

NanOlympicsMod – tutorial: how to add a new tool (e. g. m6Anet)

m6Anet, as many other tools, requires one step to run on each sample independently, and another step to combine results obtained for all the samples. For this reason, we are going to add two processes, named m6anet1 and m6anet2. In the following, we will show the steps required for including the tool m6Anet as part of the NanOlympicsMod pipeline.

1. In **pipeline.conf** file:

- a. Add one flag variable for each new process that you want to add, and set their status to true; in particular, we are going to add m6anet1 and m6anet2 variables (lines 41, 42). In case we would like to skip the processes execution, we will set these variables to false.
- b. Add options specific for each process; in particular, add the name of the image which should be downloaded from Docker Hub, mounting options and desired computational resources (lines 57-72).

2. In **pipeline.nf** file:

- a. Add channels to link processes, such as one channel from nanopolish1 to m6anet1 processes and one channel from m6anet2 to postprocessing processes (lines 101, 208)
- b. Add m6anet1 process, which should run on each sample independently (lines 128-150)
- c. Group by condition output from m6anet1 process and create a new channel for each condition (lines 152-156)
- d. Add m6anet2 process, which should analyse all samples assigned to test condition, after m6anet1 process is over for all the samples (lines 158-179)
- e. Add the path to m6Anet output as an argument to postprocessing.R script in postprocessing process (line 224)

3. In **postprocessing.R** file:

- a. Add a case to the switch code block for formatting m6Anet output into a genome-based bed file format. The final output should be a tab-separated file including chromosome, start, end, status (Mod/Unmod) and confidence parameter value. Since m6Anet works on the transcriptome, an additional code block for performing lift-over from transcriptome-based to genome-based coordinates is needed (lines 246-340).
- b. Add the name of the added tool and default/relaxed confidence parameter values (lines 345-352)

4. In **statistical_analysis.R** file:

- a. Add the tool to either listmax (as in case of m6Anet) or listmin vectors, depending on whether the confidence parameter should be maximised or minimised to obtain higher confidence calls, respectively (line 369-376)
- b. Add a line for renaming the tool with the correct case-sensitive spelling, useful for plot legends (line 393)
- c. Add the tool to the list of tools with a confidence parameter, for which PR curves should be plotted (lines 403, 409)

34 **pipeline.conf**

```
35
36 [...]
37     nanopolish1 = true
38     xpore = true
39     nanocompore1 = true
40     nanocompore2 = true
41     m6anet1 = true
42     m6anet2 = true
43     yanocomp1 = true
44     yanocomp2 = true
45     postprocessing = true
46 }
47
48 [...]
49     withName:xpore2{
50         container = 'bproject/xpore:v1'
51         containerOptions = '--bind /path/to/mount:/path/mounted'
52         cpus = { params.xpore ? 3 : 1 }
53         memory = { params.xpore ? 5.GB + (2.GB * (task.attempt-1)) : 1.GB }
54         errorStrategy = { task.exitStatus == 130 ? 'retry' : 'terminate' }
55         maxRetries = 3
56     }
57     withName:m6anet1{
58         container = 'bproject/m6anet:v1'
59         containerOptions = '--bind /path/to/mount:/path/mounted'
60         cpus = { params.m6anet1 ? 3 : 1 }
61         memory = { params.m6anet1 ? 5.GB + (2.GB * (task.attempt-1)) : 1.GB }
62         errorStrategy = { task.exitStatus == 130 ? 'retry' : 'terminate' }
63         maxRetries = 3
64     }
65     withName:m6anet2{
66         container = 'bproject/m6anet:v1'
67         containerOptions = '--bind /path/to/mount:/path/mounted'
68         cpus = { params.m6anet2 ? 3 : 1 }
69         memory = { params.m6anet2 ? 5.GB + (2.GB * (task.attempt-1)) : 1.GB }
70         errorStrategy = { task.exitStatus == 130 ? 'retry' : 'terminate' }
71         maxRetries = 3
```

```

72     }
73     withName:nanocompore1{
74         container = 'bproject/nanocompore:v1'
75         containerOptions = '--bind /path/to/mount:/path/mounted'
76         cpus = { params.nanocompore1 ? 7 : 1 }
77         memory = { params.nanocompore1 ? 10.GB + (2.GB * (task.attempt-1)) : 1.GB }
78         errorStrategy = { task.exitStatus == 130 ? 'retry' : 'terminate' }
79         maxRetries = 3
80     }
81     [...]
82
83     pipeline.nf
84
85     [...]
86
87     // Resquigling with nanopolish for each condition
88     process nanopolish1 {
89         input:
90             tuple val(condition), val(sample) from minimap2_nanopolish1
91
92             each file('transcriptome.fa') from transcriptome_fasta_nanopolish1
93             each file('transcriptome.fa.fai') from transcriptome_fai_nanopolish1
94
95             each file('genome.fa') from genome_fasta_nanopolish1
96             each file('genome.fa.fai') from genome_fai_nanopolish1
97         output:
98             tuple val(condition), val(sample) into nanopolish1_xpore
99             tuple val(condition), val(sample) into nanopolish1_nanocompore1
100            tuple val(condition), val(sample) into nanopolish1_yanocomp1
101            tuple val(condition), val(sample) into nanopolish1_m6anet1
102
103         script:
104         if(params.nanopolish1)
105         """
106             mkdir -p ${params.resultsDir}/${condition}/${sample}/nanopolish/
107             mkdir -p ${params.resultsDir}/${condition}/${sample}/nanopolish/transcriptome/
108
109         [...]

