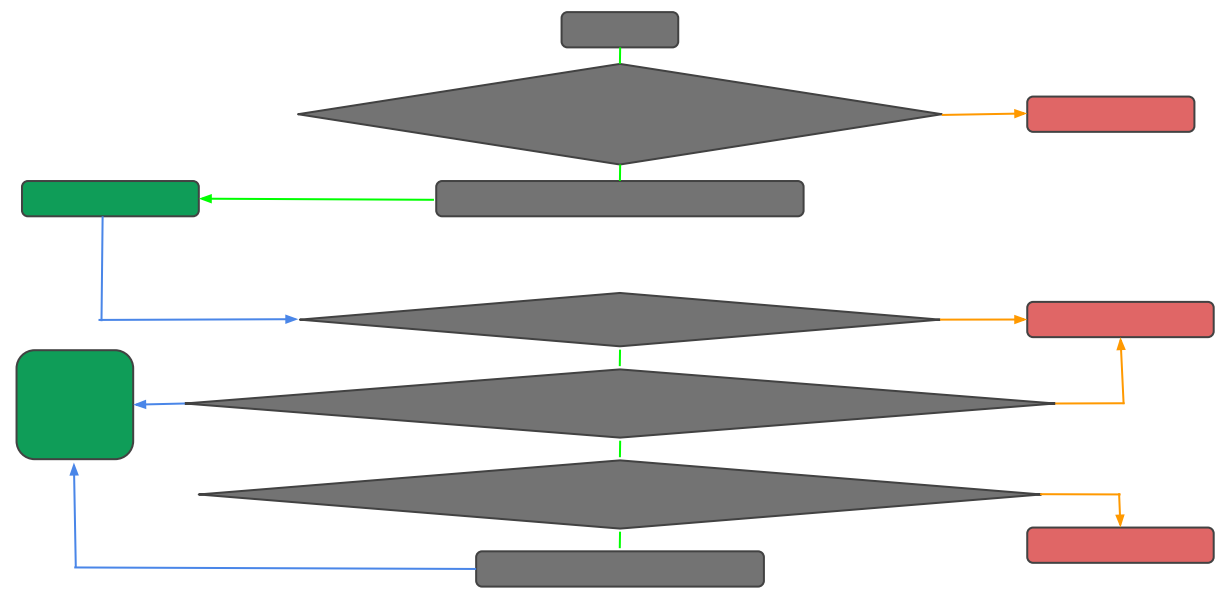


Unboxing CCS2

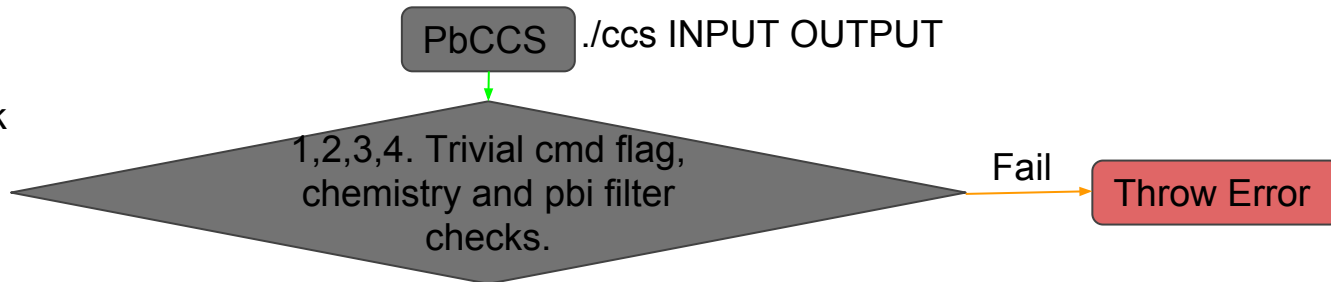
-- Avi Srivastava

- ← Single Thread
- ← Multi Thread
- ← Error / Failcheck

Flowchart



← Single Thread
← Multi Thread
← Error / Failcheck

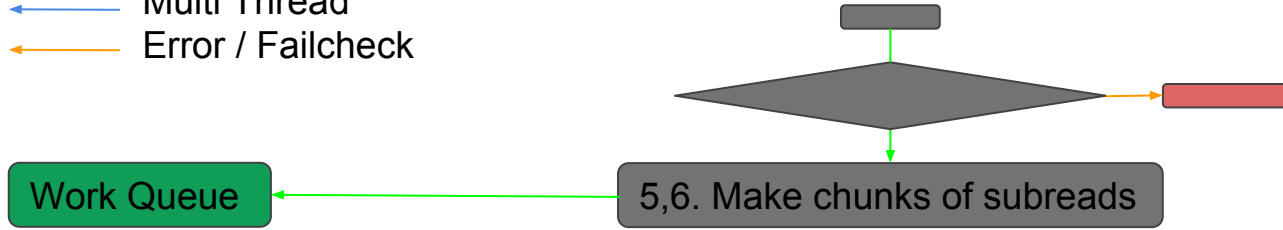


Check for command line arguments:

- minPasses ≥ 1
- Input / Output format
- Too many arguments
- Output file already exist
