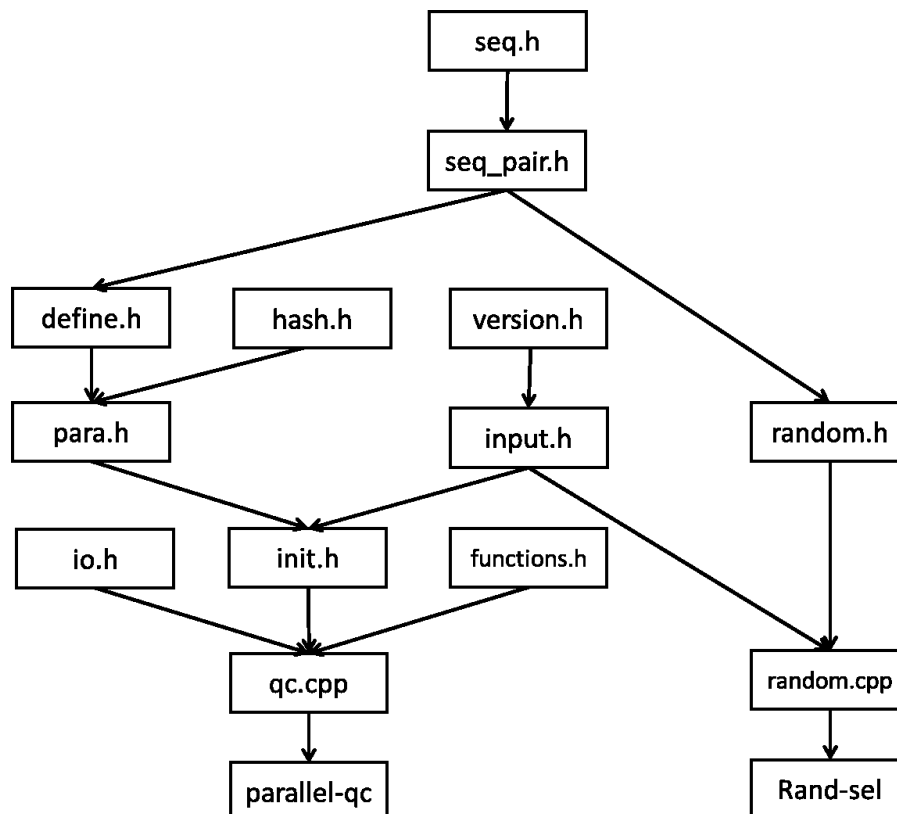


# Parallel-QC development manual

Version 1.0

Parallel-QC is implemented by Linux C++ with POSIX thread. This manual is for further development of Parallel-QC or invocated into other pipeline.

## 1 PART I. PROJECT STRUCTURE



There are 13 source code files for the project.

*seq.h*: The definition of class Seq. Details can be found in Part II.

*seq\_pair.h*: The definition of class Seq\_pair. Details can be found in Part II.

*define.h*: The extern declaration of all global variables of parallel-qc.

*hash.h*: The definition of hash function.

*version.h*: The definition of version of parallel-qc.

*para.h*: The parameter parser of parallel-qc.

*input.h*: The sequence input functions for parallel-qc and rand-sel.

*io.h*: The I/O functions for parallel-qc.

*init.h*: The initialization functions for parallel-qc.

*functions.h*: The quality control functions for parallel-qc.

*qc.cpp*: The entrance of parallel-qc.

*random.h*: The random selection functions for rand-sel.

*random.h*: The entrance of rand-sel.

## **2 PART II. CLASSES**

### **2.1 class Seq**

The basic class of a sequence with quality value.

*public*:

*Seq(string \_seq, string \_qual)*; The construction function with sequence *\_seq* and quality values *\_qual*.

*void Trim(int \_begin, int \_end)*; To get the trimmed sub-read from *\_begin* to *\_end*.

