1  **TBtools, a <u>T</u>oolkit for <u>B</u>iologists integrating various HTS-data**

2  **handling <u>tools</u> with a user-friendly interface**

3

4  Chengjie Chen[1,2,3*], Rui Xia[1,2,3], Hao Chen[4], Yehua He[2,3]*

5

6  [1]State Key Laboratory for Conservation and Utilization of Subtropical Agro-Bioresources, [2]Key

7  Laboratory of Biology and Germplasm Enhancement of Horticultural Crops in South China,

8  Ministry of Agriculture, [3]College of Horticulture, South China Agricultural University, Guangzhou,

9  510642, China, [4]Oilseed Crops Institute, Hunan Agricultural University, Changsha, 410128

10

11  *The corresponding Author.

12

13  YH: heyehua@hotmail.com

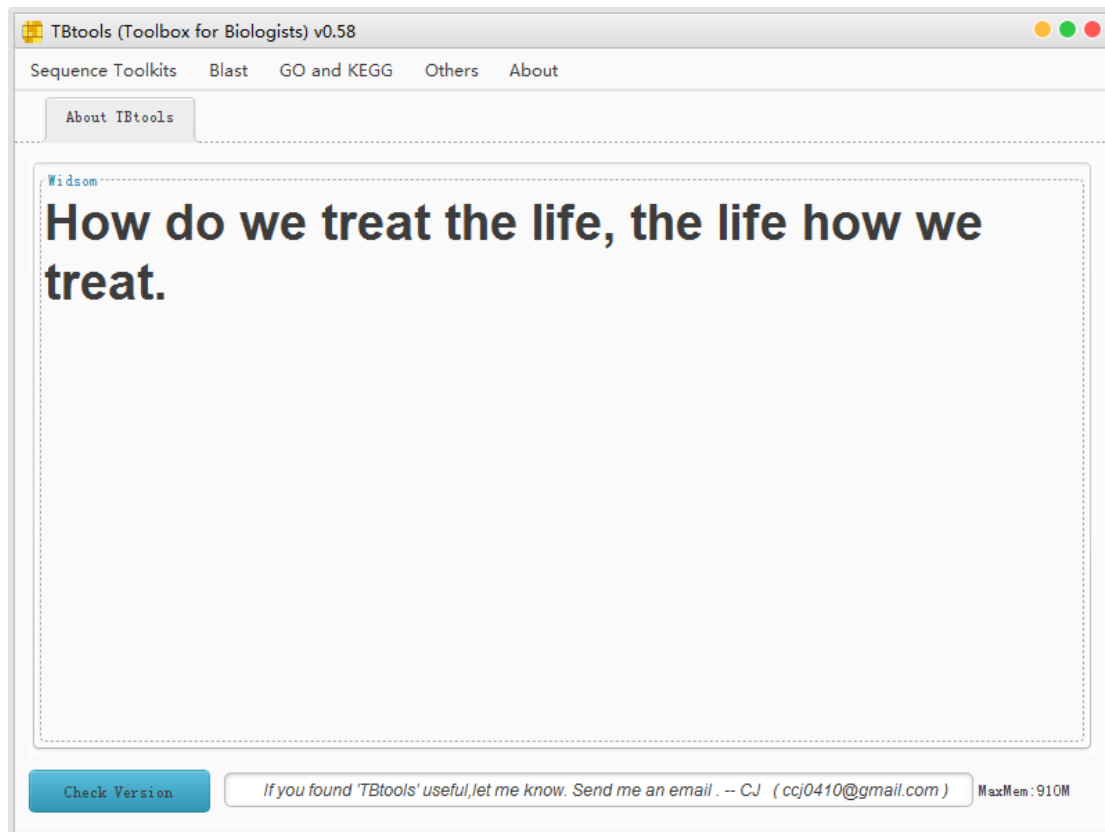14  RX: rxia@scau.edu.cn

15  CC: ccj0410@gmail.com

16

17

## Abstract

Various softwares or pipelines have been developed for biological information mining from high-throughput sequencing (HTS) data, and most of them relies on programming and command-line environment with which most biologists are unfamiliar. Bioinformatic tools with an user-friendly interface are preferred by wet-lab biologists. Here, we describe TBtools, a Toolkit for Biologists integrating various HTS-data handling tools with a user-friendly interface. It includes a large collection of functions, which facilitate many simple, routine but elaborate tasks working on HTS data, such as bulk sequence extraction, gene set functional enrichment, venn diagram and etc. TBtools can run under all operating systems with JRE1.6 and is freely available at github.com/CJ-Chen/TBtools. Since its development, it has been used by many researchers. It will be a useful toolkit for wet-lab biologists to work on all kinds of high-throughput data.

