafft.sh orthlist 2_indfa &>log_mafft &
42  #sh 4.2_pal.sh orthlist 2_indfa
43  #5. partition
44  cd 2_indfa
45  mkdir 5_part_fa
46  sh ../5_run_codon123.sh ../orthlist
47  #creat gene123.list by:
48  cd 5_part_fa
49  mkdir bak
50  mv *12.fa bak
51  ls *.fa >../gene123.list
52
53  for i in `ls 5_part_fa/condon3/*.fa`
54  do
55      Gblocks $i -t=DNA
56  done
57  !
58
59  :<<!
60  ls *.nt1.ali.fasta >condon1.list
61  ls *.nt2.ali.fasta >condon2.list
62  ls *.nt3.ali.fasta >condon3.list
63  perl ../5_join_part.pl condon1.list . geneposition1.txt align_condon1.phy
64  0
65  perl ../5_join_part.pl condon2.list . geneposition2.txt align_condon2.phy
66  0
67  perl ../5_join_part.pl condon3.list . geneposition3.txt align_condon3.phy
68  0
69  cat align_condon1.phy align_condon2.phy align_condon3.phy >align_condon12
70  3.phy
71  !
72
73  #7. mcmctree
74  ##input file: 1:phylip format sequence file;
75  ## 2:ctl file
76  ## 3.tree file with fossil node time(CIs)
77
78  #mkdir 7_mcmctree
79  #nohup mcmctree mcmctree.ctl 1>log_mcmc_test_0831 2>err_mcmc_test_0831 &
80  :<<!
81  for i in {1..2}
82  do
83      mkdir -p run$i
84      cp mcmctree.ctl mcmctree.tree align_condon123.phy run$i/
85      cd run$i
86      nohup mcmctree mcmctree.ctl 1>log_mcmc_condon123 2>err_mcmc_condon123 &

```

```
83 cd ../
84 done
85 !
```

Now you get the phylogeny topology, **remember** the branch length of tree in MCMCtree is in the unit of **100Ma**, which means you should transform it into 1Ma in order to match CAFE software.



4. Gene family expansion and contraction analysis using CAFE

```

1 #0. prepare the inputfile of cafe
2 awk 'OFS="\t" {$NF="" ;print $0}' Orthogroups.GeneCount.tsv > cafe.data.1
3 sed 's/^/null\t&/g' cafe.data.1 >cafe.data
4 #modify title of cafe.data:Desc Family ID Cil Mru Plon Pstr Rch
5 python cafetutorial_clade_and_size_filter.py -i cafe.data -o cafe.filter.d
  ata -s
6 ##. modify cafetutorial_run.sh
7
8 #1.run cafe by shell script
9 nohup cafe cafetutorial_run.sh &
10 nohup cafe cafetutorial_run_filter.sh &
11 #2. summarize the result of cafe
12 #output the rapidly changing families on each nodes
13 python2 /data/data/Juglandaceae/Platycarya/13_geneExCon/python_scripts/caf
  etutorial_report_analysis.py -i resultfile.cafe -o summary_cafe_rapidChang
  e >log_summary_cafe_rapidChange &
14 python2 /data/data/Juglandaceae/Platycarya/13_geneExCon/python_scripts/caf
  etutorial_report_analysis.py -i resultfile.largefilter.cafe -o summary_caf
  e_largefilter_rapidChange >log_summary_cafe_largefilter_rapidChange &
15 #ouput all changing families on each nodes
16 python2 /data/data/Juglandaceae/Platycarya/13_geneExCon/python_scripts/caf
  etutorial_report_analysis.py -i resultfile.cafe -o summary_cafe_allChange
  -r 0 >log_summary_cafe_allChange &
17 python2 /data/data/Juglandaceae/Platycarya/13_geneExCon/python_scripts/caf
  etutorial_report_analysis.py -i resultfile.largefilter.cafe -o summary_caf
  e_largefilter_allChange >log_summary_cafe_largefilter_allChange &
18
19 #3.1 visualize through cafe home script, the figures are quite ugly thoug
  h :C
20 python3.6 /data/data/Juglandaceae/Platycarya/13_geneExCon/python_scripts/c
  afetutorial_draw_tree.py -i summary_cafe_rapidChange_node.txt -t '(((Plo
  n:2,Pstr:2):62,Cil:64):23,Rch:87):10,0dav:97)' -d '(((Plon<0>,Pstr<2>><1
  >,Cil<4>><3>,Rch<6>><5>,0dav<8>><7>' -o summary_cafe_rapidChange_node_expa
  nd.png
21 python3.6 /data/data/Juglandaceae/Platycarya/13_geneExCon/python_scripts/c
  afetutorial_draw_tree.py -i summary_cafe_rapidChange_node.txt -t '(((Plo
  n:2,Pstr:2):62,Cil:64):23,Rch:87):10,0dav:97)' -d '(((Plon<0>,Pstr<2>><1
  >,Cil<4>><3>,Rch<6>><5>,0dav<8>><7>' -y Contractions -o summary_cafe_rapid
  Change_node_contract.png
22
23 #3.2 visualize through cafe_fig
24 python3.6 ~/software/CAFE_fig-master/CAFE_fig.py resultfile.cafe -pb 0.01
  -pf 0.01 --dump test/ -g pdf --count_all_expansions
25
26 #4. get orthologues from output of cafe

```