- Option --byStrand not compatible with --noPolish and similar-ish...

Directly throws Error without doing anything.

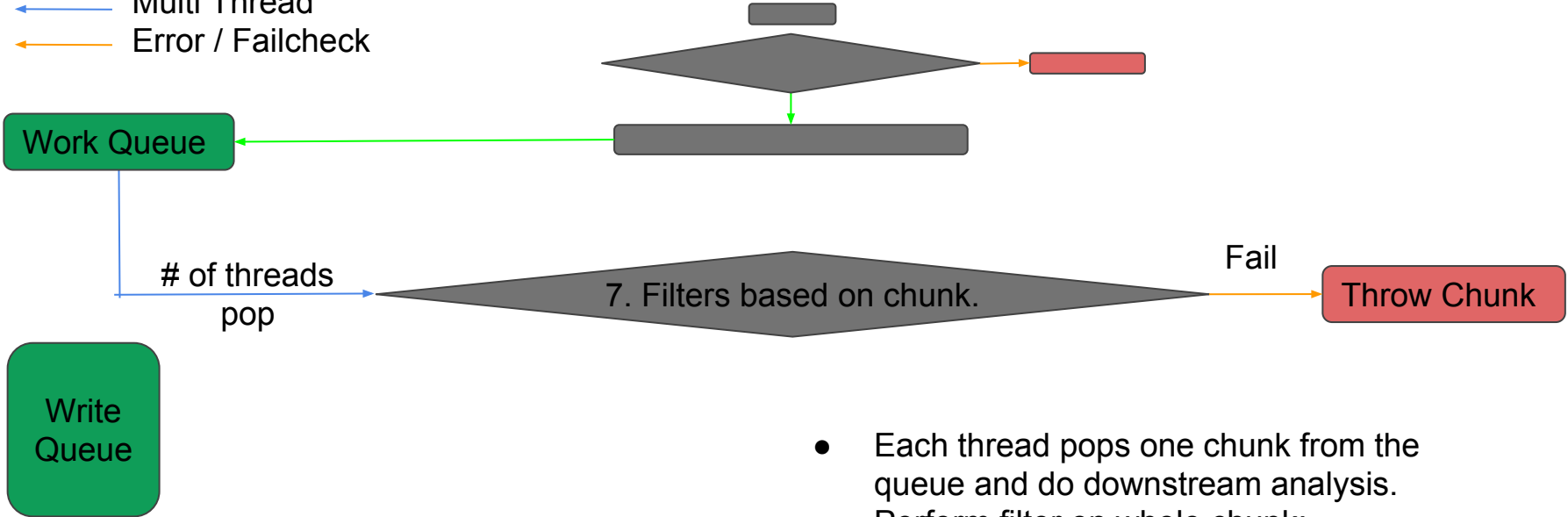
← Single Thread
← Multi Thread
← Error / Failcheck



- Initiate work / write queue based on # of threads.
- Check required output format and initialize output stream accordingly (can be FASTQ / BAM).
- Subsample all the subreads from one ZMW and push them in worker queue.

