

9jM-^M-^T 27, 25 13:53

polap-bash-fast-mtseed-ont.sh

Page 1/20

9iM-^M-^T 27, 25 13:53

polap–bash–fast–mtseed–ont.sh

Page 2/20

```

44 source "${_POLAPLIB_DIR}/polap-lib-conda.sh" || true
45 HELP="${_POLAPLIB_DIR}/fast-mtseed-ont"
46 mkdir -p "${HELP}"
47
48 # Helpers (fast-mtseed)
49 OLPY="${HELP}/polap-py-overlapness-from-paf.py"
50 TOPPY="${HELP}/polap-py-topfrac-by-wdeg.py"
51 FILTPY="${HELP}/polap-py-filter-paf-by-ids.py"
52 CDSPY="${HELP}/polap-py-cds-coverage-from-paf.py"
53 SEEDREADSPY="${HELP}/polap-py-seed-mapped-reads.py"
54 STOPPY="${HELP}/polap-py-stop-delta.py"
55 AWK_CONS="${HELP}/polap-awk-filter-conservative.awk"
56 RECRUIT_SH="${HELP}/polap-bash-recruit-count.sh"
57 PAF2MREADS_SH="${HELP}/polap-bash-paf2mreads.sh"
58 R_PICK="${HELP}/polap-r-pick-high-recruit.R"
59 THRESH_NUC_PY="${HELP}/polap-py-threshold-from-nuclear.py"
60 PY_PT_THRESH="${HELP}/polap-py-pt-ident-threshold.py"
61
62 # Helpers for PT removal & isoforms
63 PT_ISOFORM_SH="${_POLAPLIB_DIR}/polap-bash-pt-isoform.sh"
64
65 # scan+stream+edges mode
66 TOPKPY="${HELP}/paf-topk-per-q.py"
67 EDGPY="${HELP}/paf-to-edges-stream.py"
68 EDGECOMP="${HELP}/polap-py-edges-components.py" # optional
69
70 # Logging & profiling
71 # Base helper for all note levels
72 _note_base() {
73     local lvl=$1 # message verbosity level
74     shift
75     local msg="$*"
76
77     # Format with timestamp + caller
78 }

```

9|M-[M-^T 27, 25 13:53

polap-bash-fast-mtseed-ont.sh

Page 3/20

```

87 local src="${BASH_SOURCE[2]##*/}" # 2 to skip _note_base + noteX
88 local ln="${BASH_LINENO[1]}"
89 local ts
90 ts=$(date +'%F %T')
91 local line="[${ts}][${src}]:${ln}] $msg"
92
93 # Always log, but only display if verbosity allows
94 if [[ "${VERBOSE:-0}" -gt "$lvl" ]]; then
95   if [[ "${_arg_log_stderr:-off}" == "off" ]]; then
96     printf "%s\n" "$line" >&3 # real screen
97   else
98     printf "%s\n" "$line" >&2
99   fi
100 fi
101
102 # Always append to log (stdout goes to logit, or LOG if you like)
103 printf "%s\n" "$line"
104 }
105
106 # Note functions
107 note() { _note_base 1 "$@"; }
108 note0() { _note_base 0 "$@"; } # always
109 note1() { _note_base 1 "$@"; } # shown with -v
110 note2() { _note_base 2 "$@"; } # shown with -v -v
111 note3() { _note_base 3 "$@"; } # shown with -v -v -v
112
113 cmd() {
114   if ((DRY)); then
115     local p="$1"
116     shift
117     printf "[DRY] %q" "$p" | tee -a "$LOG_FD"
118     for a in "$@"; do printf "%q" "$a" | tee -a "$LOG_FD"; done
119     printf "\n" | tee -a "$LOG_FD"
120     return 0
121   fi
122   local p="$1"
123   shift
124   printf "[RUN] %q" "$p" | tee -a "$LOG_FD"
125   for a in "$@"; do printf "%q" "$a" | tee -a "$LOG_FD"; done
126   printf "\n" | tee -a "$LOG_FD"
127   "$p" "$@"
128 }
129 prof_start() { ((PROF)) && PROF_T0=$(date +%s) && PROF_STEP="$1" || true; }
130 prof_end() {
131   ((PROF)) || return 0
132   local t1=$(date +%s)
133   echo -e "${PROF_STEP}\t$((t1 - PROF_T0))" >> "$PROF_FILE"
134 }
135
136 require_tools() {
137   local miss=0

