

9IM-[M-^T 15, 25 17:43

**polap-bash-oatk-ont-sidekicks-stage1.sh**

Page 1/7

```

1 #!/usr/bin/env bash
2 # polap-bash-oatk-ont-sidekicks-stage0.sh
3 # Stage-0 prefilters ONLY for ONT organelle assembly:
4 #   trim (porechop_abi), scrub (yacrd), lenfilt (filtlong/seqkit),
5 #   HMM mt-bait, plastid drop, duplex-only, rare-mask (KMC), correction (Canu/Ratatosk).
6 #
7 # Outputs:
8 #   OUT/pre/reads.pre.fq    (or .fa/.fq.gz depending on tools used)
9 #   OUT/pre/stage0.report.tsv
10 #
11 # Notes:
12 #   - Does NOT do HPC or assembly. Feed reads.pre.fq into Stage-1/2 wrapper.
13
14 set -euo pipefail
15
16 : "${POLAP_LOG_LEVEL:=1}" # 0=quiet, 1=info, 2=verbose
17 log() {
18     local l=$1
19     shift
20     [[ $POLAP_LOG_LEVEL -ge $l ]] && echo "$@" >&2
21 }
22 die() {
23     echo "[ERR]" "$@" >&2
24     exit 1
25 }
26 need() { command -v "$1" >/dev/null 2>&1 || die "missing dependency: $1"; }
27
28 # ----- CLI defaults -----
29 READS=""
30 OUT=""
31 THREADS=32
32 OATKDB=""
33 CLADE="magnoliopsida"
34 HMMBIN="hmmannot"
35 NHMMSCAN="nhmmscan"
36
37 DO_TRIM=0
38 DO_SCRUB=0
39 DO_LENfilt=0
40 LEN_MIN=5000
41 KEEP_PCT=70
42 DO_MT_BAIT=0
43 DO_DROP_PLASTID=0
44 DO_DUPLEX_ONLY=0
45 DUPLEX_TAG="duplex"
46 DO_RARE_MASK=0
47 RARE_K=31
48 RARE_MAX=1
49 DO_CORRECT=0
50 SHORT_READS=""
51

```

9IM-[M-^T 15, 25 17:43

**polap-bash-oatk-ont-sidekicks-stage1.sh**

Page 2/7

```

52 usage() {
53   cat <<EOF
54 Usage:
55   $(basename "$0") --reads READS.fq.gz --out OUTDIR --oatkdb /path/OatkDB [options]
56
57 Required:
58   --reads FILE          ONT reads (fa/fq[.gz])
59   --out DIR             output directory
60   --oatkdb DIR          OatkDB root (must contain v20230921)
61
62 General:
63   --threads INT         (default ${THREADS})
64   --clade NAME          fam prefix under v20230921 (default ${CLADE})
65   --hmm-bin PATH        hmmannot/hmm_annotation (default ${HMMBIN})
66   --nhmmscan PATH      nhmmscan path (default ${NHMMSCAN})
67   -v|--verbose          verbose
68   -q|--quiet            quiet
69
70 Stage-0 prefilters (all optional):
71   --trim                porechop_abi
72   --scrub               minimap2 ava-ont + yacrd
73   --lenfilt              filtlong (preferred) or seqkit
74   --len-min INT         min length (default ${LEN_MIN})
75   --keep-pct INT        filtlong keep percent (default ${KEEP_PCT})
76   --mt-bait              keep reads hitting mito fam (HMM)
77   --drop-plastid         drop reads hitting plastid fam (HMM)
78   --duplex-only          keep reads with tag (default '${DUPLEX_TAG}')
79   --duplex-tag STR      header tag (default '${DUPLEX_TAG}')
80   --rare-mask            KMC rare-kmer masking (heavy placeholder)
81   --rare-k INT           (default ${RARE_K})
82   --rare-max INT         (default ${RARE_MAX})
83   --correct-canu          Canu correction-only
84   --correct-ratatosk     Ratatosk (needs --short-reads)
85   --short-reads FILE    short reads for Ratatosk
86 EOF
87 }
88
89 # ----- parse args -----
90 while [[ $# -gt 0 ]]; do
91   case "$1" in
92     --reads)
93       READS="$2"
94       shift 2
95       ;;
96     --out)
97       OUT="$2"
98       shift 2
99       ;;
100    --oatkdb)
101      OATKDB="$2"
102      shift 2