Write
Queue

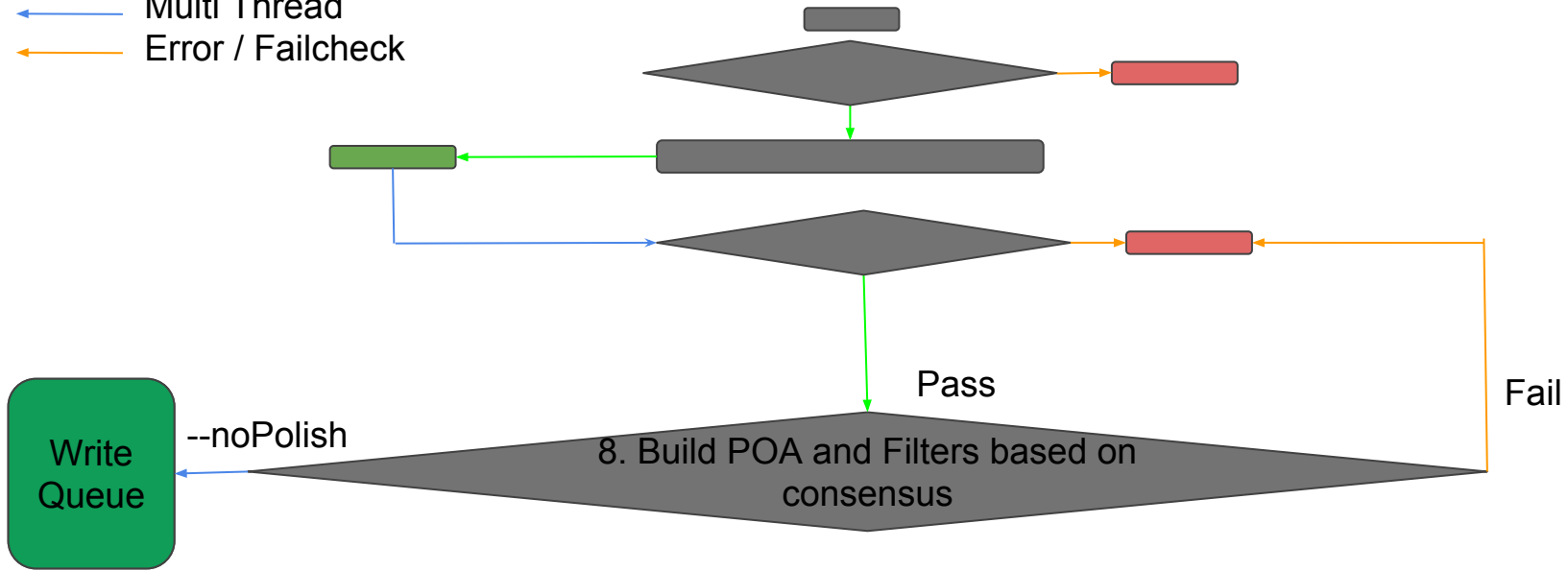
← Single Thread
← Multi Thread
← Error / Failcheck



Write
Queue

- Each thread pops one chunk from the queue and do downstream analysis.
- Perform filter on whole chunk:
 - $\text{Min SNR}(A,C,T,G) > \text{SNR thresh.}$
 - $\text{Median Read Length} < 10$
 - $\# \text{ Passes} < \text{Pass thresh.}$
- Perform filter on each read:
 - $\text{Read quality} < 0.65$
 - $\text{size}(\text{Read}) > 2 * \text{Median read len.}$
- Sort remaining reads based on deviation from median and build POA.