```

9jM-^M-^T 27, 25 13:53

polap-bash-fast-mtseed-ont.sh

Page 4/20

```

138 for t in "$@"; do
139   command -v "$t" >/dev/null 2>&1 || {
140     note "ERR missing tool: $t"
141     miss=1
142   }
143 done
144 if ((miss == 1)); then
145   note "ERR load modules/conda"
146   exit 127
147 fi
148 }

# CLI (merged options)
151 # reads=""
152 # pt_ref=""
153 outdir="fast_mtseed_ont_out"
154 threads=32
155 # PT removal knobs
156 pt_origin=""
157 nuc_origin=""
158 alen_min=3000
159 identity_min="" # override identity if set
160 fpr=0.01
161 tpr=0.95
162 profiling=0
163 # All-vs-all / selection / assembly
164 # overlap_mode="default" # default | scan+stream+edges
165 overlap_mode="scan+stream+edges"
166 scan_keep_frac="0.40"
167 scan_k=21
168 scan_w=9
169 scan_N=1
170 scan_mask="0.70"
171 scan_minocc=20
172 refine_k=19
173 refine_w=7
174 refine_N=30
175 refine_mask="0.60"
176 refine_minocc=10
177 topk_per_q=0
178 topk_edges_per_node=0

```

9|M-[M-^T 27, 25 13:53

polap-bash-fast-mtseed-ont.sh

Page 5/20

```

181 keep_scan_paf=0
182 gzip_edges=1
183 top_frac="0.20"
184 min_olen=1200
185 min_ident="0.84"
186 w_floor="0.12"
187 assembler="miniasm"
188 # intergenic / BUSCO
189 mcds=""
190 nprot=""
191 nuc_ids_opt=""
192 # iteration
193 rounds_max=1
194 delta_stop="0.05"
195 seed=13
196 steps=""
197 # nuclear overlap deplete is handled in separate tool; here we select seeds only
198
199 usage() {
200     cat <<EOF
201 polap-bash-fast-mtseed-ont.sh v0.3.0 (merged PT removal + fast mtseed)
202 Usage:
203     $0 -r reads.fq.gz -p pt_ref.fa -o outdir [options]
204
205 Required:
206     -r, --reads          ONT reads (FASTQ[.gz])
207     -p, --pt-ref         plastid reference FASTA (single circle)
208
209 PT removal:
210     --pt-origin FILE    PT-origin read IDs (weak labels) [optional]
211     --mt-origin FILE    Nuclear-origin read IDs (weak labels) [optional]
212     --alen-min INT      aligned length guard [${alen_min}]
213     --identity-min FLOAT override auto identity cutoff (e.g., 0.93)
214     --fpr FLOAT          nuclear FPR target for identity [${fpr}]
215     --tpr FLOAT          PT TPR target for identity [${tpr}]
216
217 All-vs-all:
218     --overlap-mode MODE default | scan+stream+edges [${overlap_mode}]
219     --scan-keep-frac FLOAT shortlist fraction by wdegree [${scan_keep_frac}]
220     --scan-k INT --scan-w INT --scan-N INT --scan-mask F --scan-minocc INT
221     --refine-k INT --refine-w INT --refine-N INT --refine-mask F --refine-minocc INT
222     --topk-per-q INT    cap scan hits per query (0=off)
223     --topk-edges-per-node INT cap edges per node in reducer (0=off)
224     --keep-scan-paf      store scan.paf.gz for QA
225     --no-gzip-edges      keep edges.tsv uncompressed
226     --top-frac FLOAT     top fraction if no nuclear IDs [${top_frac}]
227     --min-olen INT --min-ident F --w-floor F   edge filters [${min_olen}, ${min_ident}, ${w_floor}]
228
229 Intergenic & assembly:
230     -m, --mt-cds FILE    mito CDS proteins (miniprot) [optional]
231     -n, --nuc-prot FILE  BUSCO proteins (unused here; leave for future)