```

27 #rapidly changing OGs' list
28 sed -n '7p' summary_cafe_rapidChange_fams.txt |sed 's/,/\n/g;s/:/\n/g;s/
29 \t//g'|grep "+"|cut -c1-9 >signif_expandOG_inPS.list
30 sed -n '6p' summary_cafe_rapidChange_fams.txt |sed 's/,/\n/g;s/:/\n/g;s/
31 \t//g'|grep "+"|cut -c1-9 >signif_expandOG_inPL.list
32 sed -n '7p' summary_cafe_rapidChange_fams.txt |sed 's/,/\n/g;s/:/\n/g;s/
33 \t//g'|grep "-"|cut -c1-9 >signif_contractOG_inPS.list
34 sed -n '6p' summary_cafe_rapidChange_fams.txt |sed 's/,/\n/g;s/:/\n/g;s/
35 \t//g'|grep "-"|cut -c1-9 >signif_contractOG_inPL.list
36
37 #all changing OGs' list
38 sed -n '7p' summary_cafe_allChange_fams.txt |sed 's/,/\n/g;s/:/\n/g;s/\t//
39 g'|grep "+"|cut -c1-9 >all_expandOG_inPS.list
40 sed -n '6p' summary_cafe_allChange_fams.txt |sed 's/,/\n/g;s/:/\n/g;s/\t//
41 g'|grep "+"|cut -c1-9 >all_expandOG_inPL.list
42 sed -n '7p' summary_cafe_allChange_fams.txt |sed 's/,/\n/g;s/:/\n/g;s/\t//
43 g'|grep "-"|cut -c1-9 >all_contractOG_inPS.list
44 sed -n '6p' summary_cafe_allChange_fams.txt |sed 's/,/\n/g;s/:/\n/g;s/\t//
45 g'|grep "-"|cut -c1-9 >all_contractOG_inPL.list
46
47 #reform the OG-gene pair file
48 python split_with_one_gene.py Orthogroups.txt Orthogroups_genes.txt
49
50 #get genelist of changing OGs
51 grep -f signif_expandOG_inPS.list Orthogroups_genes.txt|grep "Pstr" >signi
52 f_expandOG_inPS.OGs.genes
53 grep -f signif_expandOG_inPL.list Orthogroups_genes.txt|grep "Plon" >signi
54 f_expandOG_inPL.OGs.genes
55 grep -f signif_contractOG_inPS.list Orthogroups_genes.txt|grep "Pstr" >sig
56 nif_contractOG_inPS.OGs.genes
57 grep -f signif_contractOG_inPL.list Orthogroups_genes.txt|grep "Plon" >sig
58 nif_contractOG_inPL.OGs.genes
59
60 cut -f2 signif_contractOG_inPL.OGs.genes >signif_contractOG_inPL.genes
61 cut -f2 signif_contractOG_inPS.OGs.genes >signif_contractOG_inPS.genes
62 cut -f2 signif_expandOG_inPL.OGs.genes >signif_expandOG_inPL.genes
63 cut -f2 signif_expandOG_inPS.OGs.genes >signif_expandOG_inPS.genes
64
65 grep -f all_expandOG_inPS.list Orthogroups_genes.txt|grep "Pstr" >all_expa
66 ndOG_inPS.OGs.genes
67 grep -f all_expandOG_inPL.list Orthogroups_genes.txt|grep "Plon" >all_expa
68 ndOG_inPL.OGs.genes
69 grep -f all_contractOG_inPS.list Orthogroups_genes.txt|grep "Pstr" >all_co
70 ntractOG_inPS.OGs.genes
71 grep -f all_contractOG_inPL.list Orthogroups_genes.txt|grep "Plon" >all_co
72 ntractOG_inPL.OGs.genes
73
74 cut -f2 all_contractOG_inPL.OGs.genes >all_contractOG_inPL.genes

```



```
59 cut -f2 all_contract0G_inPS.0Gs.genes >all_contract0G_inPS.genes
60 cut -f2 all_expand0G_inPL.0Gs.genes >all_expand0G_inPL.genes
61 cut -f2 all_expand0G_inPS.0Gs.genes >all_expand0G_inPS.genes
62
```

As for the visualization of the output, the best way is to add the number of expand/contract genefamily (or genes included) to the phylogeny by yourself.

[GenefamilyExpCon_11sps.2.pdf](#)

And you can also explore the function of genefamilies which are expand in specific taxon (taxa), by doing GO enrichment or KEGG enrichment, or both.