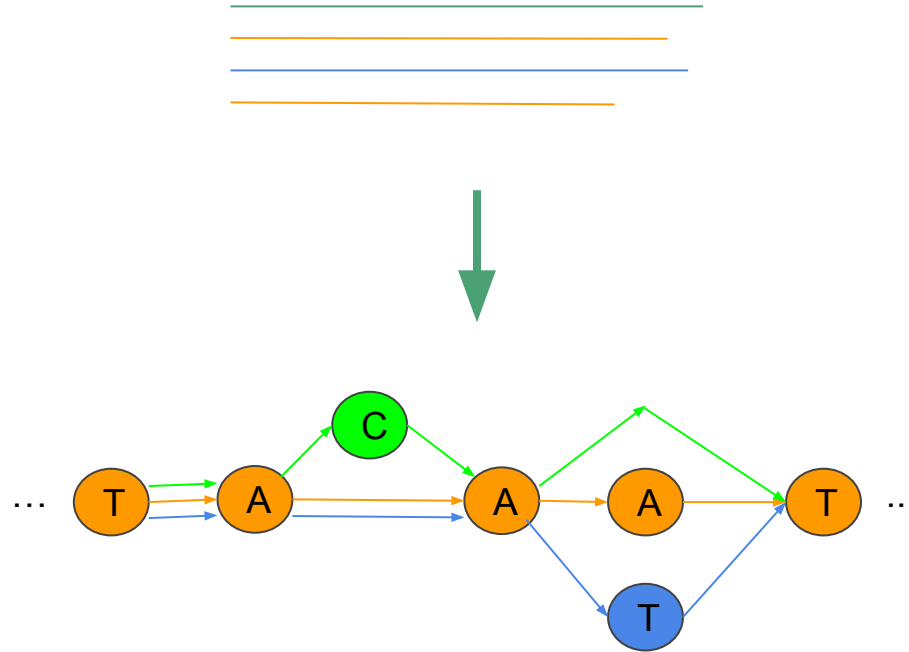
← Single Thread
← Multi Thread
← Error / Failcheck



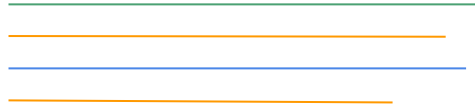
POA (Partial Order Alignment)



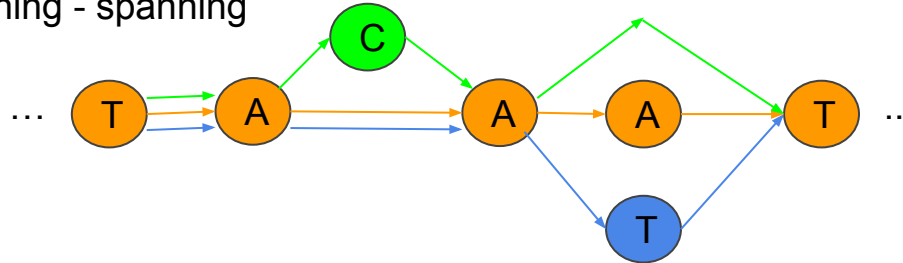
POA (Partial Order Alignment)



POA (Partial Order Alignment)



Scoring: 2*containing - spanning



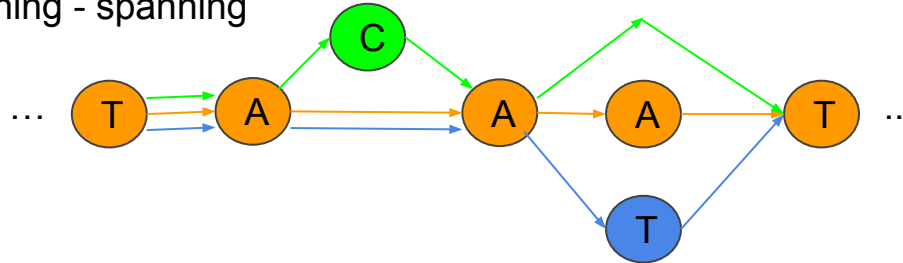
Score: $(2*3-3)+(2*2-3) + (2*1-3)+(2*1-1) + \dots$

