

7IM-^M-^T 03, 25 13:43

run-polap-r-directional.R

Page 1/9

```

1 #!/usr/bin/env Rscript
2
3 ######
4 # This file is part of polap.
5 #
6 # polap is free software: you can redistribute it and/or modify it under the
7 # terms of the GNU General Public License as published by the Free Software
8 # Foundation, either version 3 of the License, or (at your option) any later
9 # version.
10 #
11 # polap is distributed in the hope that it will be useful, but WITHOUT ANY
12 # WARRANTY; without even the implied warranty of MERCHANTABILITY or FITNESS FOR
13 # A PARTICULAR PURPOSE. See the GNU General Public License for more details.
14 #
15 # You should have received a copy of the GNU General Public License along with
16 # polap. If not, see <https://www.gnu.org/licenses/>.
17 #####
18 #####
19 #####
20 # This script selects reads using some of minimap2 output files.
21 # We have multiple scripts that selects reads.
22 # This is the 5th one.
23 #
24 # 1. run-polap-pairs.R -> the first version used in v0.2.6
25 # 2. run-polap-r-pairs.R -> used in test-reads in oga
26 # 3. run-polap-r-bridge.R -> used in test-reads in oga
27 # 4. run-polap-r-select-reads-polap.R -> the 2nd version not used
28 # 5. run-polap-r-select-reads-ptgaul.R -> a slim version of the 2nd used in dga
29 # 6. run-polap-r-directional.R -> used by an older version of dga
30 #   or _run_polap_original-directional-reads
31 #
32 # This used to be used in dga and it is not used in the latest dga.
33 #
34 # See Also:
35 # run-polap-pairs.R
36 # run-polap-r-pairs.R
37 # run-polap-r-bridge.R
38 # run-polap-r-select-reads-polap.R
39 # run-polap-r-select-reads-ptgaul.R
40 # run-polap-r-directional.R
41 # run-polap-function-oga.sh
42 # run-polap-function-dga.sh
43 #
44 # NOTE: move with directional module.
45 #####
46 #
47 # name: select long reads
48 #
49 # synopsis:
50 #   run-polap-pairs.R mt.contig.name-1 contig.tab ${MTSEEDSDIR} $SINGLE_MIN $SINGLE_MIN
51 #

```

7|M-[M-^T 03, 25 13:43

run-polap-r-directional.R

Page 2/9

```

52 # requirement: executes Flye
53 # flye --nano-raw "$LR3K" --out-dir "${_arg_outdir}" \
54 #   --threads "${_arg_threads}" \
55 #   --stop-after contigger \
56 #   --asm-coverage 30 \
57 #   --genome-size "$EXPECTED_GENOME_SIZE"
58 #
59 # input:
60 #   1. mt.contig.name-1 - contig or edge numbers
61 #   2. contig.tab - minimap2 output modified
62 #   3. seeds directory for output
63 #   4. pair minimum length V11: 3000 for MT, 1000 for PT
64 #   5. bridge minimum length V7: depends: 3000 or 5000
65 #   6. single minimum length V11: 3000 for MT, 0 or 1000 for PT
66 #
67 # output:
68 #   single.names and <contig1-contig2.name> files in the output directory
69 #
70 # MTSEEDSDIR="${_arg_outdir}"/60-mt-$STEP4/o${MR}/seeds
71 #   run-polap-pairs.R mt.contig.name-1 contig.tab ${MTSEEDSDIR} $PAIR_MIN $BRIDGE_MIN $SINGLE_MIN
72 #   run-polap-pairs.R mt.contig.name-1 contig.tab ${MTSEEDSDIR} $PAIR_MIN $BRIDGE_MIN
73
74 suppressWarnings(suppressPackageStartupMessages(library("optparse")))
75 suppressWarnings(suppressPackageStartupMessages(library("dplyr")))
76 suppressWarnings(suppressPackageStartupMessages(library("readr")))
77 suppressWarnings(suppressPackageStartupMessages(library("purrr")))
78 suppressWarnings(suppressPackageStartupMessages(library("tidyverse")))
79 suppressWarnings(suppressPackageStartupMessages(library("ggplot2")))
80
81 debug <- Sys.getenv("_POLAP_DEBUG", unset = "0")
82
83 parser <- OptionParser()
84 parser <- add_option(parser, c("-t", "--table"),
85   action = "store",
86   help = "input minimap2 PAF-like tabular format file",
87   metavar = "<FILE>")
88 )
89 parser <- add_option(parser, c("-m", "--mtcontigname"),
90   action = "store",
91   help = "mt.contig.name-1",
92   metavar = "<FILE>")
93 )
94 parser <- add_option(parser, c("-o", "--out"),
95   action = "store",
96   help = "output directory"
97 )
98 parser <- add_option(parser, c("-r", "--pair-min"),
99   type = "integer",
100  default = 3000,
101  help = "Minimum length of pair-mapping alignment",
102  metavar = "number"

