

# 1 PRINSEQ++, a multi-threaded tool for fast 2 and efficient quality control and 3 preprocessing of sequencing datasets

4 Vito Adrian Cantu<sup>1</sup>, Jeffrey Sadural<sup>2</sup>, and Robert Edwards<sup>2</sup>

5 <sup>1</sup>Computational Science Research Center, San Diego State University, San Diego,  
6 92182, USA.

7 <sup>2</sup>Department of Computer Science, San Diego State University, San Diego, 92182, USA.

8 Corresponding author:

9 Robert Edwards<sup>1</sup>

10 Email address: redwards@mail.sdsu.edu

## 11 ABSTRACT

12 PRINSEQ++ is a C++ implementation of the very popular software prinseq-lite for quality control and  
13 preprocessing of sequencing datasets. PRINSEQ++ can run multi-threaded processes, which makes  
14 it more than 10 times faster than the original version. It can read from, and write to, compressed  
15 files, drastically reducing the use of hard-drive. PRINSEQ++ can filter, trim and reformat sequences  
16 by a variety of options to improve downstream analysis. PRINSEQ++ is freely available on GitHub  
17 (<https://github.com/Adrian-Cantu/PRINSEQ-plus-plus>) and runs on all Unix-like systems.  
18

## 19 1 INTRODUCTION

20 As prices fall, high-throughput sequencing is being used in new ways and areas such as personalized  
21 medicine (Wilson et al., 2014) and recreational genomics (Evans, 2008). This brings about novel  
22 challenges for the techniques we use to analyze and draw conclusions from sequencing data, in particular  
23 speed and scalability.

24 Quality control is a crucial step in the analysis of sequencing datasets as low-quality sequences,  
25 sequence contamination and artifacts can eventually lead to erroneous conclusions. Most applications  
26 for quality control and preprocessing are written in high level programming languages such as Perl  
27 (prinseq-lite (Schmieder and Edwards, 2011)) or Java (fastQC (Andrews, 2010)) which are slower to  
28 execute and provide limited multi-threading support.

29 Since its publication in in early 2011, prinseq-lite has been cited more than 1500 times and downloaded  
30 more than 54000 times. In the same time interval, the number of bases in the Sequence Read Archive  
31 has grown 247x (from 74 Tbp to 18412 Tbp). It is clear that a new tool, one that has the usefulness of  
32 prinseq-lite while being drastically faster, is needed.

33 PRINSEQ++ implements all the functionality of the Prinseq-lite tool, adds some new features, but  
34 can run 16x times faster as it is written in C++ and can take advantage of multi-threading.

## 35 2 NEW FEATURES

### 36 2.1 Sequence duplication

37 Sequence duplication occur at different steps of the sequencing protocol (Gomez-Alvarez et al., 2009).  
38 Traditionally, duplicated sequences are hard to detect as the naive approach is to compare every single  
39 sequence to each other. This is problematic in a multi-threaded environment were each thread holds in  
40 memory a few sequences at most and cross talk between threads needs to be minimum in the interest of  
41 speed.

42 PRINSEQ++ uses a probabilistic data structure, Bloom filter (Bloom, 1970)(Partow, 2010), to identify  
43 duplicate sequences. A Bloom filter is bit-array where every sequence is transformed by several fast

44 (non-cryptographic) hash functions and the corresponding bits turned on. To see if a sequence is already  
 45 in the filter (if it is duplicated) one only need to check the corresponding bits. Reading and writing from a  
 46 bloom filter is fast and can be done asynchronously.

## 47 2.2 Parallelization

48 Most parallelization models, like OpenMP, MPI or Cuda, require the user to know the size and shape  
 49 of the input a priori. Counting the number of sequences in a FASTA or FASTQ file requires reading it  
 50 completely, which is slow and non trivial (especially so for compressed files). PRINSEQ++ uses POSIX  
 51 threads (pthreads), an application programming interface (API) designed to allow maximum freedom to  
 52 developers of multi-threaded applications.

53 With the exception of sequence duplication, the quality control and preprocessing is independent for  
 54 each sequence pair. Each thread performs all necessary operations on one sequence pair at the time.  
 55 This includes: reading from file, uncompressing if necessary, checking for duplicates and other filters,  
 56 compressing if necessary, and writing to the corresponding output file if the read pair passes all filters.  
 57 This model drastically reduces run-time, is input size agnostic, and uses little memory.

## 58 2.3 Speedup

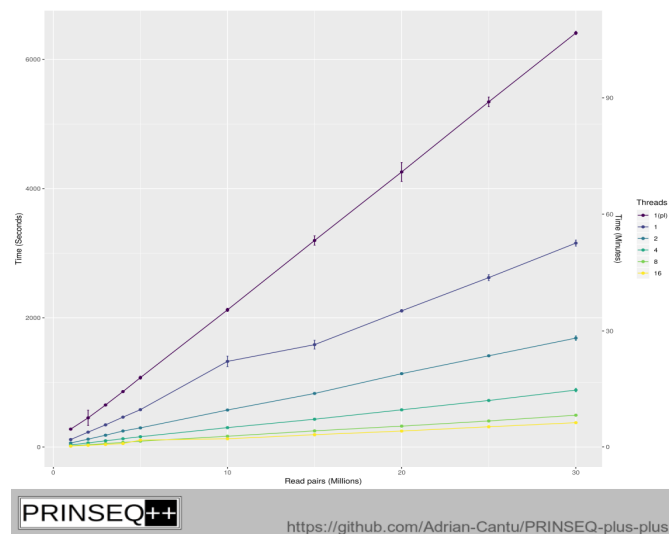
Threads	speedup over prinseq-lite	speedup over PRINSEQ++
1	1.98 x	1 x
2	3.77 x	1.93 x
4	7.26 x	3.7 x
8	12.96 x	6.62 x
16	16.47 x	8.39 x