```

```

110
111     nanocompore sampcomp --file_list1 "\${f1[*]}" --file_list2 "\${f2[*]}" \
112         --label1 ${condition1} \
113         --label2 ${condition2} \
114         --fasta transcriptome.fa \
115         --bed transcriptome.bed \
116         --outpath ${params.resultsDir}/nanocompore/ \
117         --allow_warnings \
118         --logit \
119         --nthreads ${task.cpus} \
120         --overwrite
121     """
122     else
123         """
124         echo "Skipped"
125     """
126 }
127
128 // Data formatting for m6anet for each sample
129 process m6anet1 {
130     input:
131         tuple val('condition'), val('sample') from nanopolish1_m6anet1
132
133     output:
134         tuple val(condition), val(sample), val() into m6anet1_m6anet2
135
136     script:
137         if(params.m6anet1)
138             """
139             mkdir -p ${params.resultsDir}/${condition}/${sample}/m6anet/
140
141             m6anet-dataprep --eventalign
142             ${params.resultsDir}/${condition}/${sample}/nanopolish/transcriptome/eventalign_readIndex.txt \
143             --out_dir ${params.resultsDir}/${condition}/${sample}/m6anet --n_processes ${task.cpus}
144             """
145         else
146             """
147

```

```

148         ln -sf ${params.resultsDir}/${condition}/${sample}/m6anet m6anet
149         ""
150     }
151
152     // From a single channel for all the alignments to one channel for each condition
153     ni_test_m6anet2=Channel.create()
154     ni_other_m6anet2=Channel.create()
155     m6anet1_m6anet2.groupTuple(by:0)
156         .choice( ni_test_m6anet2, ni_other_m6anet2 ) { a -> a[0] == params.test_condition ? 0 : 1 }
157
158     // RNA modifications detection with m6anet
159     process m6anet2 {
160         input:
161             tuple val('condition1'), val('sample1') from ni_test_m6anet2
162
163         output:
164             val('flagm6anet') into m6anet_postprocessing
165         script:
166             if(params.m6anet2)
167                 ""
168                 mkdir -p ${params.resultsDir}/m6anet
169                 preprocessing_dirs=$(find ${params.resultsDir}/${condition1} -maxdepth 2 -type d | grep "m6anet\$")
170                 m6anet-run_inference --input_dir \$preprocessing_dirs --out_dir ${params.resultsDir}/m6anet --
171             infer_mod_rate --n_processes ${task.cpus}
172
173             zcat ${params.resultsDir}/m6anet/data.result.csv.gz > ${params.resultsDir}/m6anet/data.result.csv
174             ""
175         else
176             ""
177             echo "Skipped"
178             ""
179     }
180
181     // Data formatting for yanocomp for each sample
182     process yanocomp1 {
183         input:
184             tuple val('condition'), val('sample') from nanopolish1_yanocomp1
185             each file('genome.gtf') from gtf_yanocomp