```

9|M-[M-^T 27, 25 13:53

polap-bash-fast-mtseed-ont.sh

Page 6/20

```

232 --assembler NAME          miniasm|raven [ ${assembler} ]
233 --rounds-max INT         iteration rounds [ ${rounds_max} ]
234 --delta-stop FLOAT       stop if |M-^T|M| < x [ ${delta_stop} ]
235
236 General:
237   -o, --outdir DIR        output directory [ ${outdir} ]
238   -t, --threads INT       threads [ ${threads} ]
239   --seed INT              RNG seed [ ${seed} ]
240   -v, --verbose           increase verbosity
241   --quiet                 quiet
242   --dry                   dry-run (no execution)
243   --profiling             per-step timings (log/profile.tsv)
244   -h, --help               help
245 EOF
246 }

247 while [[ $# -gt 0 ]]; do
248   case "$1" in
249     -r | --reads)
250       reads="$2"
251       shift 2
252       ;;
253     --step)
254       steps="$2"
255       shift 2
256       ;;
257     -p | --pt-ref)
258       pt_ref="$2"
259       shift 2
260       ;;
261     -o | --outdir)
262       outdir="$2"
263       shift 2
264       ;;
265     --nuc-ids)
266       nuc_ids_opt="$2"
267       shift 2
268       ;;
269     -t | --threads)
270       threads="$2"
271       shift 2
272       ;;
273     --pt-origin)
274       pt_origin="$2"
275       shift 2
276       ;;
277     --mt-origin)
278       mt_origin="$2"
279       shift 2
280       ;;
281     --nuc-origin)
282       ;;

```

9|M-[M-^T 27, 25 13:53

polap-bash-fast-mtseed-ont.sh

Page 7/20

```

283     nuc_origin="$2"
284     shift 2
285     ;;
286   --alen-min)
287     alen_min="$2"
288     shift 2
289     ;;
290   --identity-min)
291     identity_min="$2"
292     shift 2
293     ;;
294   --fpr)
295     fpr="$2"
296     shift 2
297     ;;
298   --tpr)
299     tpr="$2"
300     shift 2
301     ;;
302   --overlap-mode)
303     overlap_mode="$2"
304     shift 2
305     ;;
306   --scan-keep-frac)
307     scan_keep_frac="$2"
308     shift 2
309     ;;
310   --scan-k)
311     scan_k="$2"
312     shift 2
313     ;;
314   --scan-w)
315     scan_w="$2"
316     shift 2
317     ;;
318   --scan-N)
319     scan_N="$2"
320     shift 2
321     ;;
322   --scan-mask)
323     scan_mask="$2"
324     shift 2
325     ;;
326   --scan-minocc)
327     scan_minocc="$2"
328     shift 2
329     ;;
330   --refine-k)
331     refine_k="$2"
332     shift 2
333     ;;

```

9|M-[M-^T 27, 25 13:53

polap-bash-fast-mtseed-ont.sh

Page 8/20

```

334      --refine-w)
335      refine_w="$2"
336      shift 2
337      ;;
338      --refine-N)
339      refine_N="$2"
340      shift 2
341      ;;
342      --refine-mask)
343      refine_mask="$2"
344      shift 2
345      ;;
346      --refine-minocc)
347      refine_minocc="$2"
348      shift 2
349      ;;
350      --topk-per-q)
351      topk_per_q="$2"
352      shift 2
353      ;;
354      --topk-edges-per-node)
355      topk_edges_per_node="$2"
356      shift 2
357      ;;
358      --keep-scan-paf)
359      keep_scan_paf=1
360      shift
361      ;;
362      --no-gzip-edges)
363      gzip_edges=0
364      shift
365      ;;
366      --top-frac)
367      top_frac="$2"
368      shift 2
369      ;;
370      --min-olen)
371      min_olen="$2"
372      shift 2
373      ;;
374      --min-ident)
375      min_ident="$2"
376      shift 2
377      ;;
378      --w-floor)
379      w_floor="$2"
380      shift 2
381      ;;
382      -m | --mt-cds)
383      mcds="$2"
384      shift 2