```

9IM-[M-^T 15, 25 17:43

**polap-bash-oatk-ont-sidekicks-stage1.sh**

Page 3/7

```
103      ;;
104      --threads)
105      THREADS="$2"
106      shift 2
107      ;;
108      --clade)
109      CLADE="$2"
110      shift 2
111      ;;
112      --hmm-bin)
113      HMMBIN="$2"
114      shift 2
115      ;;
116      --nhmmscan)
117      NHMMSCAN="$2"
118      shift 2
119      ;;
120      -v | --verbose)
121      POLAP_LOG_LEVEL=2
122      shift
123      ;;
124      -q | --quiet)
125      POLAP_LOG_LEVEL=0
126      shift
127      ;;
128      --trim)
129      DO_TRIM=1
130      shift
131      ;;
132      --scrub)
133      DO_SCRUB=1
134      shift
135      ;;
136      --lenfilt)
137      DO_LENFILT=1
138      shift
139      ;;
140      --len-min)
141      LEN_MIN="$2"
142      shift 2
143      ;;
144      --keep-pct)
145      KEEP_PCT="$2"
146      shift 2
147      ;;
148      --mt-bait)
149      DO_MT_BAIT=1
150      shift
151      ;;
152      --drop-plastid)
153      DO_DROP_PLASTID=1
```

9IM-[M-^T 15, 25 17:43

**polap-bash-oatk-ont-sidekicks-stage1.sh**

Page 4/7

```

154     shift
155     ;;
156   --duplex-only)
157     DO_DUPLEX_ONLY=1
158   shift
159   ;;
160   --duplex-tag)
161     DUPLEX_TAG="$2"
162   shift 2
163   ;;
164   --rare-mask)
165     DO_RARE_MASK=1
166   shift
167   ;;
168   --rare-k)
169     RARE_K="$2"
170   shift 2
171   ;;
172   --rare-max)
173     RARE_MAX="$2"
174   shift 2
175   ;;
176   --correct-canu)
177     DO_CORRECT=1
178   shift
179   ;;
180   --correct-ratatosk)
181     DO_CORRECT=2
182   shift
183   ;;
184   --short-reads)
185     SHORT_READS="$2"
186   shift 2
187   ;;
188   -h | --help)
189     usage
190     exit 0
191   ;;
192   *)
193     echo "[ERR] unknown arg: $1" >&2
194     usage
195     exit 1
196   ;;
197   esac
198 done
199 [ [ -z "$READS" || -z "$OUT" || -z "$OATKDB" ] ] && {
200   usage
201   exit 1
202 }
203 # ----- prep & deps -----