```

```

186     output:
187         tuple val(condition), file('outputT.hdf5'), file('outputG.hdf5') into yanocomp1_yanocomp2
188
189     [...]
190
191     // Processing of each output to obtain bed files
192     process postprocessing {
193         input:
194             val('flagyanocomp2') from yanocomp2_postprocessing
195             val('flagdena') from dena_postprocessing
196             val('flagdrummer') from drummer_postprocessing
197             val('flagdifferr') from differr_postprocessing
198             val('flagnanom6a') from nanom6a_postprocessing
199             val('flagnanocompore') from nanocompore_postprocessing
200             val('flageligos') from eligos_postprocessing
201             val('flagmines') from mines_postprocessing
202             val('flagepinanoSVM') from epinanoSVM_postprocessing
203             val('flagepinanoError') from epinanoError_postprocessing
204             val('flagxpore') from xpore_postprocessing
205             val('flagnanodoc') from nanodoc_postprocessing
206             val('flagtombo2') from tombo2_postprocessing
207             val('flagtombo3') from tombo3_postprocessing
208             val('flagm6anet') from m6anet_postprocessing
209
210         output:
211
212         script:
213             if(params.postprocessing)
214                 """
215                     mkdir -p ${params.resultsDir}/output_bed_files/
216                     mkdir -p ${params.resultsDir}/output_statistical/
217
218                     Rscript ${params.postprocessingScript} path=${params.resultsDir} genomebed=${params.genesbed}
219                     genomegtf=${params.gtf} resultsFolder=${params.resultsDir}/output_bed_files/ mcores=${task.cpus}
220                     threshold=${params.threshold} pathdena=${params.test_condition}/dena pathdrummer=drummer pathdifferr=differr
221                     pathyanocomp=yanocomp pathmines=${params.test_condition}/mines pathnanocompore=nanocompore
222                     patheligos=eligos/merged pathepinanoError=epinanoError pathepinanoSVM=${params.test_condition}/epinanoSVM

```

```

223 pathxpore=xpore pathnanodoc=nanodoc pathnanom6a=${params.test_condition}/nanom6a/result_final
224 pathtomboComparison=tomboComparison pathm6anet=m6anet
225
226         Rscript ${params.statisticalAnalysis} bed_folder=${params.resultsDir}/output_bed_files
227 genomegtf=${params.gtf} genesbed=${params.genesbed} resultsFolder=${params.resultsDir}/output_statistical/
228 mccores=${task.cpus} peaks=${params.peaksfile} binLength=${params.binLength} genomefile=${params.genome_fasta}
229
230         ""
231         else
232         ""
233         echo "Skipped"
234         ""
235     }
236
237 Postprocessing.R
238
239     [...]
240         else {message(paste0(tool,"'s output files don't exist."))}
241     }
242     if (!file.exists(output_file)) {
243         output_tomboComparison()
244     }
245     },
246     m6anet = {output_m6anet <- function(){if (file.exists(paste0(path_folder, "/", "data.result.csv")))
247                                     && file.info(paste0(path_folder, "/",
248 "data.result.csv"))$size != 0){
249         #read results from new_tool and save relevant columns to m6anet df
250         data_m6anet <- read.table(paste0(path_folder, "/", "data.result.csv"), header = TRUE, sep=",")
251
252         m6anet <- data.frame("TranscriptID" = data_m6anet[,1],
253                             "Start" = data_m6anet[,2] + 2,
254                             "End" = data_m6anet[,2] + 2,
255                             "Status" = data_m6anet[,4],
256                             "Prob_mod" = data_m6anet[,4]
257                             )
258
259         #use m6anet$Status < filtering_parameter instead of m6anet$Status > filtering_parameter if ConfPar
260         parameter needs to be minimised