```

9|M-[M-^T 27, 25 13:53

polap-bash-fast-mtseed-ont.sh

Page 9/20

```

385      ;;
386      -n | --nuc-prot)
387          nprot="$2"
388          shift 2
389          ;;
390          --rounds-max)
391              rounds_max="$2"
392              shift 2
393              ;;
394              --delta-stop)
395                  delta_stop="$2"
396                  shift 2
397                  ;;
398                  --seed)
399                      seed="$2"
400                      shift 2
401                      ;;
402                      -v | --verbose)
403                          VERBOSE=$((VERBOSE + 1))
404                          shift
405                          ;;
406                          --quiet | -q)
407                              VERBOSE=0
408                              shift
409                              ;;
410                              --dry)
411                                  DRY=1
412                                  shift
413                                  ;;
414                                  --profiling)
415                                      PROF=1
416                                      shift
417                                      ;;
418                                      -h | --help)
419                                          usage
420                                          exit 0
421                                          ;;
422                                          *)
423                                              note "ERR unknown arg: $1"
424                                              usage
425                                              exit 2
426                                              ;;
427                                              esac
428 done
429
430 require_tools minimap2 samtools seqkit python Rscript awk sort comm cut gzip
431
432 for s in "$OLPY" "$STOPPY" "$FILTPY" "$CDSPY" "$SEEDREADSPY" \
433     "$STOPPY" "$AWK_CONS" "$RECRUIT_SH" "$PAF2MREADS_SH" \
434     "$R_PICK" "$THRESH_NUC_PY" "$PT_ISOFORM_SH" "$PY_PT_THRESH"; do
435     [[ -s "$s" ]] || {

```

9|M-[M-^T 27, 25 13:53

polap-bash-fast-mtseed-ont.sh

Page 10/20

```

436     note "ERR missing helper: $s"
437     exit 1
438   }
439 done
440 if [ "$overlap_mode" == "scan+stream+edges" ]; then
441   for s in "$STOPKPY" "$EDGPY"; do [[ -s "$s" ]] || {
442     note "ERR missing $s"
443     exit 1
444   }; done
445 fi
446 [[ -n "$reads" && -s "$reads" ]] || {
447   note "ERR missing --reads"
448   exit 2
449 }
450 [[ -n "$pt_ref" && -s "$pt_ref" ]] || {
451   note "ERR missing --pt-ref"
452   exit 2
453 }
454 }
455 mkdir -p "$outdir" "$outdir/log"
456 LOG="${outdir}/pipeline.log"
457 LOG_FD="$LOG"
458 if ((!DRY)); then
459   exec >>(tee -a "$LOG") 2>&1
460   TEE_ACTIVE=1
461 fi
462 if ((PROF)); then
463   PROF_FILE="$outdir/log/profile.tsv"
464   echo -e "step\tseconds" >"$PROF_FILE"
465 fi
466
467 note "Outdir: $outdir"
468 note "Assembler: $assembler ; threads: $threads"
469 note "Overlap mode: $overlap_mode"
470
471 # Step gating
472 _step_set=""
473 _add_steps() {
474   local spec="$1"
475   IFS=',' read -r -a arr <<<"$spec"
476   for tok in "${arr[@]}"; do
477     if [[ "$tok" =~ ^([0-9]+)-([0-9]+)$ ]]; then
478       local a=${BASH_REMATCH[1]} b=${BASH_REMATCH[2]}
479       for ((i = a; i <= b; i++)); do _step_set+=" $i"; done
480     elif [[ "$tok" =~ ^[0-9]+$ ]]; then _step_set+=" $tok"; fi
481   done
482 }
483 _should_run() {
484   local s="$1"
485   [[ -z "$steps" ]] && return 0