```

9iM-^M-^T 15, 25 17:43

**polap-bash-oatk-ont-sidekicks-stage1.sh**

Page 5/7

```

205 READS=$(readlink -f "$READS")
206 OUT=$(readlink -f "$OUT")
207 OATKDB=$(readlink -f "$OATKDB")
208 mkdir -p "$OUT/pre"
209 cd "$OUT/pre"
210
211 need seqkit
212 [[ $DO_TRIM -eq 1 ]] && need porechop_abi
213 [[ $DO_SCRUB -eq 1 ]] && {
214   need minimap2
215   need yacrd
216 }
217 [[ $DO_LENFILT -eq 1 ]] && { command -v filtlong >/dev/null 2>&1 || need seqkit; }
218 FAM_M="$OATKDB/v20230921/${CLADE}_mito.fam"
219 FAM_P="$OATKDB/v20230921/${CLADE}_pltd.fam"
220 [[ $DO_MT_BAIT -eq 0 || -s "$FAM_M" ]] || die "mito fam not found: $FAM_M"
221 [[ $DO_DROP_PLASTID -eq 0 || -s "$FAM_P" ]] || die "plastid fam not found: $FAM_P"
222 [[ $DO_MT_BAIT -eq 0 && $DO_DROP_PLASTID -eq 0 ]] || {
223   command -v "$HMMBIN" >/dev/null 2>&1 || die "missing $HMMBIN"
224   need "$NHMMSCAN"
225 }
226 [[ $DO_RARE_MASK -eq 1 ]] && {
227   need kmc
228   command -v kmc_tools >/dev/null 2>&1 || need kmc_tools
229 }
230 [[ $DO_CORRECT -eq 1 ]] && need canu
231 [[ $DO_CORRECT -eq 2 ]] && {
232   need ratatosk || need Ratatosk || true
233   [[ -n "$SHORT_READS" ]] || die "--short-reads required"
234 }
235
236 # ----- Stage-0 run -----
237 CUR="$READS"
238
239 # trim
240 if [[ $DO_TRIM -eq 1 ]]; then
241   log 1 "[trim] porechop_abi"
242   porechop_abi -i "$CUR" -o reads.trim.fq -t "$THREADS"
243   CUR="$PWD/reads.trim.fq"
244 fi
245 # scrub
246 if [[ $DO_SCRUB -eq 1 ]]; then
247   log 1 "[scrub] minimap2 ava-ont + yacrd"
248   minimap2 -x ava-ont -g 500 -t "$THREADS" "$CUR" "$CUR" >overlaps.paf
249   yacrd -i overlaps.paf -o yacrd.tsv -c 4 -n -p 0.4
250   awk '$3=="chimera"{next} $3=="unchanged"{print $1}' yacrd.tsv >keep.ids
251   seqkit grep -f keep.ids "$CUR" -o reads.scrub.fq
252   CUR="$PWD/reads.scrub.fq"
253 fi
254 # lenfilt
255 if [[ $DO_LENFILT -eq 1 ]]; then

```

9iM-[M-^T 15, 25 17:43

**polap-bash-oatk-ont-sidekicks-stage1.sh**

Page 6/7

```

256 if command -v filtlong >/dev/null 2>&1; then
257   log 1 "[lenfilt] filtlong --min_length $LEN_MIN --keep_percent $KEEP_PCT"
258   filtlong --min_length "$LEN_MIN" --keep_percent "$KEEP_PCT" "$CUR" >reads.long fq
259   CUR="$PWD/reads.long fq"
260 else
261   log 1 "[lenfilt] seqkit seq -m $LEN_MIN"
262   seqkit seq -m "$LEN_MIN" "$CUR" -o reads.len fq
263   CUR="$PWD/reads.len fq"
264 fi
265 fi
266 # HMM mt-bait
267 if [[ $DO_MT_BAIT -eq 1 ]]; then
268   log 1 "[mt-bait pre] $HMMBIN vs mito fam"
269   "$HMMBIN" -t "$THREADS" --nhmmscan "$NHMMSCAN" -o reads.mito.txt "$FAM_M" "$CUR"
270   awk -v FS='[:space:]+' '$1!~/^#/ {print $3}' reads.mito.txt | sort -u >mt.ids
271   seqkit grep -f mt.ids "$CUR" -o reads.mt fq
272   CUR="$PWD/reads.mt fq"
273 fi
274 # plastid drop
275 if [[ $DO_DROP_PLASTID -eq 1 ]]; then
276   log 1 "[plastid pre] drop plastid reads"
277   "$HMMBIN" -t "$THREADS" --nhmmscan "$NHMMSCAN" -o reads.pt.txt "$FAM_P" "$CUR"
278   awk -v FS='[:space:]+' '$1!~/^#/ {print $3}' reads.pt.txt | sort -u >pt.ids
279   seqkit grep -v -f pt.ids "$CUR" -o reads.mt.noPT fq
280   CUR="$PWD/reads.mt.noPT fq"
281 fi
282 # duplex-only
283 if [[ $DO_DUPLEX_ONLY -eq 1 ]]; then
284   log 1 "[duplex] keep header tag: $DUPLEX_TAG"
285   awk -v tag="$DUPLEX_TAG" 'BEGIN{OFS="\n"}'
286   NR%4==1{keep=index($0,tag)>0;h=$0}
287   NR%4==2{s=$0}
288   NR%4==3{p=$0}
289   NR%4==0{q=$0; if(keep){print h,s,p,q}}' "$CUR" >reads.duplex fq
290   CUR="$PWD/reads.duplex fq"
291 fi
292 # rare-mask (placeholder)
293 if [[ $DO_RARE_MASK -eq 1 ]]; then
294   log 1 "[rare-k] KMC k=$RARE_K, <=$RARE_MAX"
295   mkdir -p kmc_tmp
296   kmc -k "$RARE_K" -t "$THREADS" -ci1 -cs100000000 "$CUR" kmc.db kmc_tmp/
297   kmc_tools transform kmc.db dump -ci1 -cx"$RARE_MAX" rare_k.txt
298   cp "$CUR" reads.mask fq
299   CUR="$PWD/reads.mask fq"
300 fi
301 # correction
302 if [[ $DO_CORRECT -eq 1 ]]; then
303   log 1 "[correct] Canu correction-only"
304   mkdir -p corr
305   canu -correct -p ontcorr -d corr genomeSize=500k -nanopore-raw "$CUR" 2>corr/canu.log || log 1 "[warn] canu failed?"
306   if ls corr/ontcorr.correctedReads.fasta.gz >/dev/null 2>&1; then CUR="$PWD/corr/ontcorr.correctedReads.fasta.gz"; fi