```

```

261     m6anet$Status <- ifelse(!is.nan(m6anet$Status) & !is.na(m6anet$Status) & m6anet$Status >
262 filtering_parameter, "Mod", "Unmod")
263 #retain only Mod sites
264     m6anet <- m6anet[which(m6anet$Status == "Mod"), ]
265     #Creation Edb Database from genome GTF
266     EnsDb <- suppressWarnings(suppressMessages(ensDbFromGtf(gtf = genome_gtf)))
267     edb <- EnsDb(EnsDb)
268     # Lift-over + output bed
269     test_m6anet <- IRanges(start = m6anet[,2], end = m6anet[,3], names = c(m6anet[,1]))
270     #process num_rows_chunk at a time
271     num_rows_chunk <- 1
272     mc.cores <- as.numeric(mccores)
273     #split the reads in chunks of size num_rows_chunk
274     if (length(test_m6anet) < num_rows_chunk) {
275         test_m6anet_split <- list(test_m6anet)
276     } else {
277         test_m6anet_split <- split(test_m6anet, rep(seq(from = 1, to =
278 ceiling(length(test_m6anet)/num_rows_chunk)), each = num_rows_chunk)[1:length(test_m6anet)])
279     }
280     #initialize vectors
281     tmp1 <- vector(mode = "list", length = length(test_m6anet_split))
282     names(tmp1) <- 1:length(test_m6anet_split)
283     tmp <- vector(mode = "list", length = length(test_m6anet_split))
284     names(tmp) <- 1:length(test_m6anet_split)
285     ind_retry <- 1:length(test_m6anet_split)
286     #keep retrying on chunks which failed unexpectedly
287     while(any(unlist(lapply(tmp, is.null)))) {
288         cat(sprintf("Starting new iteration for m6Anet; %d sites missing\n",
289 length(which(unlist(lapply(tmp, is.null))))))
290         tmp1 <- tmp1[ind_retry]
291         #if lift-over fails, print "Warning" or "Error"
292         tmp1 <- mclapply(test_m6anet_split[ind_retry], function(x) {
293             tryCatch({
294                 coordinate_m6anet_unlisted <- unlist(transcriptToGenome(x, edb))
295                 return(coordinate_m6anet_unlisted)
296             }, warning = function(w) {
297                 print("Warning")
298                 return(NULL)

```



```

299     }, error = function(e) {
300         print("Error")
301         return(NULL)
302     }
303     }}, mc.cores = mc.cores)
304     #retrieve chunks which failed
305     ind_retry <- names(which(unlist(lapply(tmp1, function(x) is.null(x))))))
306     ind_ok <- names(which(unlist(lapply(tmp1, function(x) !is.null(x))))))
307     tmp[ind_ok] <- tmp[ind_ok]
308     if (length(ind_retry) > 0) {
309         tmp1 <- tmp1[ind_retry]
310     }
311 }
312 #unlist the output of lift-over
313 coordinate_m6anet_unlisted <- unlist(as(tmp, "GRangesList"))
314 #convert the lifted-over hits to dataframe
315 df_m6anet <- as.data.frame(unname(coordinate_m6anet_unlisted[,c(0,2,4,5)]))[,c(1:3,5,6,7,8)]
316 #assign rownames and delete duplicate hits (may be due to hits falling on exon-intron junctions)
317 tmp_rownames <- paste0(df_m6anet[, 6], "_", df_m6anet[, 7], "_", df_m6anet[, 5])
318 dup_names <- names(which(table(tmp_rownames) > 1))
319 ind_dup <- which(tmp_rownames %in% dup_names)
320 ind_dup_rm <- ind_dup[which(duplicated(df_m6anet[ind_dup, c(5,6,7)]))]
321 if (length(ind_dup_rm) > 0) {
322     df_m6anet <- df_m6anet[-ind_dup_rm, ]
323 }
324 rownames(df_m6anet) <- paste0(df_m6anet[, 6], "_", df_m6anet[, 7], "_", df_m6anet[, 5])
325 rownames(m6anet) <- paste0(m6anet[, 2], "_", m6anet[, 3], "_", m6anet[, 1])
326 #assign colnames and write to file
327 df_m6anet$Status <- m6anet[rownames(df_m6anet), 4]
328 df_m6anet$Prob_Mod <- m6anet[rownames(df_m6anet), 5]
329 df_m6anet_final <- df_m6anet[,c(1,2,3,4,8,9)]
330 colnames(df_m6anet_final) <- c("Chr", "Start", "End", "Strand", "Status", "Prob_mod")
331 write.table(df_m6anet_final, file = output_file, quote = F, sep = "\t", row.names = F)
332 }
333     else {message(paste0(tool,"'s output files don't exist."))}
334 }
335 if (!file.exists(output_file)) {
336     output_m6anet()

