

Run GENESPACE vignettes - Cotton genomes treated as tetraploids

JTLovell

13-Feb 2022

1. Global Parameters

```
baseDir <- "/Users/jlovell/Desktop/GENESPACE_data/results/"
rawAnnotationDir <- "/Users/jlovell/Desktop/GENESPACE_data"
mcscanDir <- "/Users/jlovell/Documents/comparative_genomics/programs/MCScanX"
nThreads <- 6
path2of <- "orthofinder" # since orthofinder is in the path via conda
```

2. Set parameters

NOTE: Wheat, Maize and Switchgrass WGD occurred after MRCA of all genomes below. Therefore, treat them as polyploid.

Since some ploidy > 1, run orthofinder again in syntenic blocks.

```
cotton <- list(
  wd = file.path(baseDir, "cotton4x"),
  speciesIDs = "cotton",
  genomes = data.table(do.call("rbind", list(
    c(genome = "Gbarbadense", version = "Gbarbadense", ploidy = 2),
    c("Gdarwinii", "Gdarwinii", 2),
    c("Gtomentosum", "Gtomentosum", 2))))))

cottonParams <- list(
  pepString = "fa",
  orthofinderInBlk = T,
  blkSize = 10,
  nGaps = 10)
```

3 Initialize the rho run with the above specified parameters

```
gparCotton4x <- with(cotton, init_genespace(
  genomeIDs = genomes$genome,
  versionIDs = genomes$version,
  ploidy = genomes$ploidy,
  speciesIDs = rep(speciesIDs, length(genomes$genome)),
```

```

orthofinderInBlk = TRUE,
pepString = "fa",
wd = wd,
nCores = nThreads,
path2orthofinder = path2of,
path2mcscanx = mcscanDir,
rawGenomeDir = rawAnnotationDir))

## set working directory to /Users/jlovell/Desktop/GENESPACE_data/results/cotton4x
##
## found raw gff files:
##   /Users/jlovell/Desktop/GENESPACE_data/cotton/Gbarbadense/annotation/Gbarbadense.gff3.gz
##   /Users/jlovell/Desktop/GENESPACE_data/cotton/Gdarwinii/annotation/Gdarwinii.gff3.gz
##   /Users/jlovell/Desktop/GENESPACE_data/cotton/Gtomentosum/annotation/Gtomentosum.gff3.gz
##
## found raw peptide files:
##   /Users/jlovell/Desktop/GENESPACE_data/cotton/Gbarbadense/annotation/Gbarbadense.fa.gz
##   /Users/jlovell/Desktop/GENESPACE_data/cotton/Gdarwinii/annotation/Gdarwinii.fa.gz
##   /Users/jlovell/Desktop/GENESPACE_data/cotton/Gtomentosum/annotation/Gtomentosum.fa.gz
##
## Can't find all parsed annotation files ... need to run parse_annotations, parse_ncbi or parse_phytoz
##
## GENESPACE run initialized:
##   Initial orthofinder database generation method: default inside R
##   Orthology graph method: inBlock

```

4 Parse the raw annotations.

```

parse_phytozome(gsParam = gparCotton4x)

## Parsing annotations: Gbarbadense
##   Reading gff ... found 74561 protein coding genes
##   Reading peptide fasta ... found 74561 / 74561 total and unique entries
##   Merging fa and gff found 74561 matching entires
## Parsing annotations: Gdarwinii
##   Reading gff ... found 78303 protein coding genes
##   Reading peptide fasta ... found 78303 / 78303 total and unique entries
##   Merging fa and gff found 78303 matching entires
## Parsing annotations: Gtomentosum
##   Reading gff ... found 78338 protein coding genes
##   Reading peptide fasta ... found 78338 / 78338 total and unique entries
##   Merging fa and gff found 78338 matching entires

```

5 Get orthofinder results

```

gparCotton4x <- set_syntenyParams(gsParam = gparCotton4x)
gparCotton4x <- run_orthofinder(gsParam = gparCotton4x)

## Running 'default' genespace orthofinder method