```

7IM-[M-^T 03, 25 13:43

run-polap-r-directional.R

Page 3/9

```

103  )
104  parser ← add_option(parser, c("-w", "--single-min"),
105   type = "integer",
106   default = 3000,
107   help = "Minimum length of single-mapping alignment",
108   metavar = "number"
109  )
110 parser ← add_option(parser, c("-x", "--bridge-min"),
111   type = "integer",
112   default = 3000,
113   help = "Minimum length of bridging reads",
114   metavar = "number"
115  )
116 parser ← add_option(parser, c("--intra-base-ratio"),
117   type = "double",
118   default = 0.7,
119   help = "Minimum of V10/V11",
120   metavar = "number"
121  )
122 parser ← add_option(parser, c("--intra-read-ratio"),
123   type = "double",
124   default = 0.7,
125   help = "Minimum of (V4-V3)/V2",
126   metavar = "number"
127  )
128 parser ← add_option(parser, c("--inter-base-ratio"),
129   type = "double",
130   default = 0.7,
131   help = "Minimum of V10/V11",
132   metavar = "number"
133  )
134 parser ← add_option(parser, c("--use-strand"),
135   action = "store_true",
136   default = FALSE, help = "will create ptgaul.names"
137  )
138 parser ← add_option(parser, c("--use-reverse"),
139   action = "store_true",
140   default = FALSE, help = "will use the reads mapped on the reverse complement"
141  )
142 parser ← add_option(parser, c("--create-ptgaul"),
143   action = "store_true",
144   default = FALSE, help = "will create ptgaul.names"
145  )
146 parser ← add_option(parser, c("--create-single"),
147   action = "store_true",
148   default = FALSE, help = "will create single.names"
149  )
150 parser ← add_option(parser, c("--create-pair"),
151   action = "store_true",
152   default = FALSE, help = "will create pair.names"
153  )

```

7|M-[M-^T 03, 25 13:43

run-polap-r-directional.R

Page 4/9

```

154 parser ← add_option(parser, c("--create-combined"),
155   action = "store_true",
156   default = FALSE, help = "will create combined.names"
157 )
158 parser ← add_option(parser, c("--all"),
159   action = "store_true",
160   default = FALSE, help = "will create all the 4 names files."
161 )
162 parser ← add_option(parser, c("--outptgaul"),
163   action = "store",
164   default = "ptgaul.names",
165   help = "ptgaul output file"
166 )
167 parser ← add_option(parser, c("--outsingle"),
168   action = "store",
169   default = "single.names",
170   help = "single output file"
171 )
172 parser ← add_option(parser, c("--outpair"),
173   action = "store",
174   default = "pair.names",
175   help = "pair output file"
176 )
177 parser ← add_option(parser, c("--outcombined"),
178   action = "store",
179   default = "combined.names",
180   help = "combined output file"
181 )
182
183 args1 ← parse_args(parser)
184
185 if (is_null(args1$table)) {
186   s ← "bioprojects"
187   o ← "PRJNA817235-Canavalia_ensiformis"
188
189   input_dir0 ← file.path(".")
190   input1 ← file.path(input_dir0, "contig.tab")
191   input2 ← file.path(input_dir0, "mt.contig.name-1")
192   output1 ← file.path(input_dir0, "02-reads")
193
194   args1 ← parse_args(parser, args = c(
195     "--table", input1,
196     "--mtcontigname", input2,
197     "--o", output1,
198     "--r", 10000,
199     "--w", 10000,
200     "--x", 10000,
201     "--all",
202     "--use-strand",
203     "--use-reverse",
204     "--intra-base-ratio", 0.7,