```

```

337         }},
338         stop("Enter a valid tool as input!")
339     )}
340 }
341
342 # Definining a set of parameters used to filter the results for those tools which give as output all the sites
343 rrach <- c("AAACA","AAACT","AAACC","GAACA","GAACT","GAACC","GGACA","GGACT","GGACC","GAACA","GAACT","GAACC")
344 tools <- c("dena", "drummer", "differr", "yanocomp", "nanocompore", "eligos", "mines"
345           , "epinanoErr", "epinanoSvm", "xpore", "nanodoc", "nanom6a", "tomboComparison", "m6anet")
346
347 pathTools <- c(pathdena, pathdrummer, pathdifferr, pathyanocomp, pathnanocompore, patheligos
348               , pathmines, pathepinanoError, pathepinanoSVM, pathxpore, pathnanodoc, pathnanom6a
349               , pathtomboComparison, pathm6anet)
350
351 default <- c(0.1, 0.05, 0.05, 0.05, 0.01, 0.0001, NA, 0.1, 0.5, 0.05, 0.02, NA, 0.05, 0.9)
352 relaxed <- c(0, NA, NA, NA, 1, 1, NA, NA, 0, 1, 0, NA, 1, 0)
353 value <- rep(threshold, length(default))
354
355 names(pathTools) <- names(default) <- names(relaxed) <- names(value) <- tools
356
357 parameters_list <- list("default" = unname(default), "relaxed" = unname(relaxed), "value" = unname(value))
358
359 [...]
360
361 statistical_analysis.R
362
363 [...]
364
365 ## 1) Load all the output files from the output directory
366 files <- list.files(bed_folder, full.names = TRUE, pattern = "\\\\.bed") # output_directory parameter from
367 outside
368
369 listmax <- paste0(c("DENA", "EpiNano-Error", "EpiNano-SVM", "NanoDoc", "m6Anet"), collapse = "|") # for these
370 tools we need to maximize the filtering parameter when there are more than 1 in a bin
371 listmin <- paste0(c("DiffErr", "DRUMMER", "Yanocomp", "Nanocompore", "ELIGOS", "xPore", "Tombo"), collapse =
372 "|") # for these tools we need to minimize the filtering parameter
373
374 threshold_default <- c(0.1, 0.05, 0.05, 0.05, 0.01, 0.0001, 0.1, 0.5, 0.05, 0.02, 0.05, 0.9)

```

```

375 names(threshold_default) <- c("DENA", "DiffErr", "DRUMMER", "Yanocomp", "Nanocompore", "ELIGOS",
376                               "EpiNano-Error", "EpiNano-SVM", "xPore", "NanoDoc", "Tombo", "m6Anet")
377
378 chrs <- readDNAStringSet(genomefile, format="fasta")
379 RRACH_plus <- GRanges(vmatchPattern(pattern = "RRACH", subject = chrs, fixed = "subject"), strand = "+")
380 RRACH_minus <- GRanges(vmatchPattern(pattern = "DGTYY", subject = chrs, fixed = "subject"), strand = "-")
381 RRACH <- c(RRACH_plus, RRACH_minus)
382
383 [...]
384
385 Run_statistical_analysis <- function(genesBins_par, peaks_par, files_par, notes = "", w) {
386   tools <- basename(files_par)
387   tools[grep(pattern = "dena", x = tools)] <- "DENA"
388   tools[grep(pattern = "drummer", x = tools)] <- "DRUMMER"
389   tools[grep(pattern = "differr", x = tools)] <- "DiffErr"
390   tools[grep(pattern = "eligos", x = tools)] <- "ELIGOS"
391   tools[grep(pattern = "epinanoErr", x = tools)] <- "EpiNano-Error"
392   tools[grep(pattern = "epinanoSvm", x = tools)] <- "EpiNano-SVM"
393   tools[grep(pattern = "m6anet", x = tools)] <- "m6Anet"
394   tools[grep(pattern = "mines", x = tools)] <- "MINES"
395   tools[grep(pattern = "nanocompore", x = tools)] <- "Nanocompore"
396
397   [...]
398
399   # Overlap between m6A detected site of each tool and genome binned
400   overlap <- as.matrix(findOverlaps(query = granges, subject = genesBins_par, minoverlap = 1, type = "any"))
401   if (grepl(x, pattern = paste0(c("DENA", "DRUMMER", "DiffErr", "Yanocomp", "Nanocompore", "ELIGOS", "EpiNano-
402 Error",
403                               "EpiNano-SVM", "xPore", "NanoDoc", "Tombo", "m6Anet"), collapse = "|")) {
404     # Add column of filtering parameter
405     filtering_parameter <- bed_file[overlap[, "queryHits"], 6]
406     overlap_w_parameter <- cbind(overlap, filtering_parameter)
407     default <- unname(threshold_default[grep(x, pattern =
408 paste0(c("DENA", "DRUMMER", "DiffErr", "Yanocomp", "Nanocompore", "ELIGOS", "EpiNano-Error", "EpiNano-SVM",
409                               "xPore", "NanoDoc", "Tombo", "m6Anet"),
410 collapse = "|"), value = T)])
411   [...]

```