**Table 1.** Speedup of multi-threaded PRINSEQ++ over prinseq-lite and single-threaded PRINSEQ++ on the same computer, for the same dataset, using different numbers of threads.

59 To assess the effect of increasing the number of threads on speed and the speedup of PRINSEQ++  
 60 over prinseq-lite, we measured run-time of prinseq-lite and PRINSEQ++ on several FASTQ pair files of  
 61 different sizes. A pair of FASTQ files from a metagenomic sample were downloaded from the sequence  
 62 read archive (Run:SRR7091319). The FASTQ files were cut into files of 1, 2, 3, 4, 5, 10, 15, 20, 25, 30  
 63 millions read pairs. PRINSEQ++ and prinseq-lite were run on those files with equivalent filtering options  
 64 ("min\_len 100 -min\_gc 40 -max\_gc 60 -lc\_method entropy -lc\_threshold 90" for prinseq-lite and "-min\_len  
 65 100 -min\_gc 40 -max\_gc 60 -lc\_entropy=0.9" for PRINSEQ++).

66 Run-time was measured using GNU time 1.7 on a 24 cores Intel Xeon CPU X5650 running at  
 67 2.67GHz with 189Gb of RAM. Each measurement was done three times and the mean time and 0.95  
 68 confidence interval were plotted on figure 1. Table 1 shows the speedup of multi-threaded PRINSEQ++  
 69 over prinseq-lite and over single-threaded PRINSEQ++.

70 There are two main reasons why the speedup don't scales linearly with the number of threads for  
 71 PRINSEQ++. There is a small overhead in creating a thread and, more importantly, input and output files  
 72 need to be accessed synchronously. A thread cannot write or read if another thread is doing so, and must  
 73 wait for it to finish. As the number of threads increases this happens more often and more time is spent  
 74 waiting for access to files. Additionally, there is little advantage in using more threads than the number of  
 75 cores in the CPU, as this will cause multiple threads to run on the same core and share execution cycles.



**Figure 1.** Runtime comparison, Run-time of prinseq-lite and PRINSEQ++ was measured on several FASTQ pair files of different sizes with equivalent options. PRINSEQ++ was run with different number of threads, prinseq-lite single-threaded. Mean speedup of PRINSEQ++ over prinseq-lite(pl) is 1.98x for single core and 16.47x for 16 cores. Error bars use a 0.95 confidence interval.

### 76 3 CONCLUSION

77 PRINSEQ++ is a fast and efficient and can significantly reduce the run-time of sequencing datasets  
 78 analysis. This is critical in applications that are time-sensitive or where the amount of data is so large that  
 79 slower method are not feasible. PRINSEQ++ has the capacity or reading from, and writing to compressed  
 80 files without ever uncompressing the whole file, this drastically reduces use hard-drive use. PRINSEQ++  
 81 emulates prinseq-lite syntax, thus it can be easily added to any pipeline currently using prinseq-lite.

### 82 ACKNOWLEDGMENTS

83 We thank the PRINSEQ++ users for comments and suggestions and Dr. James Otto for computational  
 84 support.

### 85 REFERENCES

- 86 Andrews, S. (2010). A quality control tool for high throughput sequence data.  
 87 <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>.
- 88 Bloom, B. H. (1970). Space/time trade-offs in hash coding with allowable errors. *Commun. ACM*,  
 89 13(7):422–426.
- 90 Evans, J. P. (2008). Recreational genomics; what's in it for you? *Genetics in medicine : official journal of*  
 91 *the American College of Medical Genetics*, 10(10):709–710.
- 92 Gomez-Alvarez, V., Teal, T. K., and Schmidt, T. M. (2009). Systematic artifacts in metagenomes from  
 93 complex microbial communities. *The ISME Journal*, 3(11):1314–1317.
- 94 Partow, A. (2010). General purpose hash function algorithms.  
 95 <http://www.partow.net/programming/hashfunctions/index.html>.
- 96 Schmieder, R. and Edwards, R. (2011). Quality control and preprocessing of metagenomic datasets.  
 97 *Bioinformatics*, 27(6):863–864.
- 98 Wilson, M. R., Naccache, S. N., Samayoa, E., Biagtan, M., Bashir, H., Yu, G., Salamat, S. M., Somasekar,  
 99 S., Federman, S., Miller, S., and et al. (2014). Actionable diagnosis of neuroleptospirosis by next-  
 100 generation sequencing. *New England Journal of Medicine*, 370(25):2408–2417.