POA (Partial Order Alignment)

Subreads



Scoring: 2*containing - spanning

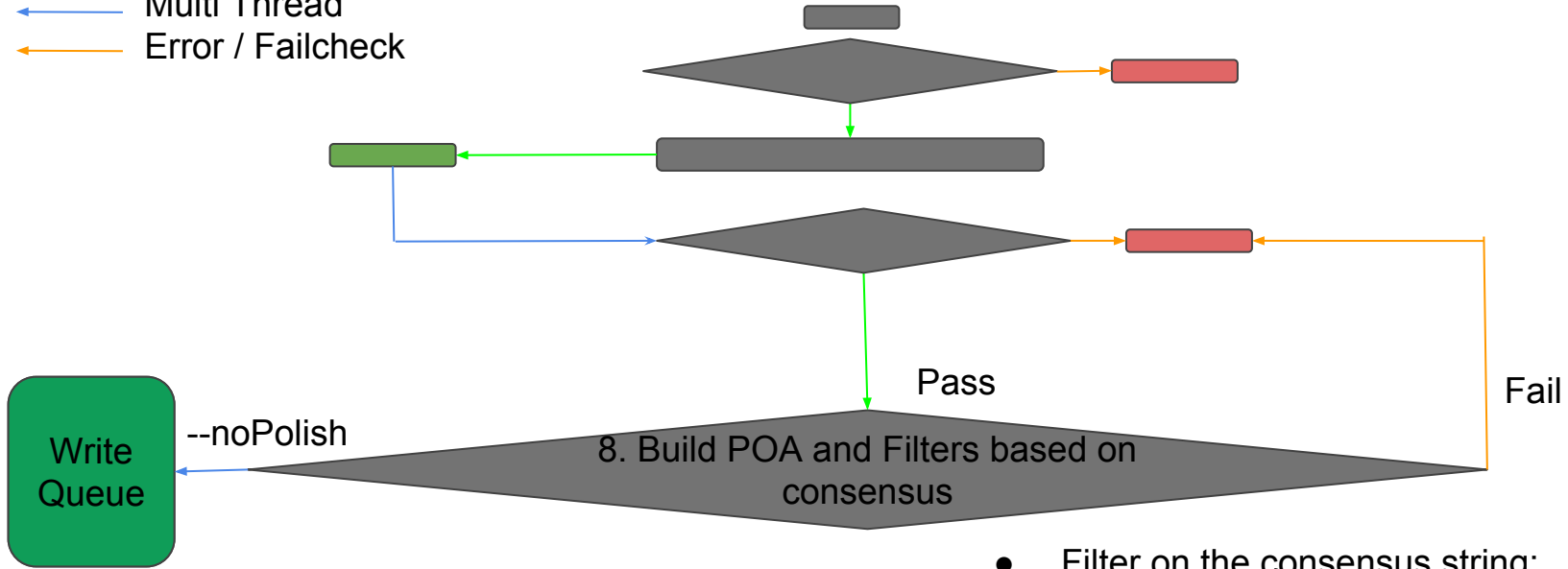


Score: $(2*3-3)+(2*1-3)$ + $(2*1-3)$ + ...

Score: $(2*3-3)+(2*2-3)$ + $(2*1-3) + (2*1-1)$ + ...