```

7|M-[M-^T 03, 25 13:43

run-polap-r-directional.R

Page 5/9

```

205     "--intra-read-ratio", 0.7,
206     "--inter-base-ratio", 0.7
207   })
208 }
209
210 # mt.contig.name-x
211 x1 <- as_tibble(read.table(args1$mtcontigname))
212
213 # https://lh3.github.io/minimap2/minimap2.html
214 # Li 2008 - "OUTPUT FORMAT: Minimap2 outputs mapping positions in the
215 # Pairwise mAPPING Format (PAF) by default. PAF is a TAB-delimited text format
216 # with each line consisting of at least 12 fields as are described in the
217 # following table:
218 #
219 # Col Type Description
220 # 1 string Query sequence name
221 # 2 int Query sequence length
222 # 3 int Query start coordinate (0-based)
223 # 4 int Query end coordinate (0-based)
224 # 5 char âM-^@M-^X+âM-^@M-^Y if query/target on the same strand; âM-^@M-^X-âM-^@M-^Y if opposite
225 # 6 string Target sequence name
226 # 7 int Target sequence length
227 # 8 int Target start coordinate on the original strand
228 # 9 int Target end coordinate on the original strand
229 # 10 int Number of matching bases in the mapping
230 # 11 int Number bases, including gaps, in the mapping
231 # 12 int Mapping quality (0-255 with 255 for missing)
232 #
233 # Read in the data with assigned column names
234 data <- read_tsv(args1$table,
235   show_col_types = FALSE,
236   col_names = c(
237     "rname",
238     "rlen",
239     "rstart",
240     "rend",
241     "strand",
242     "cname",
243     "clen",
244     "cstart",
245     "cend",
246     "match",
247     "base"
248   )
249 )
250
251 if (args1$`use-strand`) {
252   data <- data |>
253     filter(strand == "+")
254 }
255

```

7|M-[M-^T 03, 25 13:43

run-polap-r-directional.R

Page 6/9

```
256 mtdir ← args1$out
257 pair_min ← as.numeric(args1$'pair-min')
258 brigde_min ← as.numeric(args1$'bridge-min')
259 single_min ← as.numeric(args1$'single-min')
260 a ← as.numeric(args1$'intra-read-ratio')
261 b ← as.numeric(args1$'intra-base-ratio')
262 c ← as.numeric(args1$'inter-base-ratio')
263
264 ptgaul_option_base ← paste("ptgaul",
265   a,
266   b,
267   single_min,
268   c,
269   pair_min,
270   brigde_min,
271   sep = "-"
272 )
273 ptgaul_option_base ← paste0(ptgaul_option_base, ".names")
274 ptgaul_option_base ← "ptgaul.names"
275
276 single_option_base ← paste("single",
277   a,
278   b,
279   single_min,
280   c,
281   pair_min,
282   brigde_min,
283   sep = "-"
284 )
285 single_option_base ← paste0(single_option_base, ".names")
286
287 pair_option_base ← paste("pair",
288   a,
289   b,
290   single_min,
291   c,
292   pair_min,
293   brigde_min,
294   sep = "-"
295 )
296 pair_option_base ← paste0(pair_option_base, ".names")
297
298 combined_option_base ← paste("combined",
299   a,
300   b,
301   single_min,
302   c,
303   pair_min,
304   brigde_min,
305   sep = "-"
306 )
```

7|M-[M-^T 03, 25 13:43

run-polap-r-directional.R

Page 7/9

```

307 combined_option_base <- paste0(combined_option_base, ".names")
308
309 # 0. Check preconditions
310 # rstart is not greater than rend.
311 # match is not greater than base.
312 result <- data %>%
313   summarize(all_rows_meet_condition = all(rstart ≤ rend & match ≤ base))
314
315 # Check the result
316 stopifnot(result$all_rows_meet_condition)
317 stopifnot(pair_min ≥ 0, brigde_min ≥ 0, single_min ≥ 0)
318
319 # 1. ptGAUL: we use ptGAUL for the case of single edge reference
320 if (args1$'create-ptgaul' || args1$all) {
321   ptgaul_file <- file.path(mtdir, args1$outputgaul)
322   ptgaul_mapped_reads <- data |>
323     filter(match / base > args1$'intra-base-ratio', base > single_min) |>
324     select(rname) |>
325     distinct(rname)
326
327   ptgaul_mapped_reads |>
328     write.table(
329       ptgaul_file,
330       row.names = FALSE,
331       col.names = FALSE,
332       quote = FALSE
333     )
334   ptgaul_option_file <- file.path(mtdir, ptgaul_option_base)
335   print(ptgaul_file)
336   print(ptgaul_option_file)
337   file.copy(ptgaul_file, ptgaul_option_file)
338   # file.copy(ptgaul_file, ptgaul_option_file, showWarnings = FALSE)
339 }
340
341 # 2. intra-contig mapping
342 if (args1$'create-single' || args1$all) {
343   single_file <- file.path(mtdir, args1$outputsingle)
344   intra_contig_mapped_reads <- data |>
345     filter(
346       (rend - rstart) / rlen > args1$'intra-read-ratio',
347       match / base > args1$'intra-base-ratio',
348       base > single_min
349     ) |>
350     select(rname) |>
351     distinct(rname)
352
353   intra_contig_mapped_reads |>
354     write.table(single_file, row.names = FALSE, col.names = FALSE, quote = FALSE)
355
356   single_option_file <- file.path(mtdir, single_option_base)
357   print(single_file)