```

9jM-^M-^T 27, 25 13:53

polap–bash–fast–mtseed–ont.sh

Page 11/20

9|M-^|[M-^T 27, 25 13:53

polap-bash-fast-mtseed-ont.sh

Page 12/20

```

530     if ln.startswith('>'):
531         if name is None: name=ln[1:].strip().split()[0]
532         else: break
533     else: seq.append(ln.strip())
534 S=''.join(seq)
535 with open(outp,'w') as w:
536     w.write('>pt.double\n'); w.write(S+S+'\n')
537 PY
538 }
539 DBL_A="$PANEL_DIR/pt_isomerA.double.fa"
540 DBL_B="$PANEL_DIR/pt_isomerB.double.fa"
541 prof_start "pt_double"
542 ((!DRY)) && dbld "$ISO_A" "$DBL_A"
543 ((!DRY)) && { [[ -s "$ISO_B" ]] && dbld "$ISO_B" "$DBL_B" || true; }
544 prof_end
545 note1 "A: $DBL_A"
546 [[ -s "$DBL_B" ]]; && note1 "B: $DBL_B"
547
548 # 1c) Map readsâM-^FM-^R doubled A/B
549 note0 "1c) Map readsâM-^FM-^R doubled A/B"
550 PAF_A="$outdir/map/formA.paf"
551 PAF_B="$outdir/map/formB.paf"
552 mkdir -p "$outdir/map"
553 prof_start "map_ptA"
554 note "Step1: map readsâM-^FM-^R isomerA.double > $PAF_A"
555 cmd minimap2 -x map-ont --secondary=yes -N 50 -t "$threads" "$DBL_A" "$reads" >"$PAF_A"
556 prof_end
557 if [[ -s "$DBL_B" ]]; then
558     prof_start "map_ptB"
559     note "Step1: map readsâM-^FM-^R isomerB.double > $PAF_B"
560     cmd minimap2 -x map-ont --secondary=yes -N 50 -t "$threads" "$DBL_B" "$reads" >"$PAF_B"
561     prof_end
562 fi
563
564 # 1d) Determine identity cutoff (MT-guided) & emit pt.ids
565 note0 "1d) Determine identity cutoff (MT-guided) & emit pt.ids"
566 PT_VARS="$outdir/pt_thresh.vars"
567 PT_DIAG="$outdir/pt_thresh.diag.tsv"
568 PT_IDS="$outdir/pt.ids"
569 if [[ -n "$identity_min" ]]; then
570     note "Step1: using --identity-min=$identity_min ; alen_min=$alen_min"
571     prof_start "emit_pt_ids_override"
572     if ((!DRY)); then
573         awk -v ID="$identity_min" -v AL="$alen_min" 'BEGIN{FS=OFS="\t"}'
574         NF>=12{ id=($11>0?$10/$11:0); if(id>=ID && $11+0>=AL) print $1 }' \
575             "$PAF_A" | sort -u >"$outdir/ptA.ids"
576     if [[ -s "$PAF_B" ]]; then
577         awk -v ID="$identity_min" -v AL="$alen_min" 'BEGIN{FS=OFS="\t"}'
578         NF>=12{ id=($11>0?$10/$11:0); if(id>=ID && $11+0>=AL) print $1 }' \
579             "$PAF_B" | sort -u >"$outdir/ptB.ids"
580         sort -u "$outdir/ptA.ids" "$outdir/ptB.ids" >"$PT_IDS"

```

9|M-^|[M-^T 27, 25 13:53

polap-bash-fast-mtseed-ont.sh

Page 13/20

```

581     else
582         mv "$outdir/ptA.ids" "$PT_IDS"
583     fi
584     printf "ident_min=%s\nalen_min=%d\n" "$identity_min" "$alen_min" >"$PT_VARS"
585     : >"$PT_DIAG"
586   fi
587   prof_end
588 else
589   prof_start "pt_ident_threshold"
590   if ((!DRY)); then
591     python "$PY_PT_THRESH" \
592       --paf "$PAF_A" ${PAF_B:+ $PAF_B} \
593       --pt-ids "${pt_origin:-/dev/null}" \
594       --mt-ids "${mt_origin:-/dev/null}" \
595       --alen-min "$alen_min" \
596       --fpr "$fpr" --tpr "$tpr" \
597       --diag "$PT_DIAG" \
598       --emit-pt-ids "$PT_IDS" >"$PT_VARS"
599   note "PT IDs: $PT_IDS"
600   note "PT Vars: $PT_VARS"
601   fi
602   prof_end
603 fi
604
605 # 1e) Remove PT reads âM-^FM-^R R1
606 note0 "1e) Remove PT reads âM-^FM-^R R1"
607 ident_min_out="$identity_min"
608 if [[ -s "$PT_VARS" ]]; then
609   while IFS== read -r kv; do
610     case "$kv" in
611       ident_min=*) ident_min_out="${kv#ident_min=}" ;;
612       alen_min=*) alen_min="${kv#alen_min=}" ;;
613     esac
614   done <"$PT_VARS"
615 fi
616 note "Step1: PT thresholds âM-^FM-^R ident_min=${ident_min_out} alen_min=${alen_min}"
617 prof_start "write_nonPT"
618 if ((!DRY)); then
619   seqkit fx2tab -ni "$reads" | sort -u >"$outdir/map/all.ids"
620   sort -u "$PT_IDS" >"$outdir/map/pt.ids.sorted"
621   comm -23 "$outdir/map/all.ids" "$outdir/map/pt.ids.sorted" >"$outdir/keep.nonpt.ids"
622   note "Step1: write $outdir/reads.nonpt.fq.gz"
623   seqkit grep -f "$outdir/keep.nonpt.ids" "$reads" -o "$outdir/reads.nonpt.fq.gz"
624 fi
625 prof_end
626 R1="$outdir/reads.nonpt.fq.gz"
627 fi
628 [[ -s "$R1" ]] || {
629   note "ERR missing R1 ($R1)"
630   exit 3
631 }