```

```

## #####
## Cleaning out orthofinder directory and prepping run
## Calculating blast results and running OrthoFinder
## #####
## #####
##
##
## OrthoFinder version 2.5.4 Copyright (C) 2014 David Emms
##
## 2022-02-13 07:16:36 : Starting OrthoFinder 2.5.4
## 6 thread(s) for highly parallel tasks (BLAST searches etc.)
## 1 thread(s) for OrthoFinder algorithm
##
## Checking required programs are installed
## -----
## Test can run "mcl -h" - ok
## Test can run "fastme -i /Users/jlovell/Desktop/GENESPACE_data/results//cotton4x/orthofinder/Results
##
## Dividing up work for BLAST for parallel processing
## -----
## 2022-02-13 07:16:38 : Creating diamond database 1 of 3
## 2022-02-13 07:16:39 : Creating diamond database 2 of 3
## 2022-02-13 07:16:40 : Creating diamond database 3 of 3
##
## Running diamond all-versus-all
## -----
## Using 6 thread(s)
## 2022-02-13 07:16:40 : This may take some time....
## 2022-02-13 07:16:41 : Done 0 of 9
## 2022-02-13 08:10:25 : Done all-versus-all sequence search
##
## Running OrthoFinder algorithm
## -----
## 2022-02-13 08:10:26 : Initial processing of each species
## 2022-02-13 08:11:19 : Initial processing of species 0 complete
## 2022-02-13 08:12:15 : Initial processing of species 1 complete
## 2022-02-13 08:13:10 : Initial processing of species 2 complete
## 2022-02-13 08:13:24 : Connected putative homologues
## 2022-02-13 08:13:30 : Written final scores for species 0 to graph file
## 2022-02-13 08:13:36 : Written final scores for species 1 to graph file
## 2022-02-13 08:13:42 : Written final scores for species 2 to graph file
## 2022-02-13 08:14:23 : Ran MCL
##
## Writing orthogroups to file
## -----
## OrthoFinder assigned 223378 genes (96.6% of total) to 53265 orthogroups. Fifty percent of all genes
##
## 2022-02-13 08:15:01 : Done orthogroups
##
## Analysing Orthogroups
## =====
##
## Calculating gene distances
## -----

```

```

## 2022-02-13 08:17:34 : Done
## 2022-02-13 08:17:37 : Done 0 of 24863
## 2022-02-13 08:17:40 : Done 1000 of 24863
## 2022-02-13 08:17:43 : Done 2000 of 24863
## 2022-02-13 08:17:46 : Done 3000 of 24863
## 2022-02-13 08:17:49 : Done 4000 of 24863
## 2022-02-13 08:17:53 : Done 5000 of 24863
## 2022-02-13 08:17:56 : Done 6000 of 24863
## 2022-02-13 08:17:59 : Done 7000 of 24863
## 2022-02-13 08:18:03 : Done 8000 of 24863
## 2022-02-13 08:18:06 : Done 9000 of 24863
## 2022-02-13 08:18:09 : Done 10000 of 24863
## 2022-02-13 08:18:12 : Done 11000 of 24863
## 2022-02-13 08:18:16 : Done 12000 of 24863
## 2022-02-13 08:18:19 : Done 13000 of 24863
## 2022-02-13 08:18:22 : Done 14000 of 24863
## 2022-02-13 08:18:25 : Done 15000 of 24863
## 2022-02-13 08:18:29 : Done 16000 of 24863
## 2022-02-13 08:18:32 : Done 17000 of 24863
## 2022-02-13 08:18:35 : Done 18000 of 24863
## 2022-02-13 08:18:39 : Done 19000 of 24863
## 2022-02-13 08:18:42 : Done 20000 of 24863
## 2022-02-13 08:18:46 : Done 21000 of 24863
## 2022-02-13 08:18:49 : Done 22000 of 24863
## 2022-02-13 08:18:53 : Done 23000 of 24863
## 2022-02-13 08:18:56 : Done 24000 of 24863
##
## Inferring gene and species trees
## -----
##
## Best outgroup(s) for species tree
## -----
## 2022-02-13 08:19:00 : Starting STRIDE
## 2022-02-13 08:19:05 : Done STRIDE
## Observed 1067 well-supported, non-terminal duplications. 466 support the best root and 601 contradic
## Best outgroup for species tree:
##   Gtomentosum
##
## Reconciling gene trees and species tree
## -----
## Outgroup: Gtomentosum
## 2022-02-13 08:19:05 : Starting Recon and orthologues
## 2022-02-13 08:19:05 : Starting OF Orthologues
## 2022-02-13 08:19:06 : Done 0 of 24863
## 2022-02-13 08:19:11 : Done 1000 of 24863
## 2022-02-13 08:19:14 : Done 2000 of 24863
## 2022-02-13 08:19:18 : Done 3000 of 24863
## 2022-02-13 08:19:20 : Done 4000 of 24863
## 2022-02-13 08:19:23 : Done 5000 of 24863
## 2022-02-13 08:19:26 : Done 6000 of 24863
## 2022-02-13 08:19:29 : Done 7000 of 24863
## 2022-02-13 08:19:32 : Done 8000 of 24863
## 2022-02-13 08:19:35 : Done 9000 of 24863
## 2022-02-13 08:19:38 : Done 10000 of 24863

