

A compression algorithm of fastq file based on distribution characteristics analysis

Abstract—With the continuous development of sequencing technology scientists in the cost of DNA sequencing in reduce gradually, it also makes the number of DNA sequencing data to increase substantially. While the genome data is need to store, the traditional computer room has not enough to store such large data. Therefore, more and more genome data need to be uploaded to the cloud. Due to the speed of growth of communication have been much faster than the growth of the genomic data, so it is particularly important for genome data compression to reduce the cost of scientific research institutions and it is of great significance to speed up the sharing of genomic data. Fastq file is an important format of genomic data, and now the compression algorithm for fastq files is mainly include of DSRC, FQC, etc. These algorithms are also compressed based on the characteristics of fastq files. In order to improve the rate of compression, we propose a new algorithm of DDSRC and establish the statistical models for the distribution characteristics of strings in fastq files to make a more efficient compression algorithm. This paper accordingly explain the algorithm based on their distribution and characteristics analysis and compare the results with other compression algorithms.

The EOZComb algorithm has the ability to reduce the ratio of fastq files. **Keywords**—genome data; fastq; Compression algorithm; DDSRC;

I.

High-throughput DNA sequence data generated by next generation sequencing (NGS) technologies have brought tremendous stress in data storage and transmission, which poses new problems in memory and storage capabilities [1]. Most fastq file compression tools are running on a Linux platform, according to whether based on specific genetic template to compress or not, it can be divided into two kinds in general. One kind is the non-reference genome compression tool, such as FQZ_Comp, FQC and DSRC, the other is a reference genome compression tool, such as DNAZip. This project designed a common compression tool for fastq file, and the reference genome compression algorithm facing a big problem is that there is no ready reference genome often for comparison. Due to the lack of target genome sequence which is similar to reference genome, lead to the situation that there is a great compression effect but not to be a common solution. Therefore, we adopted the compression strategy of non-reference genome, and compared with the non-reference genome algorithms of DSRC, FQZ_Comp and FQC.

The DSRC method treats the FASTQ file as a file consists of sequence, quality, header and trailer. It is a hierarchical structure, which divides the data into several record block and 512 record blocks as super blocks. These parameters can be adjusted according to specific files. And each superblock compressed by independent super block as a unit will also record the statistical information of each block. It is that if you want to unzip a block, the first location to the record belongs to super block. It computes the block addresses in the superblock and read all records of sharing information, then can unzip record data in a block. The more Randomness of the DSRC for file to be written to improve the speed of compression. Therefore, to some extent at the expense of the compression rate. In general, the theory of DSRC maximum compression rate in the 5-1, sacrifice the part of the compression efficiency to improve the speed of compression.

ting of four parts, a line as a part [2]. The four parts are: the title line, the sequence of bases in the third row and line quality, especially the third lines without compression processing. The following is the main process of the DSRC algorithm. Firstly the header row as several serial string with a delimiter, according to the characteristics of the title string, using the corresponding coding technology that can get the maximum compression including Huffman coding or differential coding. The next is the transfer of additional symbols from DNA data to quality data in addition to the four bases of ATCG. Then analyze the statistical model of statistic sequence, and according to the statistical model and the range of the quality of the base, adopt corresponding encoding, mainly including Huffman coding and coding directly. After that can be used to make code for the secondary compression. Last is the compression of line quality, according to the characteristics of the quality of line of the data, that is quality line from 0 to 40 plus the offset of the data of 33 scope, taking some processing, such as removing possible character that may appear in the end. After processing the data and reunification of Huffman encoding, or directly using the run-length coding. At the same time, the DSRC proposes a hierarchical structure, which divides the data into several record block and 512 record blocks as super blocks. These parameters can be adjusted according to specific files. And each superblock compressed by independent super block as a unit will also record the statistical information of each block. It is that if you want to unzip a block, the first location to the record belongs to super block. It computes the block addresses in the superblock and read all records of sharing information, then can unzip record data in a block. The more Randomness of the DSRC for file to be written to improve the speed of compression. Therefore, to some extent at the expense of the compression rate. In general, the theory of DSRC maximum compression rate in the 5-1, sacrifice the part of the compression efficiency to improve the speed of compression.

LZMA to compress, using different compression coding method to compress with target in order to get the best rate of compression . The final output three streams of data will be archived into a single compressed file.

"0.786811": "enii biiiii3 qo to oiii the (3)

The difference between the DNAZIL and the standard reference genome is 80 bp. The difference between the DNAZIL and the standard reference genome is 80 bp.

The DNA^{1b} compression algorithm uses the compressed file as input to decompress it.