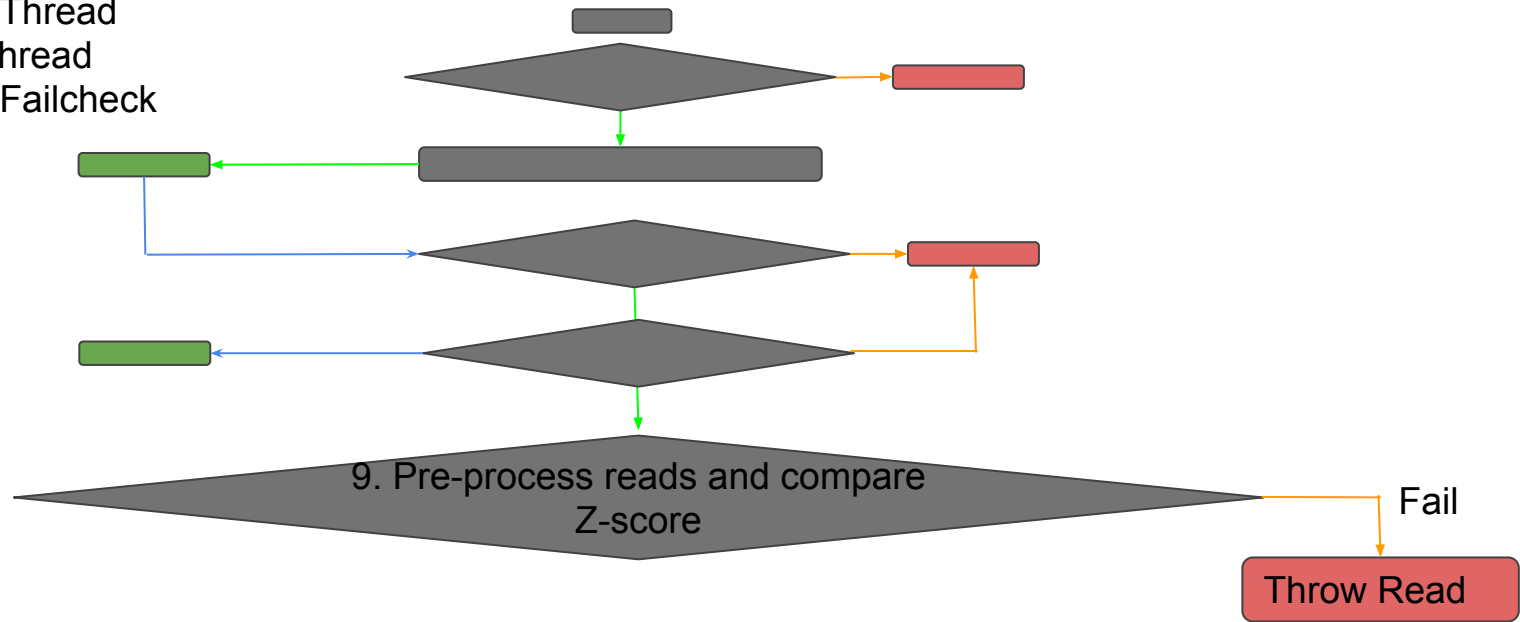
Score: $(2*3-3)+(2*2-3)$ + $(2*1-3) + (2*1-1)$ + ...

← Single Thread
← Multi Thread
← Error / Failcheck



- Filter on the consensus string:
 - Throw chunk if consensus is too short or long.**
- Based on --noPolish push the consensus to write queue and exit.
- Otherwise go to next stage of Polishing.

← Single Thread
← Multi Thread
← Error / Failcheck



- **Zscore:** How likely is the probability of; the read must be coming from the consensus sequence under pacbio model.

Z-Score

$P(T / R)$: Posterior probability of template given a set of reads.