```

9|M-^M-^T 27, 25 13:53

polap-bash-fast-mtseed-ont.sh

Page 14/20

9|M-[M-^T 27, 25 13:53

polap-bash-fast-mtseed-ont.sh

Page 15/20

```

675     awk 'NR>1{print $1"\t"$3}' "$SCAN_DIR/scan.overlapness.tsv" |
676         sort -k2,2nr |
677         awk -v f="$scan_keep_frac" '{a[NR]=$1} END{lim=int(NR*f); for(i=1;i<=lim;i++) print a[i]}' \
678             >"$SCAN_DIR/shortlist.ids"
679     seqkit grep -f "$SCAN_DIR/shortlist.ids" "$R1" -o "$SCAN_DIR/R1.short fq.gz"
680
fi
681 prof_end
682
683 # 2b) refine streamed  $\hat{M}-^F M-^R$  edges + overlapness
684 note0 "2b) refine streamed  $\hat{M}-^F M-^R$  edges + overlapness"
685 prof_start "refine"
686 if ((!DRY)); then
687     minimap2 -t "$threads" -x ava-ont \
688         -k"$refine_k" -w"$refine_w" --secondary=yes -N "$refine_N" \
689         --mask-level "$refine_mask" --min-occ-floor "$refine_minocc" \
690         "$SCAN_DIR/R1.short fq.gz" "$SCAN_DIR/R1.short fq.gz" |
691         python "$EDGPY" \
692             --min-olen "$min_olen" --min-ident "$min_ident" --w-floor "$w_floor" \
693             --topk-node "$topk_edges_per_node" \
694             --edges "$ST1/edges.tsv" --overlapness "$OTSV_GLOBAL"
695     ((gzip_edges)) && gzip -f "$ST1/edges.tsv"
696
fi
697 prof_end
698 EDGES_GLOBAL="$ST1/edges.tsv${gzip_edges:+.gz}"
699
700 # optional graph summary
701 if ((!DRY)) && [[ -s "$EDGECOMP" ]]; then
702     COMP_DIR="$ST1/comp"
703     mkdir -p "$COMP_DIR"
704     python "$EDGECOMP" "$EDGES_GLOBAL" \
705         --summary "$COMP_DIR/summary.tsv" \
706         --sizes "$COMP_DIR/comp_sizes.tsv" \
707         --membership "$COMP_DIR/node2comp.tsv" \
708         --deg "$COMP_DIR/degree.tsv" \
709         --giant-nodes "$COMP_DIR/giant_nodes.txt" \
710         --giant-edges "$COMP_DIR/giant_edges.tsv" \
711         --print || true
712
fi
713
fi
714
715 # 2c) BUSCO on a 10% length-stratified subsample (QC only)
716 note0 "2c) BUSCO on a 10% length-stratified subsample (QC only)"
717 if [[ -n "$nprot" ]]; then
718     note "Step2: miniprot BUSCO on 10% subsample of PT-depleted reads (QC)"
719     SAMP_IDS="$ST1/nonpt.sample.ids"
720     SAMP_FQ="$ST1/reads.nonpt.sample.fq.gz"
721     # simple random 10% by seed (keeps it fast and reproducible)
722     ((!DRY)) && seqkit sample -s "$seed" -p 0.10 "$R1" -o "$SAMP_FQ"
723     ((!DRY)) && miniprot -d "$ST1/nonpt.sample.mpi" "$SAMP_FQ"
724     ((!DRY)) && miniprot -t "$threads" -S -N 3 --outc 0.4 \
725         "$ST1/nonpt.sample.mpi" "$nprot" > "$ST1/nonpt.sample.busco.paf"

