

```

1  import os.path
2
3  from glob import glob
4  import atexit
5  from pprint import pprint
6
7  from tools import *
8  from techs import *
9
10
11 ##########
12 ## SETUP
13 #########
14 dataDir = setup.dataDir
15 outDir = setup.outDir
16 print "OUT DIR",outDir
17 print "DATA DIR",dataDir
18
19 #########
20 ## CHECKING
21 #########
22 env      = Environment()
23
24 if not os.path.exists(outDir):
25     print "CRETING OUTPUT DIR"
26     Mkdir(outDir)
27     print "OUT DIR",outDir
28     Decider('timestamp-newer')
29
30
31 #########
32 ## FUNCTIONS
33 #########
34 #########
35 ##fastq2jf_obj = Builder( action      = 'jellyfish $setup.JELLYPARAMS --output $TARGET $SOURCE',
36 ##fastq2jf_obj = Builder( action      = runners.fastq2jf,
37 ##                           suffix     = '.jf',
38 ##                           src_suffix = '.fastq')
39
40 ##
41 ##env['BUILDERS']['fastq2jf_B'] = fastq2jf_obj
42
43 ##F5_Illumina_GOG20_matepair_2000_110401_SN365_A_s_4_2_seq_GOG-20.RD30.NotEmpty.NotLink.fastq.histo
44 ##F5_Illumina_GOG20_matepair_2000_110401_SN365_A_s_4_2_seq_GOG-20.RD30.NotEmpty.NotLink.fastq.histo.desc
45 ##F5_Illumina_GOG20_matepair_2000_110401_SN365_A_s_4_2_seq_GOG-20.RD30.NotEmpty.NotLink.fastq.histo.png
46 ##F5_Illumina_GOG20_matepair_2000_110401_SN365_A_s_4_2_seq_GOG-20.RD30.NotEmpty.NotLink.fastq.jf
47 ##F5_Illumina_GOG20_matepair_2000_110401_SN365_A_s_4_2_seq_GOG-20.RD30.NotEmpty.NotLink.fastq.nfo
48 ##F5_Illumina_GOG20_matepair_2000_110401_SN365_A_s_4_2_seq_GOG-20.RD30.NotEmpty.NotLink.fastq.stats
49 ##F5_Illumina_GOG20_matepair_2000.histo
50 ##F5_Illumina_GOG20_matepair_2000.histo.desc
51 ##F5_Illumina_GOG20_matepair_2000.histo.png
52 ##F5_Illumina_GOG20_matepair_2000.jf
53 ##F5_Illumina_GOG20_matepair_2000.stats
54 ##F5_Illumina.histo
55 ##F5_Illumina.histo.desc
56 ##F5_Illumina.histo.png
57 ##F5_Illumina.jf
58 ##F5_Illumina.stats
59
60
61
62 def checkFq(env, infiles, outfiles=[], baseDir=outDir, dataset=None, library=None, pair=None, object=None):
63     for file in infiles:
64         bn = os.path.basename(str(file))
65
66         if library is not None:
67             bn = library + '_' + bn
68
69         if dataset is not None:
70             bn = dataset + '_' + bn
71
72         bn = os.path.join(baseDir, bn)
73
74         nfo    = File(bn + ".nfo")
75         nfoOut = env.Command( [nfo], [file], runners.makeNfo )
76         Precious( nfoOut )
77         NoClean( nfoOut )
78
79         outJf   = bn + ".jf"
80         outStat = bn + ".stats"
81         outHisto= bn + ".histo"
82         outPng  = bn + ".histo.png"
83         outDesc = bn + ".histo.desc"
84         okFile  = bn + ".ok"
85
86         ##fastqc  = bn +
87         ##LISTCMDQCF[COQCF]
88         ##CMD="$QCP -o '$OUTQCDS' $fastq";
89
90         ##LISTCMDQCS[COQCS]
91
92         ##histoProjectDoQcDoFastqc "\$PROJFOLDER\" \"\$TECHFOLDER\" \"\$LIBFOLDER\" \"\$LOGFILE\" \"\$fastq
93         \
94         \
95         ##CMD="$QCP -o '$OUTQCDS' $fastq";
96
97         ##LISTCMDQCS[COQCS]
98
99         ##histoProjectDoQcDoSolexaqa "\$PROJFOLDER\" \"\$TECHFOLDER\" \"\$LIBFOLDER\" \"\$LOGFILE\" \"\$fas
tq\""
= "./histoProjectDoQcDoFastqc "\$PROJFOLDER\" \"\$TECHFOLDER\" \"\$LIBFOLDER\" \"\$LOGFILE\" \"\$fastq
\
"
= "./histoProjectDoQcDoSolexaqa "\$PROJFOLDER\" \"\$TECHFOLDER\" \"\$LIBFOLDER\" \"\$LOGFILE\" \"\$fas
tq\""

```

```

91      ##CMD="$SOLEXAQA -s $SOLEXAQA_SAMPLESIZE $JELLYIN -o $OUTQCDS";
92
93      if object is not None:
94          if 'addOutput' in dir(object) and callable(getattr(object, 'addOutput')):
95              object.addOutput('jellyDb' , outJf)
96              object.addOutput('jellyStats' , outStat)
97              object.addOutput('jellyHisto' , outHisto)
98              object.addOutput('jellyPng' , outPng)