```

```
## 2022-02-13 08:19:41 : Done 11000 of 24863
## 2022-02-13 08:19:44 : Done 12000 of 24863
## 2022-02-13 08:19:47 : Done 13000 of 24863
## 2022-02-13 08:19:50 : Done 14000 of 24863
## 2022-02-13 08:19:53 : Done 15000 of 24863
## 2022-02-13 08:19:56 : Done 16000 of 24863
## 2022-02-13 08:19:59 : Done 17000 of 24863
## 2022-02-13 08:20:02 : Done 18000 of 24863
## 2022-02-13 08:20:05 : Done 19000 of 24863
## 2022-02-13 08:20:08 : Done 20000 of 24863
## 2022-02-13 08:20:10 : Done 21000 of 24863
## 2022-02-13 08:20:13 : Done 22000 of 24863
## 2022-02-13 08:20:15 : Done 23000 of 24863
## 2022-02-13 08:20:17 : Done 24000 of 24863
## 2022-02-13 08:20:19 : Done 0F Orthologues
##
## Writing results files
## =====
## 2022-02-13 08:20:36 : Done orthologues
##
## Results:
##   /Users/jlovell/Desktop/GENESPACE_data/results/cotton4x/orthofinder/Results_Feb13/
##
## CITATION:
##   When publishing work that uses OrthoFinder please cite:
##   Emms D.M. & Kelly S. (2019), Genome Biology 20:238
##
##   If you use the species tree in your work then please also cite:
##   Emms D.M. & Kelly S. (2017), MBE 34(12): 3267-3278
##   Emms D.M. & Kelly S. (2018), bioRxiv https://doi.org/10.1101/267914
```

Also set the synteny parameters as default.

```
gparCotton4x <- set_syntenyParams(
  gsParam = gparCotton4x,
  blkSize = 10,
  nGaps = 10)
```

6 Build synteny data

```
gparCotton4x <- find_orthofinderResults(gsParam = gparCotton4x)
gparCotton4x <- synteny(gsParam = gparCotton4x, overwrite = T)
```

```
## Parsing the gff files ...
## Reading the gffs and adding orthofinder IDs ... Done!
## Found 61089 global OGs for 231202 genes
## QC-ing genome to ensure chromosomes/scaffolds are big enough...
##   Genome: n. chrs PASS/FAIL, n. genes PASS/FAIL, n. OGs PASS/FAIL
##   Gbarbadense: 102/775, 72970/1591, 71090/1437
##   Gdarwinii: 41/129, 77967/336, 75457/255
##   Gtomentosum: 37/90, 78145/193, 75367/146
## All look good!
## Defining collinear orthogroup arrays ...
```

```

## Found the following counts of arrays / genome:
##   Gbarbadense: 3031 genes in 1288 collinear arrays
##   Gdarwinii: 3636 genes in 1503 collinear arrays
##   Gtomentosum: 4010 genes in 1685 collinear arrays
## Pulling synteny for 6 unique pairwise combinations of genomes
## Running 1 chunks of 6 combinations each:
## Chunk 1 / 1 (08:20:55) ... Done!
## Gtomento-Gtomento: 1604647 (tot), 257083/178 (reg), 143124/410 (blk)
## Gtomento-Gdarwini: 1603797 (tot), 249751/105 (reg), 139383/360 (blk)
## Gdarwini-Gdarwini: 1616753 (tot), 257928/226 (reg), 143910/420 (blk)
## Gtomento-Gbarbade: 1566227 (tot), 234204/102 (reg), 134192/376 (blk)
## Gdarwini-Gbarbade: 1572604 (tot), 236164/102 (reg), 136647/306 (blk)
## Gbarbade-Gbarbade: 1546417 (tot), 233015/933 (reg), 136510/1129 (blk)
## Defining synteny-constrained orthogroups ...
## Found 72353 synteny-split OGs for 231202 genes
## Running orthofinder by region ...
## genome combinat. : n. non-self genes, nOGs global/syntenic/inblk
## Gbarbade-Gbarbade: 60668 genes, 41568 / 42176 / 32198
## Gdarwini-Gbarbade: 138947 genes, 52628 / 54721 / 42068
## Gdarwini-Gdarwini: 62983 genes, 43360 / 44000 / 33525
## Gtomento-Gbarbade: 136383 genes, 52536 / 54568 / 41241
## Gtomento-Gdarwini: 141780 genes, 54286 / 56862 / 43249
## Gtomento-Gtomento: 63275 genes, 42366 / 43091 / 33603
## Combining synteny-constrained and inblock orthogroups ...
## syn OGs: 72353, inblk OGs: 70893, combined OGs: 56066
## Found the following counts of arrays / genome:
##   Gbarbadense: 4629 genes in 1943 collinear arrays
##   Gdarwinii: 5556 genes in 2272 collinear arrays
##   Gtomentosum: 5910 genes in 2427 collinear arrays
## Pulling synteny for 6 unique pairwise combinations of genomes
## Running 1 chunks of 6 combinations each:
## Chunk 1 / 1 (08:46:53) ... Done!
## Gtomento-Gtomento: 1604647 (tot), 258027/177 (reg), 139802/414 (blk)
## Gtomento-Gdarwini: 1603797 (tot), 250591/99 (reg), 131865/359 (blk)
## Gdarwini-Gdarwini: 1616753 (tot), 258824/222 (reg), 140474/425 (blk)
## Gtomento-Gbarbade: 1566227 (tot), 235017/99 (reg), 127044/390 (blk)
## Gdarwini-Gbarbade: 1572604 (tot), 236567/101 (reg), 129750/319 (blk)
## Gbarbade-Gbarbade: 1546417 (tot), 233749/933 (reg), 133810/1147 (blk)
## Found 56066 OGs across 231202 genes. gff3-like text file written to:
## /Users/jlovell/Desktop/GENESPACE_data/results//cotton4x/results/gffWithOgs.txt.gz
## Calculating syntenic block breakpoints ...
## Found 4934 blocks. Text file written to:
## /Users/jlovell/Desktop/GENESPACE_data/results//cotton4x/results/syntenicBlocks.txt.gz:

```

7 Print session info

```

save(gparCotton4x, file = file.path(baseDir, "gparCotton4x.rda"))
sessionInfo()

```

```

## R version 4.1.2 (2021-11-01)
## Platform: x86_64-apple-darwin17.0 (64-bit)
## Running under: macOS Big Sur 10.16

```

```

##
## Matrix products: default
## BLAS: /Library/Frameworks/R.framework/Versions/4.1/Resources/lib/libRblas.0.dylib
## LAPACK: /Library/Frameworks/R.framework/Versions/4.1/Resources/lib/libRlapack.dylib
##
## locale:
## [1] en_US.UTF-8/en_US.UTF-8/en_US.UTF-8/C/en_US.UTF-8/en_US.UTF-8
##
## attached base packages:
## [1] stats      graphics  grDevices  utils      datasets  methods   base
##
## other attached packages:
## [1] GENESPACE_0.9.3  data.table_1.14.2
##
## loaded via a namespace (and not attached):
## [1] Rcpp_1.0.8                compiler_4.1.2
## [3] restfulr_0.0.13          GenomeInfoDb_1.30.1
## [5] XVector_0.34.0           MatrixGenerics_1.6.0
## [7] R.methodsS3_1.8.1        bitops_1.0-7
## [9] R.utils_2.11.0           tools_4.1.2
## [11] zlibbioc_1.40.0          digest_0.6.29
## [13] lattice_0.20-45          evaluate_0.14
## [15] pkgconfig_2.0.3          rlang_1.0.0
## [17] Matrix_1.4-0             igraph_1.2.11
## [19] DelayedArray_0.20.0      cli_3.1.1
## [21] rstudioapi_0.13          yaml_2.2.2
## [23] parallel_4.1.2           xfun_0.29
## [25] fastmap_1.1.0            GenomeInfoDbData_1.2.7
## [27] rtracklayer_1.54.0       stringr_1.4.0
## [29] knitr_1.37               Biostrings_2.62.0
## [31] S4Vectors_0.32.3         IRanges_2.28.0
## [33] grid_4.1.2              stats4_4.1.2
## [35] Biobase_2.54.0           BiocParallel_1.28.3
## [37] XML_3.99-0.8             rmarkdown_2.11
## [39] magrittr_2.0.2           matrixStats_0.61.0
## [41] GenomicAlignments_1.30.0 Rsamtools_2.10.0
## [43] GenomicRanges_1.46.1    htmltools_0.5.2
## [45] BiocGenerics_0.40.0      SummarizedExperiment_1.24.0
## [47] stringi_1.7.6            RCurl_1.98-1.5
## [49] rjson_0.2.21             crayon_1.4.2
## [51] dbscan_1.1-10           BiocIO_1.4.0
## [53] R.oo_1.24.0

```