



Quack: A quality assurance tool for high throughput sequence data

Adam Thrash, Mark Arick II, Daniel G. Peterson*

Institute for Genomics, Biocomputing & Biotechnology, Mississippi State University, Mississippi State, MS 39762, USA

ARTICLE INFO

Keywords:

Quality assurance
Visualization
FASTQ
Quality check
DNA

ABSTRACT

The quality of data generated by high-throughput DNA sequencing tools must be rapidly assessed in order to determine how useful the data may be in making biological discoveries; higher quality data leads to more confident results and conclusions. Due to the ever-increasing size of data sets and the importance of rapid quality assessment, tools that analyze sequencing data should quickly produce easily interpretable graphics. *Quack* addresses these issues by generating information-dense visualizations from FASTQ files at a speed far surpassing other publicly available quality assurance tools in a manner independent of sequencing technology.

Introduction

Twenty-first century sequencing techniques have increased biomolecular data production at a rate outpacing Moore's Law [9]. However, with the staggering proliferation of sequence data comes many issues including the need for faster, more intuitive quality control assessment prior to downstream data analysis [4]. Also, with the introduction of Pacific Biosciences and Oxford Nanopore Technologies sequencing platforms, the quality and length of sequence reads have become more varied, making platform specific tools designed for short read sequences [e.g. [3]] of limited utility.

Several principles underlie good quality assessment tools [2,13]. First, sequence quality data needs to be presented in a graphical format to allow rapid assessment. While numerical summary statistics can be extremely useful, simple graphs can expose fundamental problems that would not be detected using summary statistics alone (Fig. 1). Second, the visualizations need to be data-dense, meaning a large amount of information needs to be displayed in a relatively small area [13]. By displaying all of the information within the common eyespan, the data can be assessed more quickly and ergonomically. Third, the tool should be sequencing technology independent. Finally, the tool should be fast and should scale to increasingly larger data sets.

There are several available tools designed to assess the quality of sequence data. Ref [1] notes, “*The undisputed champion of quality control visualization is a tool named FastQC developed by Babraham Institute.*” However, Ref [1] adds, “*Even though it is a de-facto standard of visualization, its results are not always the simplest to interpret.*” While FastQC does follow most of the design principles outlined above, it is still

relatively labor intensive, especially when examining the quality results for large datasets [1,5]. FastQC is written in Java [3].¹

Fastqp [11] is another tool for quality assessment of FASTQ files. Though not as widely adopted as FastQC, Fastqp generates high quality graphics. Fastqp is written in Python.

To improve the speed and ease of quality assessment of FASTQ files, we developed the *Quack* algorithm. *Quack* is written in C which accounts for much of its speed. We demonstrate the value of *Quack* by comparing its speed and output to FastQC and Fastqp. Our findings indicate that *Quack* readily outperforms the other tools.

Methods

Implementation

We chose to write *Quack* in the programming language C for several reasons. First of all, C has been in use for > 40 years, yet, if anything, its relevance is increasing; e.g., TIOBE (www.tiobe.com), a company that accesses the quality and usage of software, named C its 2017 Programming Language of the Year [12]. Operating systems that are based, in full or part, on C include UNIX, Windows, Linux, GNU, and Macintosh, while mobile devices that use C kernels include iOS, Android, and Windows Phone. The most popular databases – Oracle, MySQL, MS SQL Server, and PostgreSQL – are primarily written in C, most embedded systems (e.g., programs that run appliances, vehicles, etc.) are written in C, and the majority of powerful supercomputers run the Linux (C-based) kernel. C has a relatively small, standardized vocabulary and a short runtime which have made it the language of choice

* Corresponding author.

E-mail address: dp127@msstate.edu (D.G. Peterson).

¹ Recently, a tool called FQC was developed to parse FastQC's HTML output into a more interactive format [5]. However, while FQC increases the utility of FastQC, it is ultimately dependent upon (and ostensibly limited by) the FastQC algorithm, and thus it is not considered separately from FastQC.

متن کامل مقاله

دریافت فوری ←

ISIArticles

مرجع مقالات تخصصی ایران

- ✓ امکان دانلود نسخه تمام متن مقالات انگلیسی
- ✓ امکان دانلود نسخه ترجمه شده مقالات
- ✓ پذیرش سفارش ترجمه تخصصی
- ✓ امکان جستجو در آرشیو جامعی از صدها موضوع و هزاران مقاله
- ✓ امکان دانلود رایگان ۲ صفحه اول هر مقاله
- ✓ امکان پرداخت اینترنتی با کلیه کارت های عضو شتاب
- ✓ دانلود فوری مقاله پس از پرداخت آنلاین
- ✓ پشتیبانی کامل خرید با بهره مندی از سیستم هوشمند رهگیری سفارشات