## 1. Introduction

The rapid development of high-throughtput sequencing techniques leads to the boom of sequencing data for biologists, all kinds of softwares, packages, and pipelines have been developed to meet various analysis needs. Most of them relies on programming and command-line environment with which most biologists are unfamiliar. Few user-friendly tools focus on common but elaborate tasks such as bulk sequences extraction, gene set functional enrichment and blast alignments visualization[1]. Many web-based applications, R packages or pipelines require user to either upload big data or work under command-line environment. Running small scripts or a few commands, which might be easy for computation specialists or bioinformaticians to accomplish in seconds, may take non-specialist users much time for the first try and even more time in future. Thus, we here present TBtools, a toolset that sets out to save time for wet-lab biologists from daily sequence analysis work. To date, it includes more than 40 little functions classified into five groups.

44

45  **Figure 1. Overview of the main menu of TBtools.** More than 40 functions are classified into five

46  groups, which are showed on the main menu.

47

## 48  2. Manual of TBtools

49  TBtools is an easy-to-use toolset containing two running mode: GUI mode for out-of-the-box

50  usage and command-line mode for batch analyses. A bundle of functions have been developed in

51  the current version of TBtools (Figure 1, GUI mode) for frequent needs in daily molecular

52  analyses, and more tools can be added with the increase of demands. The "Check Version" could

53  be used to check whether an update of TBtools is available.

54

## 55  2.1 Sequence Toolkits

56  **Fasta Tools are used for the** manipulation of sequence files in fasta format.

57      **Amazing Fasta Extrator** can accomplish almost all fasta sequences extraction tasks with the

58      utilization of pre-built index.

59      **Quick Fasta Extractor or Filter** extracts/filters sequences in a quick mode, without pre-built

60      index.

61      **Fasta Extractor** extract sequences from a fasta file in a slow fasta record parsing mode.
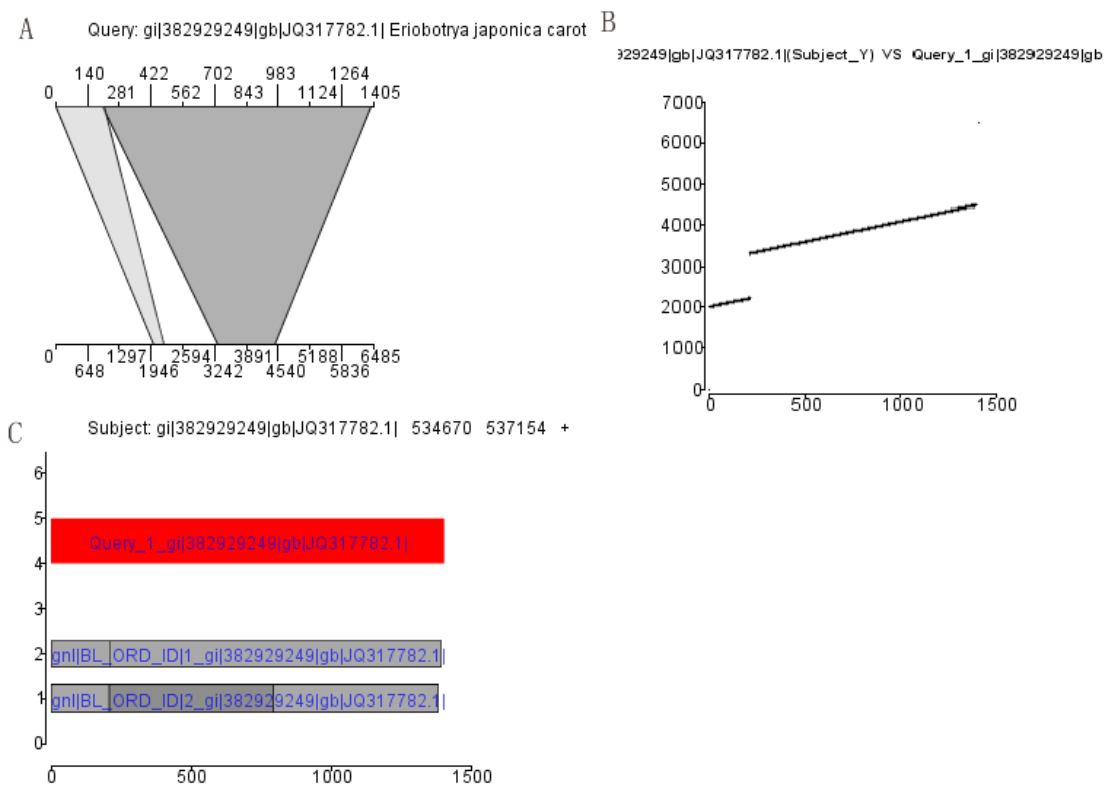
62      (Decrepited)

63      **Fasta Subseq** is similar to Fasta Extrator, but is used to extract subsequences of fasta records.

64      (Decrepited)

65      **Fasta Stat** generates summary information of each fasta record and the whole fasta set, such

66      as number of record, length of sequences, GC content and N50 and so on.

67    **Fasta simplifier** simplifies ID of all fasta records.

68    **Fasta Table Converter** converts sequence file between fasta format and tab-delimited table.

69    **Fasta Merge And Split** merges fasta records from several files or splits fasta records into

70    several files.