*void Qual\_trim(int \_qual, float \_p)*; To do the quality trim with quality value *\_qual* and proportion of *\_p*.

*void Primer\_trim(string \*\_primers, int n)*; To do the tag-sequence trim with *n* tag sequences *\*\_primers*.

*void GC\_trim(float \_p\_min, float \_p\_max)*; To do the GC proportion trim with minimum GC proportion *\_p\_min* and maximum GC proportion *\_p\_max*.

*bool Is\_qual\_trim()*; If this sequence is trimmed by quality trimming.

*bool Is\_primer\_trim()*; If this sequence is trimmed by tag-sequence trimming.

*bool Is\_gc\_trim()*; If this sequence is trimmed by GC proportion trimming.

*bool Is\_dup\_trim();* If this sequence is trimmed by duplication trimming.

*bool Is\_seq;* If this sequence is trimmed.

## **2.2 class Seq\_pair**

The pair-ended sequences based on class Seq.

*public:*

*Seq\_pair(Seq \_seq\_1, Seq \_seq\_2);* The construction function with 2 Seq objects.

*void Add\_seq\_1(Seq \_seq\_1);* Add sequence\_1 to this sequence pair.

*void Add\_seq\_2(Seq \_seq\_2);* Add sequence\_2 to this sequence pair.

*void Add\_qual\_1(string qual\_1);* Add the quality values to sequence\_1.

*void Add\_qual\_2(string qual\_2);* Add the quality values to sequence\_2.

*void Out\_put(string \_label, ostream \*out\_1, ostream \*out\_2, bool format);*

Output the sequence to stream \*out\_1 for sequence\_1 and \*out\_2 for

sequence\_2. If format= True the output format is in Fastq, else in Fasta.

## **3 PART III. VARIABLES AND FUNCTIONS**

### **3.1 Global variables for parallel-qc**

*bool \_in\_pair;* If the input is pair-ended.

*bool \_single\_input ;* If the input is in two files.

*bool \_keep\_pair;* If the output is kept in pair-ended.

*bool \_dup;* If the duplication trimming is processed.

*bool \_ext\_qual;* If there are extern quality values.

*bool \_in\_format ;* Input format, true for Fastq, and false for Fasta.

*bool \_out\_format = true;* Output format, true for Fastq, and false for Fasta.

*int \_thread* ; Thread number.

*vector <string> order*; Vector to keep the sequence order by recording the label.

*hash\_map <string, Seq\_pair, std\_string\_hash> seq\_pairs*; Hash\_map to store all sequences indexed by the sequence label.

### **3.2 Functions for quality control**

*void Input\_reads\_pair()*; Input the raw reads of pair-ended format.

*void Input\_reads()*; Input the raw reads of single-ended format.

*void Input\_ext\_qual\_pair()*; Input the quality values of pair-ended format for Fasta.

*void Input\_ext\_qual()*; Input the quality values of single-ended format for Fasta.

*void Output\_analysis\_report()*; Output the analysis report

*void Output\_reads\_pair()*; Output the analysis results in pair-ended format.

*void Output\_reads()*; Output the analysis results in single-ended format.

*void Para\_args()*; Parse the arguments.

*void Init()*; Initialization.

*void Trim()*; Trimming processes.

*void Drop\_dup()*; Drop the duplication reads.

## **4 PARALLEL COMPUTING**

Parallel-QC parallelizes the trimming processes and distribute into different threads based on POSIX thread technology.

*void \* Trim\_parallel(void \* args)*; The function pointer for multi-thread trimming processes.

*void \* Drop\_dup\_1(void \*)*; The function pointer for parallel duplication trimming of sequence\_1 for pair-ended data.

*void \* Drop\_dup\_2(void \*)*; The function pointer for parallel duplication trimming of sequence\_2 for pair-ended data.

## **5 CONTACT US**

Any problem please contact

Xiaoquan Su

suxq@qibebt.ac.cn