The characteristics of the first filter and then use the EOC to filter itself to get a feedback loop with compensation effect in one-time. The progress of the EOC is EOC. Since it is the first time, the EOC does not aim to the best

The early history of the ROC confessions is divided into

We can use the differentials to approximate derivatives.

translate into an integer data records in a file, so that we can implement that use the previous four characters to compress to an integer data accounts for only a byte, the compression rate of our method is 75%.

Figure 2: The statistical results of the weights

Ниже мы видим, что для векторов, имеющих одинаковые коэффициенты, то есть для векторов, состоящих из единичных элементов, коэффициенты равны единице.

Все остальные коэффициенты равны нулю. Это означает, что для векторов, имеющих одинаковые коэффициенты, то есть для векторов, состоящих из единичных элементов, коэффициенты равны единице.

Но для векторов, имеющих различные коэффициенты, то есть для векторов, состоящих из различных элементов, коэффициенты равны нулю. Это означает, что для векторов, имеющих различные коэффициенты, то есть для векторов, состоящих из различных элементов, коэффициенты равны нулю.

Например, для вектора $\begin{pmatrix} 1 & 2 & 3 & 4 \end{pmatrix}$ коэффициенты равны нулю, кроме первого, который равен единице. Для вектора $\begin{pmatrix} 1 & 2 & 3 & 4 \end{pmatrix}$ коэффициенты равны нулю, кроме первого, который равен единице.

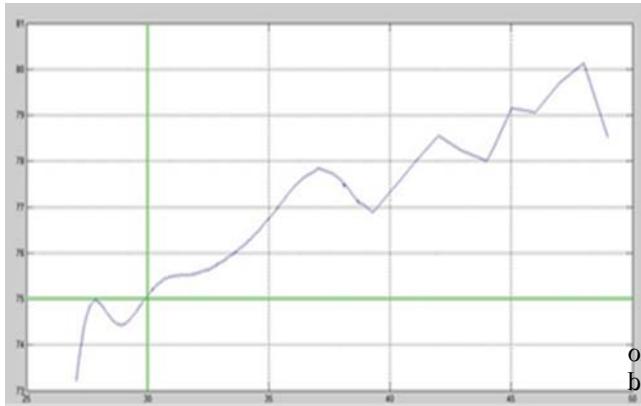
Например, для вектора $\begin{pmatrix} 1 & 2 & 3 & 4 \end{pmatrix}$ коэффициенты равны нулю, кроме первого, который равен единице. Для вектора $\begin{pmatrix} 1 & 2 & 3 & 4 \end{pmatrix}$ коэффициенты равны нулю, кроме первого, который равен единице.

Например, для вектора $\begin{pmatrix} 1 & 2 & 3 & 4 \end{pmatrix}$ коэффициенты равны нулю, кроме первого, который равен единице. Для вектора $\begin{pmatrix} 1 & 2 & 3 & 4 \end{pmatrix}$ коэффициенты равны нулю, кроме первого, который равен единице.

Например, для вектора $\begin{pmatrix} 1 & 2 & 3 & 4 \end{pmatrix}$ коэффициенты равны нулю, кроме первого, который равен единице. Для вектора $\begin{pmatrix} 1 & 2 & 3 & 4 \end{pmatrix}$ коэффициенты равны нулю, кроме первого, который равен единице.

Например, для вектора $\begin{pmatrix} 1 & 2 & 3 & 4 \end{pmatrix}$ коэффициенты равны нулю, кроме первого, который равен единице. Для вектора $\begin{pmatrix} 1 & 2 & 3 & 4 \end{pmatrix}$ коэффициенты равны нулю, кроме первого, который равен единице.

Например, для вектора $\begin{pmatrix} 1 & 2 & 3 & 4 \end{pmatrix}$ коэффициенты равны нулю, кроме первого, который равен единице. Для вектора $\begin{pmatrix} 1 & 2 & 3 & 4 \end{pmatrix}$ коэффициенты равны нулю, кроме первого, который равен единице.



:swallow as if it had nothing

from 10 to 150 [1].

use *Krus* *Gen&th#t* *Sodin&th#t* to find characters which reflect these aspects of words [9]. The effect of the secondarily compression reflects the initial primary compression rule. At the same time, in order to continue to use a rule to find secondary compression to reflect effect of compression. After this *Krus* *Gen&th#t* *Sodin&th#t* can encode as *AZBZ* with *Krus* *Gen&th#t* *Sodin&th#t* to reflect the characters that reflect. For example, *AAAABBBB* can be compressed into characters to replace it with the number of occurrences of characters. The *Krus* *Gen&th#t* *Sodin&th#t* refers to the number of

The Kullback-Leibler Coding refers to the number of

memories of your family coming to compressibly score

This sequence is only # reflect many times, but specieses

We can see a sequence of a difficut score line and it is not

the biography in my book about the construction rate.
Focus on the analysis of the characteristics of the relationship between the construction rate and the construction rate of fixed assets.
So in terms of fixed asset construction, we can choose because for example scores from 1 to 100, it is the maximum value of the construction rate of fixed assets.
General, there is a lot of time to calculate the average construction rate of fixed assets.
[] will be calculated as follows:
construction rate = $\frac{\text{fixed asset construction rate}}{\text{construction rate}}$

Figure 3: Layout of the workflow

outright file is less than 500MB, and as the increase of outright file becomes more in the speed of compression when the size of compression rate. The DDSKC algorithm has almost the same test case, while it becomes not as well as other algorithms in speed, it compresses faster than the other 3 algorithms in all.