99              object.addOutput('jellyOk' , okFile)
100             object.addOutput('desc' , outDesc)
101             object.addOutput('nfo' , str(nfo))
102
103
104     JFout = env.Command([outJf, outStat, outHisto ], [file], runners.fastq2jf)
105     ImgOut = env.Command([outPng, outDesc], [outStat, outHisto], runners.histo2png)
106     Depends( ImgOut, nfoOut )
107     Depends( JFout, nfoOut )
108     Depends( ImgOut, JFout )
109
110     all = []
111     all.append( file )
112     all.extend( nfoOut )
113     all.extend( JFout )
114     all.extend( ImgOut )
115
116     ok = env.Command([ okFile ], all, runners.checkOk)
117     outfiles = [ outJf, okFile ]
118
119     print "COUNT KMERS ON FASTQ FILES\n\t", "\n\t".join(map(lambda x: str(x), infilenames)),'\\nto\\n\\t','\\n\\t'.join
120     (map(lambda x:
121         str(x), outfiles)),'\non\\n\\t',baseDir
122     print "OUT FQ\\n\\t","\\n\\t".join(map(lambda x: str(x), outfiles)), "\\n\\n"
123     return outfiles
124
125
126     def joinjf(env, infilenames, outfiles=
127     [], baseDir=outDir, dataset=None, library=None, pair=None, object=None, name=""):
128         file = ""
129         print "MERGING JF FILES [",name,"]\\n\\t", "\\n\\t".join(map(lambda x: str(x), infilenames)),'\\nto\\n\\t','\\n\\t'.join
130         (map(lambda x
131             : str(x), outfiles)),'\non\\n\\t',baseDir
132
133         if len(outfiles) == 0:
134             ##infilenamesShort = [ os.path.basename(str(x)) for x in infilenames ]
135             ##print "IN FILES SHORT",str(infilenamesShort)
136             ##file = os.path.commonprefix(infilenamesShort)
137             ##print " COMMON",file
138             ##file = file.rstrip("._")
139             ##if file == "": file = "mergedRuns"
140             ##exit(1)
141
142         else:
143             file = outfiles[0]
144             if file[-3:] == ".jf":
145                 file = file[:-3]
146             else:
147                 pass
148
149             ##infilenames = [file]
150             ##print "INFILES",map(lambda x: str(x), infilenames),FILE,file
151
152             if os.path.dirname(os.path.abspath(file)) != baseDir:
153                 bn = os.path.basename(file)
154                 if library is not None:
155                     bn = library + '_' + bn
156                 if dataset is not None:
157                     bn = dataset + '_' + bn
158                 file = os.path.join(baseDir, bn)
159
160             outJf = file + ".jf"
161             outStat = file + ".stats"
162             outHisto = file + ".histo"
163             outPng = file + ".histo.png"
164             outDesc = file + ".histo.desc"
165             okFile = file + ".ok"
166
167             if object is not None:
168                 if 'addOutput' in dir(object) and callable(getattr(object, 'addOutput')):
169                     object.addOutput('jellyDb' , outJf)
170                     object.addOutput('jellyStats' , outStat)
171                     object.addOutput('jellyHisto' , outHisto)
172                     object.addOutput('jellyHistoDesc' , outDesc)
173                     object.addOutput('jellyPng' , outPng)
174                     object.addOutput('jellyOk' , okFile)
175
176             JFout = env.Command([outJf, outStat, outHisto], [x for x in infilenames if str(x).endswith
177 ('.jf')], runners.j
178             joinjf)
179             ImgOut = env.Command([outPng, outDesc], [outStat, outHisto], runners.histo2png)
180             Depends( JFout, infilenames )
181             ##Depends( ImgOut, JFout )
182
183             all = []

```

```

179     all.append( infiles )
180     all.extend( JFout )
181     all.extend( ImgOut )
182     ok = env.Command([ okFile ], all, runners.checkOk)
183     ##print "OK",okFile
184     outfiles = [ outJf , okFile ]