```

9|M-^M-^T 27, 25 13:53

polap-bash-fast-mtseed-ont.sh

Page 16/20

```

726  (!DRY) ) && awk 'BEGIN{FS=OFST="t"} NF>=12 && $11+0>=150 {print $1}' \
727    "$ST1/nonpt.sample.busco.paf" | sort -u >"$ST1/nuc.ids.sample"
728
729 fi
730
731 # Step 3: overlapness + selection (organized)
732 # Step 3: overlapness + selection (organized)
733 # Step 3: overlapness + selection (organized)
734 RDIR="$outdir/round_1"
735 mkdir -p "$RDIR"
736 if _should_run 3; then
737
738   OTSV="$RDIR/overlapness.tsv"
739   note0 "3a) Prepare $OTSV"
740   if [[ "$overlap_mode" == "scan+stream+edges" ]]; then
741     note "Step3: reuse refined overlapness $RDIR/$OTSV"
742     (!DRY) && cp -f "$OTSV_GLOBAL" "$OTSV"
743   else
744     note "Step3: compute overlapness from PAF $RDIR/$OTSV"
745     (!DRY) && python "$OLPY" "$PAFALL" \
746       --min_olen "$min_olen" --min_ident "$min_ident" --w_floor "$w_floor" \
747       >"$OTSV"
748 fi
749
750 # --- Step 3: selection (prefer BUSCO sample labels from Step 2c) ---
751 note0 "3b) selection (prefer BUSCO sample labels from Step 2c)"
752 SELECT_IDS="$RDIR/select_ids.txt"
753
754 # Prefer the miniprot BUSCO sample labels from Step 2c, else fallback to --nuc-ids
755 NUC_SAMPLE="$ST1/nuc.ids.sample"
756 NUC_FOR_THRESH=""
757 if [[ -s "$NUC_SAMPLE" ]]; then
758   NUC_FOR_THRESH="$NUC_SAMPLE"
759 elif [[ -n "$nuc_ids_opt" && -s "$nuc_ids_opt" ]]; then
760   NUC_FOR_THRESH="$nuc_ids_opt"
761 fi
762
763 if [[ -n "$NUC_FOR_THRESH" ]]; then
764   note "Step3: nuclear-guided selection using $NUC_FOR_THRESH -> $SELECT_IDS"
765   if ((!DRY)); then
766     python "$THRESH_NUC_PY" "$OTSV" "$NUC_FOR_THRESH" \
767       --mode fpr --fpr 0.01 --diag "$RDIR/threshold_from_nuclear.tsv" \
768       >"$RDIR/nuc thresh.vars"

```

9|M-^M-^T 27, 25 13:53

polap-bash-fast-mtseed-ont.sh

Page 17/20

```

769
770     wdeg_min=0
771     deg_min=0
772     while IFS== read -r kv; do
773         case "$kv" in
774             wdeg_min=*) wdeg_min="${kv#wdeg_min=}" ;;
775             deg_min=*) deg_min="${kv#deg_min=}" ;;
776         esac
777     done <"$RDIR/nuc_thresh.vars"
778
779     awk -v W="$wdeg_min" -v D="$deg_min" 'BEGIN{FS=OFS="\t"}'
780     NR>1 && ($3+0)>=W && ($2+0)>=D {print $1}' "$OTSV" |
781         sort -u >"$RDIR/organelle.ids"
782
783     # Exclude labeled nuclear reads (from sample or full set)
784     comm -23 <(sort -u "$RDIR/organelle.ids") \
785         <(sort -u "$NUC_FOR_THRESH") \
786         >"$SELECT_IDS"
787 fi
788 else
789     note "Step3: select top-fraction ($top_frac) by wdegree"
790     (!DRY) && python "$TOPPY" "$OTSV" --top_frac "$top_frac" >"$SELECT_IDS"
791 fi
792
793 fi
794
795 # Step 4: assemble selected reads (organized)
796 # Step 4: assemble selected reads (organized)
797 # Step 4: assemble selected reads (organized)
798 SEEDS_ROUND="""
799 if _should_run 4; then
800     if [[ "$Assembler" == "miniasm" ]]; then
801         SUBPAF="$RDIR/top_overlaps.paf.gz"
802         # In streamed mode we do not have a refined PAF; miniasm will still run with -f
803         if [[ "$overlap_mode" != "scan+stream+edges" ]]; then
804             note "Step4: filter PAF by selected ids M-^FM-^R $SUBPAF"
805             (!DRY) && python "$FILTPY" "$PAFALL" "$RDIR/select_ids.txt" -o "$SUBPAF"
806         fi
807         TOPFA="$RDIR/top_reads.fa.gz"
808         note "Step4: extract selected reads M-^FM-^R $TOPFA"
809         (!DRY) && seqkit fq2fa "$R1" | seqkit grep -f "$RDIR/select_ids.txt" -o "$TOPFA"
810         GFA="$RDIR/miniasm.gfa"
811         note "Step4: run miniasm on the selected reads"