```

9iM-^M-^T 15, 25 17:43

**polap-bash-oatk-ont-sidekicks-stage1.sh**

Page 7/7

```

307 elif [[ $DO_CORRECT -eq 2 ]]; then
308   log 1 "[correct] Ratatosk"
309   need ratatosk || need Ratatosk || true
310   [[ -n "$SHORT_READS" ]] || die "--short-reads required for Ratatosk"
311   ratatosk -s "$SHORT_READS" -l "$CUR" -o reads.rtk.fq -t "$THREADS" || log 1 "[warn] ratatosk failed?"
312   CUR="$PWD/reads.rtk.fq"
313 fi
314
315 # emit final
316 FINAL="$PWD/reads.pre fq"
317 cp -f "$CUR" "$FINAL"
318
319 # report
320 {
321   echo -e "key\tvalue"
322   echo -e "reads_in\t$READS"
323   echo -e "out\t$OUT"
324   echo -e "reads_pre\t$FINAL"
325   echo -e "trim\t$([[ $DO_TRIM -eq 1 ]]) && echo ON || echo OFF)"
326   echo -e "scrub\t$([[ $DO_SCRUB -eq 1 ]]) && echo ON || echo OFF)"
327   echo -e "len_min\t$([[ $DO_LENFILT -eq 1 ]]) && echo $LEN_MIN || echo .)"
328   echo -e "keep_pct\t$([[ $DO_LENFILT -eq 1 ]]) && $(
329     command -v filtlong >/dev/null
330     echo $?
331   ) -eq 0 ])&& echo $KEEP_PCT || echo .)"
332   echo -e "mt_bait_prefilter\t$([[ $DO_MT_BAIT -eq 1 ]]) && echo ON || echo OFF)"
333   echo -e "drop_plastid_prefilter\t$([[ $DO_DROP_PLASTID -eq 1 ]]) && echo ON || echo OFF)"
334   echo -e "duplex_only\t$([[ $DO_DUPLEX_ONLY -eq 1 ]]) && echo $DUPLEX_TAG || echo OFF)"
335   echo -e "rare_mask\t$([[ $DO_RARE_MASK -eq 1 ]]) && echo K=$RARE_K,max=$RARE_MAX || echo OFF)"
336   echo -e "corrector\t$([[ $DO_CORRECT -eq 1 ]]) && echo CANU || ([[ $DO_CORRECT -eq 2 ]]) && echo RATATOSK || echo OFF))"
337 } > "$OUT/pre/stage0.report.tsv"
338
339 log 1 "[Stage-0 done] Preprocessed reads: $FINAL"

```