```

7|M-[M-^T 03, 25 13:43

run-polap-r-directional.R

Page 8/9

```

358   file.copy(single_file, single_option_file)
359   # file.copy(single_file, single_option_file, showWarnings = FALSE)
360 }
361
362 if (nrow(x1) > 1) {
363   y <- t(combn(x1$V1, 2))
364
365   orders <- function(y) {
366     stopifnot(length(y) == 2)
367     z <- as.numeric(gsub("\\D", "", y))
368     bname <- paste0(z[1], "_", z[2])
369     oname <- paste0(mtdir, "/", bname, ".name")
370     y1 <- data |>
371       filter(
372         cname == y[1],
373         match / base > args1$`inter-base-ratio`,
374         base > pair_min,
375         rlen > brigde_min
376       ) |>
377       select(rname) |>
378       distinct(rname)
379     y2 <- data |>
380       filter(
381         cname == y[2],
382         match / base > args1$`inter-base-ratio`,
383         base > pair_min,
384         rlen > brigde_min
385       ) |>
386       select(rname) |>
387       distinct(rname)
388     intersect(y1, y2) |>
389     write.table(ongame, row.names = FALSE, col.names = FALSE, quote = FALSE)
390   }
391
392
393 # 3. inter-contig mapping: individual pair
394 if (args1$all) {
395   apply(y, 1, orders)
396 }
397
398 # 4. inter-contig mapping
399 #
400 # Filter data based on the given criteria
401
402 if (args1$`create-pair` || args1$all) {
403   inter_contig_mapped_reads <- data |>
404     # Calculate the match/base ratio
405     mutate(match_base_ratio = match / base) |>
406     # Filter for match/base > 0.7 and base > 3000
407     filter(
408       match_base_ratio > args1$`inter-base-ratio`,

```

7|M-[M-^T 03, 25 13:43

run-polap-r-directional.R

Page 9/9

```

409     base > pair_min,
410     rlen > bridge_min
411   ) |>
412   # Count the number of contigs for each read
413   group_by(rname) |>
414   # Keep only those reads that are mapped to at least two contigs
415   filter(n_distinct(cname) ≥ 2) |>
416   # Select unique reads
417   distinct(rname) |>
418   ungroup()
419
420 pair_file ← file.path(mtdir, args1$outpair)
421 inter_contig_mapped_reads |>
422   write.table(
423     pair_file,
424     row.names = FALSE,
425     col.names = FALSE,
426     quote = FALSE
427   )
428 pair_option_file ← file.path(mtdir, pair_option_base)
429 print(pair_file)
430 file.copy(pair_file, pair_option_file)
431 # file.copy(pair_file, pair_option_file, showWarnings = FALSE)
432 }
433
434 # 5. combined: single + pair
435 if (args1$`create-combined` || args1$all) {
436   combined_mapped_reads ← bind_rows(
437     inter_contig_mapped_reads,
438     intra_contig_mapped_reads
439   ) |>
440   distinct(rname, .keep_all = TRUE)
441
442 combined_file ← file.path(mtdir, args1$outcombined)
443 combined_mapped_reads |>
444   write.table(
445     combined_file,
446     row.names = FALSE,
447     col.names = FALSE,
448     quote = FALSE
449   )
450 combined_option_file ← file.path(mtdir, combined_option_base)
451 print(combined_file)
452 file.copy(combined_file, combined_option_file)
453 # file.copy(combined_file, combined_option_file, showWarnings = FALSE)
454 }
455 }
```