We have to maximize this posteriori for an unknown template given a set of subreads.

$$\text{Bayes: } P(T / R) = P(R / T) * P(T) / P(R)$$

To maximize $P(T / R)$ we have to maximize $P(R / T)$

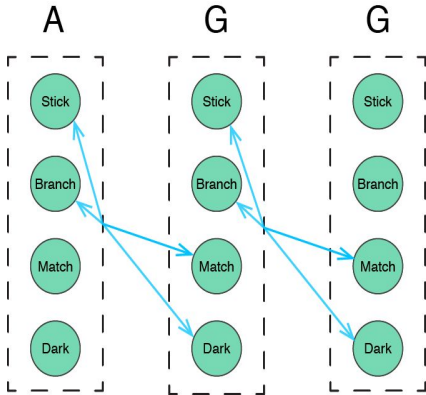
$P(R / T) = \text{Product}(P(R_i / T))$ where R_i is a read from set R

$$P(R_i / T) = \text{Sum}(P(R_i, A / T))$$

Where A represents all the variations of the read

Instead of Summing we can just look for most probable alignment A for each read using Forward Backward Algorithm.

Based on the chemistry we have pre-learned emission and transition probabilities.



Polish

POA consensus

Subreads

Likelihood: Sum of ZScores of all the subreads. Say here it is **LL**

Look for all the mutation of the positions in all reads in the set of subreads.

111

112

113

Set of Mutations
 $\|x\| > LL$

Sort Mutations on ZScore and remove overlapping mutation in window of 10.