185
186     print "OUT JF\n\t", "\n\t".join(map(lambda x: str(x), outfiles)), "\n\n"
187     return outfiles
188 ##env.AddMethod(fastq2jf_def, "fastq2jf")
189 AddMethod(Environment, joinjf)
190
191
192
193 def fastqc(env, infiles, outfiles=[], baseDir=outDir, dataset=None, library=None, pair=None, object=None):
194     print "RUNNING FASTQC\n\t", "\n\t".join(map(lambda x: str(x), infiles)), '\n\t', '\n\t'.join(map(
195     lambda x: str(x), out
196     files)), '\n\t', baseDir
197
198     if len(outfiles) == 0:
199         exit(1)
200
201     file = outfiles[0]
202     if os.path.dirname(os.path.abspath(file)) != baseDir:
203         bn = os.path.basename(file)
204         if library is not None:
205             bn = library + '_' + bn
206         if dataset is not None:
207             bn = dataset + '_' + bn
208         file = os.path.join(baseDir, bn)
209
210     outZip = file + ".zip"
211     outPath = file + "_fastqc"
212     outData = os.path.join(outPath, 'fastqc_data.txt')
213     outReport = os.path.join(outPath, 'fastqc_report.html')
214     outSummary = os.path.join(outPath, 'summary.txt')
215
216     okFile = file + ".fqc.ok"
217
218     if object is not None:
219         if 'addOutput' in dir(object) and callable(getattr(object, 'addOutput')):
220             ##TODO: ADD INDIVIDUAL FILES INSIDE DIR ?
221             object.addOutput('fastQCZip', outZip)
222             object.addOutput('fastQCDir', outPath)
223             object.addOutput('fastQCData', outData)
224             object.addOutput('fastQCReport', outReport)
225             object.addOutput('fastQCSummary', outSummary)
226             object.addOutput('fastQCOK', okFile)
227
228     #print "\tFASTQC\n\t", outZip, "\n\t", outDir
229 FQCout = env.Command([outZip, outReport, outSummary], infiles, runners.fastqc)
230 #print "\n\tOUT\n\t", FQCout
231 Clean(FQCout, outPath)
232
233 all = []
234 all.append( infiles )
235 all.extend( FQCout )
236 ok = env.Command([ okFile ], all, runners.checkOk)
237     ##print "OK",okFile
238     outfiles = ok
239
240     print "OUT FQC\n\t", "\n\t".join(map(lambda x: str(x), outfiles)), "\n\n"
241     return outfiles
242 ##env.AddMethod(fastq2jf_def, "fastq2jf")
243 AddMethod(Environment, fastqc)
244
245
246 def solexaqa(env, infiles, outfiles=[], baseDir=outDir, dataset=None, library=None, pair=None, object=None):
247     print "RUNNING SOLEXAQA\n\t", "\n\t".join(map(lambda x: str(x), infiles)), '\n\t', '\n\t'.join(map(
248     lambda x: str(x), out
249     files)), '\n\t', baseDir
250
251     if len(outfiles) == 0:
252         exit(1)
253
254     file = outfiles[0]
255     bn = os.path.basename(file)
256     if os.path.dirname(os.path.abspath(file)) != baseDir:
257         bn = os.path.basename(file)
258         if library is not None:
259             bn = library + '_' + bn
260         if dataset is not None:
261             bn = dataset + '_' + bn
262         file = os.path.join(baseDir, bn)
263
264     #QUALP="$OUTQCDS/$fastqName$fastqExt.quality.pdf"
265     #QUALI="$OUTQCDS/$fastqName$fastqExt.quality.pdf.png"
266     #HISP="$OUTQCDS/$fastqName$fastqExt.segments.hist.pdf"
267     #HISI="$OUTQCDS/$fastqName$fastqExt.segments.hist.pdf.png"
268
269     outPath = file + '_solexaqa'
270     outMtx = os.path.join(outPath, bn) + ".matrix"
271     outQualPdf = os.path.join(outPath, bn) + ".quality.pdf"
272     outHistPdf = os.path.join(outPath, bn) + ".segments.hist.pdf"

```

```

270     outMtxPng = outMtx + ".png"
271     outQualPng = outQualPdf + '.png'
272     outHistPng = outHistPdf + '.png'
273     okFile = file + ".sqa.ok"
274
275     if object is not None:
276         if 'addOutput' in dir(object) and callable(getattr(object, 'addOutput')):
277             ##TODO: ADD INDIVIDUAL FILES INSIDE DIR ?