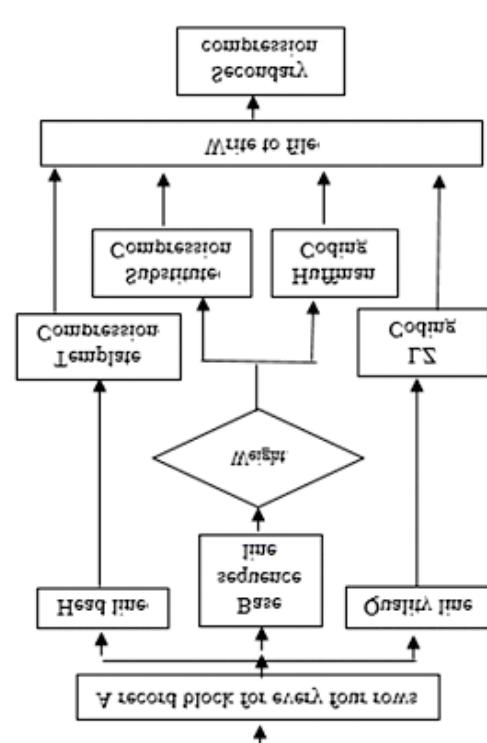
means the time used to compresses the original file size, so lower rate means the better performance. The time note: the rate is equal to compressed file size divide original file

Date		Time		Event Type		Location		Participants		Notes	
2023-09-15	15:00	2023-09-15	15:00	Meeting	Video Call	Headquarters	Online	John Doe	Jane Smith	Discussed project A	✓
2023-09-16	10:00	2023-09-16	10:00	Meeting	In-person	Office A	In-person	John Doe	Jane Smith	Presented findings from research	✓
2023-09-17	14:00	2023-09-17	14:00	Meeting	Video Call	Headquarters	Online	John Doe	Jane Smith	Reviewed financial reports	✓
2023-09-18	09:00	2023-09-18	09:00	Meeting	Video Call	Headquarters	Online	John Doe	Jane Smith	Planned next steps for project A	✓
2023-09-19	13:00	2023-09-19	13:00	Meeting	Video Call	Headquarters	Online	John Doe	Jane Smith	Discussed potential partners	✓
2023-09-20	08:00	2023-09-20	08:00	Meeting	Video Call	Headquarters	Online	John Doe	Jane Smith	Reviewed market trends	✓
2023-09-21	11:00	2023-09-21	11:00	Meeting	Video Call	Headquarters	Online	John Doe	Jane Smith	Planned future meetings	✓
2023-09-22	16:00	2023-09-22	16:00	Meeting	Video Call	Headquarters	Online	John Doe	Jane Smith	Discussed operational efficiency	✓
2023-09-23	09:00	2023-09-23	09:00	Meeting	Video Call	Headquarters	Online	John Doe	Jane Smith	Reviewed performance metrics	✓
2023-09-24	14:00	2023-09-24	14:00	Meeting	Video Call	Headquarters	Online	John Doe	Jane Smith	Planned resource allocation	✓
2023-09-25	08:00	2023-09-25	08:00	Meeting	Video Call	Headquarters	Online	John Doe	Jane Smith	Discussed strategic planning	✓
2023-09-26	12:00	2023-09-26	12:00	Meeting	Video Call	Headquarters	Online	John Doe	Jane Smith	Reviewed competitive analysis	✓
2023-09-27	17:00	2023-09-27	17:00	Meeting	Video Call	Headquarters	Online	John Doe	Jane Smith	Planned stakeholder engagement	✓
2023-09-28	09:00	2023-09-28	09:00	Meeting	Video Call	Headquarters	Online	John Doe	Jane Smith	Discussed regulatory requirements	✓
2023-09-29	13:00	2023-09-29	13:00	Meeting	Video Call	Headquarters	Online	John Doe	Jane Smith	Reviewed legal documents	✓
2023-09-30	17:00	2023-09-30	17:00	Meeting	Video Call	Headquarters	Online	John Doe	Jane Smith	Planned final review	✓

LARGE compression results

Performance Inquiries

—Любимые места для пикника—



file size, it compresses faster than FQZ and FQC while has the similar compression rate. The implement of DDSRC has only single thread when I/O.