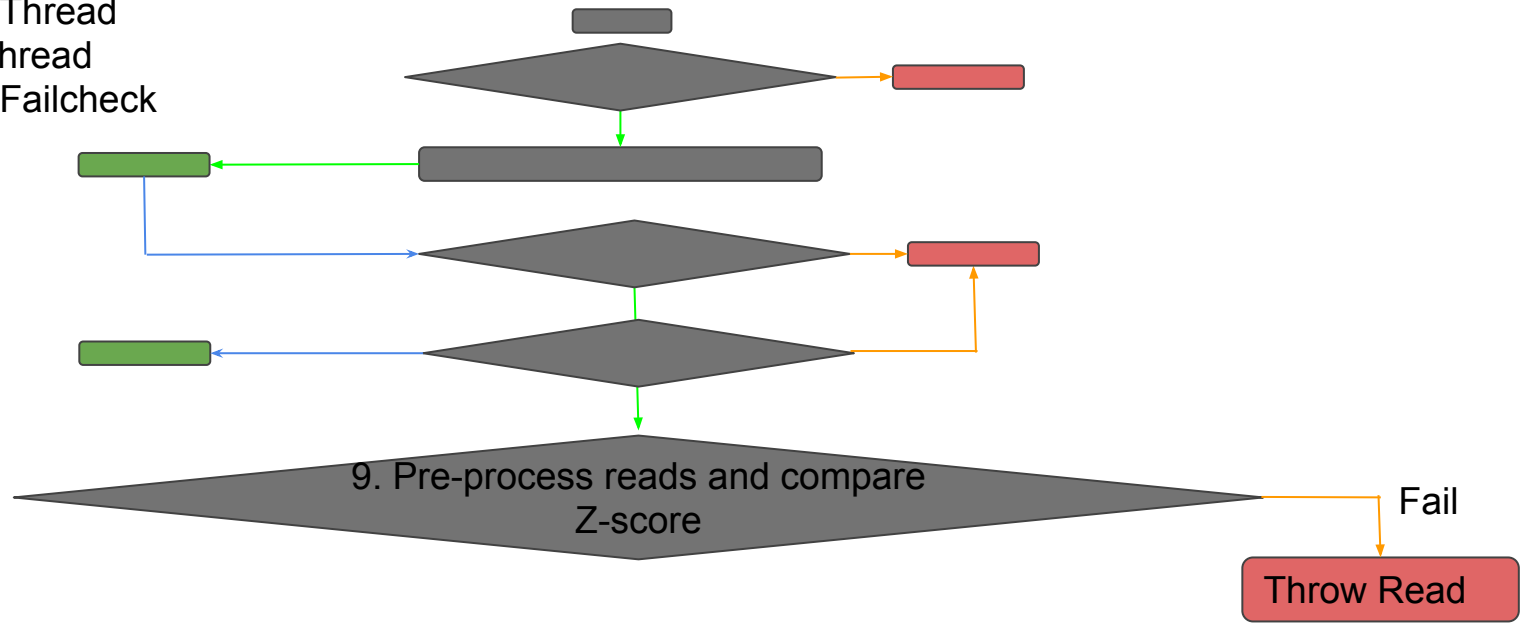
Apply

1122 + 1199

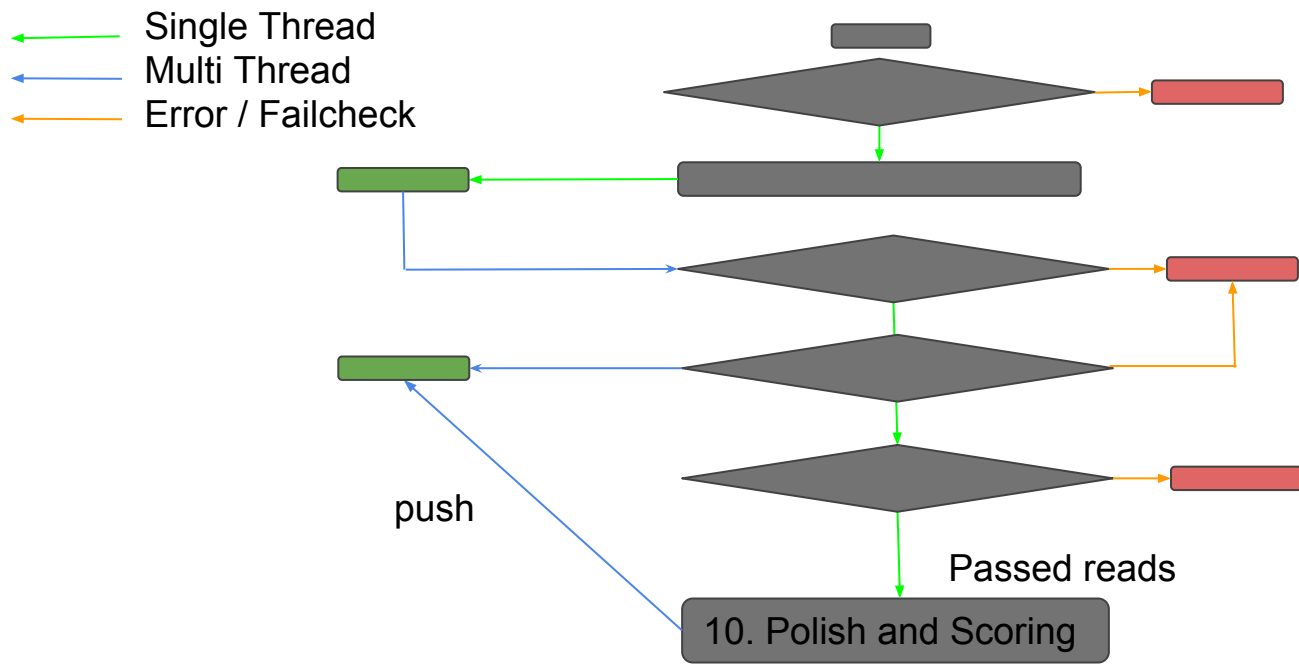
Look for
mutation in
window of 20
and iterate 40

time

- ← Single Thread
- ← Multi Thread
- ← Error / Failcheck



- **Zscore:** How likely is the probability of; the read must be coming from the consensus sequence under pacbio model.
- Throw read:
 - $Zscore < 12.5$
- Throw chunk:
 - $\# \text{ passes} < \text{thresh. Passes}$



← Single Thread
← Multi Thread
← Error / Failcheck