278             object.addOutput('solexaqaMatrix', outMtx)
279             object.addOutput('solexaqaMatrixPng', outMtxPng)
280             object.addOutput('solexaqaQualPng', outQualPng)
281             object.addOutput('solexaqaHistPdf', outHistPdf)
282             object.addOutput('solexaqaHistPng', outHistPng)
283             object.addOutput('solexaqaDir', outPath)
284             object.addOutput('solexaqaOk', okFile)
285
286     #print "\tFASTQC\n\t\t",outZip,"\\n\t\t",outDir
288     SQAout = env.Command
289     #[outMtx, outMtxPng, outQualPng, outHistPdf, outHistPng], infilenames, runners.solexaqa)
290     #print "\\n\tOUT\\n\\t\\t",FQCout
291     Clean(SQAout, outPath)
292
293     all = []
294     all.append( infilenames )
295     all.extend( SQAout )
296     ok = env.Command([ okFile ], all, runners.checkOk)
297     ##print "OK",okFile
298     outfiles = ok
299
300     print "OUT FQC\\n\\t", "\\n\\t".join(map(lambda x: str(x), outfiles)), "\\n\\n"
301     return outfiles
302     ##env.AddMethod(fastq2jf_def, "fastq2jf")
303     AddMethod(Environment, solexaqa)
304
305     ##### GLOBAL TARGETS #####
306     #####
307     def getTargets():
308         for dataset in setup.allSpecies:
309             datasetName = dataset.getName()
310
311             ##path = os.path.abspath(os.path.join(dataDir,datasetName))
312             ##bn = os.path.basename(path)
313             ##if not os.path.exists(path) or not os.path.isdir(path):
314             ##    print "DATASET",datasetName,"DOES NOT EXISTS ON",dataDir
315             ##    continue
316             sppok = os.path.join(os.path.join(outDir, datasetName), datasetName+'.ok')
317             ##print "tgt",datasetName,>,sppok
318             yield [datasetName, sppok]
319
320     def checkTargets(spp):
321         if spp not in valid_targets:
322             print " Invalid species '",spp,"'"
323             print " Valid species are:"
324             print "\\t" + "\\n\\t".join(valid_targets)
325             exit(1)
326
327     ##yield [None, 'all', os.path.join(outDir,'all')+'.ok']
328
329
330     allName = setup.allSpecies.getName()
331     allAlias = []
332     for spp in getTargets():
333         print 'ADDING ALIAS',spp[0],'TO FILE',spp[1]
334         env.Alias(spp[0], spp[1])
335         allAlias.append(env.Alias(spp[0]))
336         ##print "FINISHED " + str( finished )
337     ##
338     ##
339     Default(env.Alias(allName))
340     env.Alias( allName, [ os.path.join(outDir,allName+'.ok') ] )
341
342     valid_targets = map(lambda x: x[0], getTargets())
343
344     ##env['DISTDIR'] = setup.outDir
345
346
347
348
349
350     ##### CALLERS #####
351     ## CALLERS
352     #####
353     def doSpecies(datasetName, dataset):
354         print "      RUNNING DATASET '" + datasetName + "'"
355         baseDir = os.path.join(outDir,datasetName)
356         ##str(dataset)
357
358         print "      DATASET",dataset.getName()," ",dataset
359
360         datasetOut = []
361         for key in dataset.getIlluminaNames():
362             print "          ILLUMINA DATASET NAME",key
363             illDataset = dataset.getIllumina(key)

```

```

364         #print illDataset
365
366         libs = []
367         for lib in illDataset:
368             libName           = lib.getName()
369             print '           LIB ' + libName
370
371         pairs = []
372         for pair in lib:
373             pairName   = pair.getName()
374             print '           PAIR ' + pairName + ' TYPE ' + pair.getType()
375
376         runFiles  = []
377         pairFiles = []
378         for run in pair:
379             ##print '           RUN ' + run.getShortName() + ' FN ' + run.getFileName()
380             pairFiles.append(run.getFileName())
381         runFiles.extend(env.checkFq([run.getFileName
()]), [], baseDir=baseDir, dataset=datasetName, library=libNa
me, pair=pairName, object=run))
382             ##iffinished = env.Command(env.Alias(spp), ends, runners.checkOk)
383
384
385             if ( setup.runFastqc ): datasetOut.extend(env.fastqc
( pairFiles, [ pairName ], baseDir=baseDir, dataset=
datasetName, library=libName, pair=pairName, object=pair))
386                 if ( setup.runSolexaqa ): datasetOut.extend(env.solexaqa
(pairFiles, [ pairName ], baseDir=baseDir, dataset=