```

9|M-^M-^T 27, 25 13:53

polap-bash-fast-mtseed-ont.sh

Page 18/20

9|M-[M-^T 27, 25 13:53

polap-bash-fast-mtseed-ont.sh

Page 19/20

```

855     SUBPAF="$RDIR/top_overlaps.paf.gz"
856     note "Step4: filter full PAF to selected read overlaps ÂM-^FM-^R $SUBPAF"
857     ((!DRY)) && python "$FILTPY" "$PAFALL" "$SELECT_IDS" -o "$SUBPAF"
858
859     note "Step4: run miniasm with filtered overlaps"
860     ((!DRY)) && miniasm -f "$TOPFA" <(zcat -f "$SUBPAF") >"$GFA"
861
862 else
863     # STREAMED MODE:
864     # In Step 2 we *did not* save a full PAF (we streamed directly to reducers).
865     # Miniasm still needs overlaps, so we regenerate them ÂM-^@M-^T but only among
866     # the smaller set of selected reads.
867     SELPAF="$RDIR/selected_allvsall.paf.gz"
868     note "Step4: generate overlaps among selected reads ÂM-^FM-^R $SELPAF"
869     ((!DRY)) && minimap2 -x ava-ont -t "$threads" --secondary=yes -N 30 \
870         --mask-level 0.60 --min-occ-floor 10 \
871         "$TOPFA" "$TOPFA" | gzip -1 >"$SELPAF"
872
873     note "Step4: run miniasm with selected-read overlaps"
874     ((!DRY)) && miniasm -f "$TOPFA" <(zcat -f "$SELPAF") >"$GFA"
875 fi
876
877 # Extract contigs from the GFA into a FASTA for downstream use
878 SEEDS_ROUND="$RDIR/m_seeds_raw.fa"
879 note "Step4: extract contigs from GFA ÂM-^FM-^R $SEEDS_ROUND"
880 ((!DRY)) && awk '/^S/{print ">$2\n"$3}' "$GFA" >"$SEEDS_ROUND" || :
881
882 else
883     # RAVEN MODE:
884     # Raven overlaps + assemblies internally. We only need to provide the reads.
885     TOPFQ="$RDIR/top_reads.fq.gz"
886     note "Step4: extract selected reads (fq) ÂM-^FM-^R $TOPFQ"
887     ((!DRY)) && seqkit grep -f "$SELECT_IDS" "$R1" -o "$TOPFQ"
888
889     SEEDS_ROUND="$RDIR/raven_contigs.fasta"
890     note "Step4: raven --disable-polishing ÂM-^FM-^R $SEEDS_ROUND"
891     ((!DRY)) && raven --threads "$threads" --disable-polishing -o "$SEEDS_ROUND" "$TOPFQ"
892 fi
893
894 # Save path for finalize step
895 SEEDS_CUR="$SEEDS_ROUND"
896 fi
897
898 if _should_run 5; then
899     note "polap readassemble"
900
901     # annotate miniasm assembly contigs
902     #
903     contigger_dir="${RDIR}/30-contigger"
904     mkdir -p "${contigger_dir}"
905     sed 's/LN:i:dp:i/' "${RDIR}/miniasm.gfa" >"${contigger_dir}/graph_final.gfa"

```

9|M-[M-^T 27, 25 13:53

polap-bash-fast-mtseed-ont.sh

Page 20/20

```
906
907     # polap annotate -o "${outdir}" -i round_1
908
909     # select seed contigs
910     #
911     bash "${_POLAPLIB_DIR}/../.polap2.sh" readassemble annotated \
912         -o "${outdir}" -i round_1 \
913         -l "${reads}"
914
915 fi
916
917     # Finalize
918     if ((!DRY)); then
919         FINAL_DIR="$(dirname ${SEEDS_CUR})"
920         FINAL_SEEDS="$FINAL_DIR/m_seeds_final.fa"
921         note "Finalize: cp ${SEEDS_CUR} ${FINAL_SEEDS}"
922         cp "${SEEDS_CUR}" "$FINAL_SEEDS"
923         note "Done. Final seeds: $FINAL_SEEDS"
924     else
925         note "--dry finished (no commands executed)."
926     fi
```