datasetName, library=libName, pair=pairName, object=pair))
387                 if ( setup.runJellyfish ): pairs.extend(   env.joinjf
( runFiles, [ pairName ], baseDir=baseDir, dataset=
datasetName, library=libName, pair=pairName, object=pair, name="PAIR_" + datasetName + "_" + libName + "_" + pai
)
388
389                 ##print "PAIRS",map(lambda x: str(x), pairs)
390                 if ( setup.runJellyfish ): libs.extend(env.joinjf
(pairs, [ libName ], baseDir=baseDir, dataset=datasetName, obje
ct=lib, name="LIB_" + datasetName + "_" + libName))
391
392                 ##print "LIBS",map(lambda x: str(x), libs)
393                 if ( setup.runJellyfish ): datasetOut.extend(env.joinjf
(libs, [ datasetName ], baseDir=baseDir, object=illDataset, n
ame="DATASET_" + datasetName))
394
395
396
397
398
399         for key in dataset.getRocheNames():
400             print "           ROCHE DATASET NAME",key
401             rocheDataset = dataset.getRoche(key)
402             #print rocheDataset
403             ##datasetOut.extend(env.joinjf(libs, [ datasetName ], baseDir=baseDir, object=dataset))
404
405         for key in dataset.getPacBioNames():
406             print "           PACBIO DATASET NAME",key
407             pacbioDataset = dataset.getPacBio(key)
408             #print pacbioDataset
409             ##datasetOut.extend(env.joinjf(libs, [ datasetName ], baseDir=baseDir, object=dataset))
410
411         sppOk = env.Command(env.Alias(spp), datasetOut, runners.checkOk)
412         print 'SPP OK ALIAS',sppOk,'FROM\n\t'+"\n\t".join(map(lambda x: str(x), datasetOut))
413         return sppOk
414
415         #return env.joinjf(libs, [datasetName])
416
417
418
419
420
421 #####
422 ## MAIN
423 #####
424 print "BUILD TARGET",map(str, BUILD_TARGETS)
425 if allName in [str(x).strip() for x in BUILD_TARGETS ]:
426     print "SPP "+allName+""
427
428
429         envs = []
430         for spp in valid_targets:
431             print "           SUB SPP '",spp,"'"
432             dataset  = setup.allSpecies.get(spp)
433             print "           DATASET '",dataset.getName(),"'"
434             ends    = doSpecies(spp, dataset)
435             ##iffinished = env.Command(env.Alias(spp), ends, runners.checkOk)
436             env.Alias(allName, env.Alias(spp))
437             envs.extend(ends)
438
439         print "ENVS", str([str(x) for x in envs])
440         ##print "ALIAS",str([str(x) for x in env.Alias('all')]) )
441         finished = env.Command([ os.path.join(outDir,allName+'.ok') ], [os.path.join(os.path.join(outDir, str
(x)), str(x)+".ok")]
442             for x in envs], runners.checkOk)
443             ##Depends(env.Alias('all'), finished)
444

```

```
444     else:
445         for spp in [str(x).strip() for x in BUILD_TARGETS]:
446             if spp == allName: continue
447
448             checkTargets(spp)
449             print "  Valid species", spp
450             dataset = setup.allSpecies.get(spp)
451             ends = doSpecies(spp, dataset)
452             ##finished = env.Command(env.Alias(spp), ends, runners.checkOk)
453
454             print "ENDS " + str(map(lambda x: str(x), ends))
455             ##env.Depends(env.Alias(spp), ends)
456
457
458
459     setup.exporter()
460     ##exit(0)
461
462
463
464
465
466
467
468
469
470
471
472 # Make the build fail if we pass fail=1 on the command line
473 if ARGUMENTS.get('fail', 0):
474     Command('target', 'source', ['/bin/false'])
475
476 def bf_to_str(bf):
477     """Convert an element of GetBuildFailures() to a string
478     in a useful way."""
479     import SCons.Errors
480     if bf is None: # unknown targets product None in list
481         return '(unknown tgt)'
482     elif isinstance(bf, SCons.Errors.StopError):
483         return str(bf)
484     elif bf.node:
485         return str(bf.node) + ': ' + bf.errstr
486     (bf.executor) + "\nACTION: "
487     " + str(bf.action)
488     elif bf.filename:
489         return bf.filename + ': ' + bf.errstr + "\nCOMMAND: " + bf.command + "\nEXECUTOR: " + str
490     (bf.executor) + "\nACTION:
491     " + str(bf.action)
492
493     return 'unknown failure: ' + bf.errstr
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525 Help("""
526 Type: 'run' to build the genome
527 """)
```