71    **Fasta Sequence Manipulator** can do Reverse sequence, Complement sequence, and DNA to

72    RNA conversion.

73    **NCBI Seq Downloader** downloads sequences from NCBI via their API in Fasta/GenBank format.

74    **Get Complete ORF (Open Reading Frame)** predicts complete ORF from input sequence(s).

75    **Check Primers** (Simple e-PCR) can be used to check the specificity of primers via hamming

76    distance.

77    **Gtf/Gff3 Sequences Extractor** extracts sequences from genome sequences according to the

78    gtf/gff3 file, like complete set of cds sequences.

79

## 80    2.2 Blast Wrapper

81    **Auto Blast Tools** contains three auto-blast wrapper, **Blast Several Sequences 2** (to) **A Big**

82    **Database**, **Blast Compare 2 Seq [Sets] and Blast Compare 2 Seq [Sets] <Big File>**, which can

83    automatically determine the molecular type of input sequences.

84    **Reciprocal Blast** is used for conducting reciprocal blast between two fasta files. It can also be

85    used for ID matching between two separate transcriptome assembly results.

86    **Blast Result Visualization** contained three little application to visualize blast results in different

87    format "Alignment Graph", "Dot plot" and "Pileup Graph" (Figure 2).

88    **Blast XML File Converter** converts xml file into tab-delimited table, same as Blast+ "-outfmt 6"

89    and the other collapsed all HSP into a single line.

90

91

92 **Figure 2 Three ways of visualization of Blast XML file.** A) Alignment Graph; B) Dot
93 plot; C) Pileup Graph
94

## 95 2.3 GO and KEGG

96 **GO Annotation** is an application for mapping NCBI gi/accession and Uniprot ID to Gene Ontology
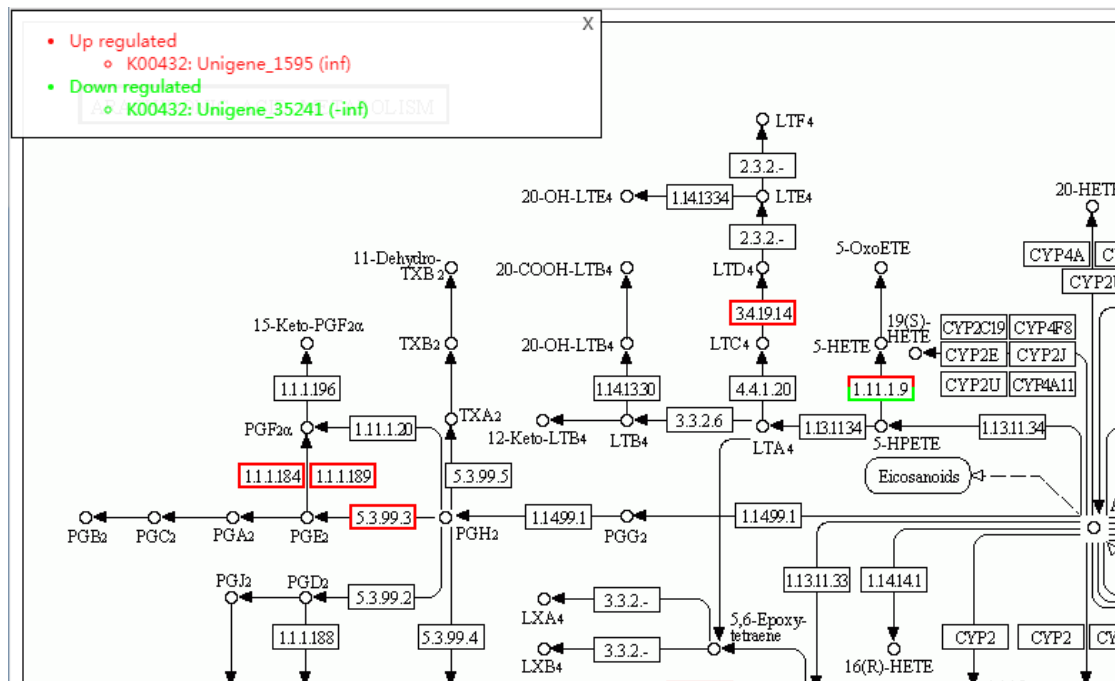97 ID based on ID matchning information from
98 ftp://ftp.pir.georgetown.edu/databases/idmapping/idmapping.tb.gz.
99 **GO Enrichment** conducts GO term enrichment analysis based on hypergemetric distribution.
100 **KEGG Enrichment** conducts KEGG pathway enrichment analysis based on hypergemetric
101 distribution.
102 **KEGG Pathway Map Drawer** colors pathway maps and generates interactive pathway file (.html,
103 Figure 3).
104

105


106 **Figure 3 Interactive Graph generated by TBtools.**
107

## 108 2.4 Others

109 **Table Manipulator** provides two applications (Table ID Extractor and Filter) to manipulate **big**
110 tab-delimited file, like extraction, filtering or sorting of rows and columns.
111 **Wonderful Venn** generates interactive venn plot allowing up to six ID sets (Figure 4A).
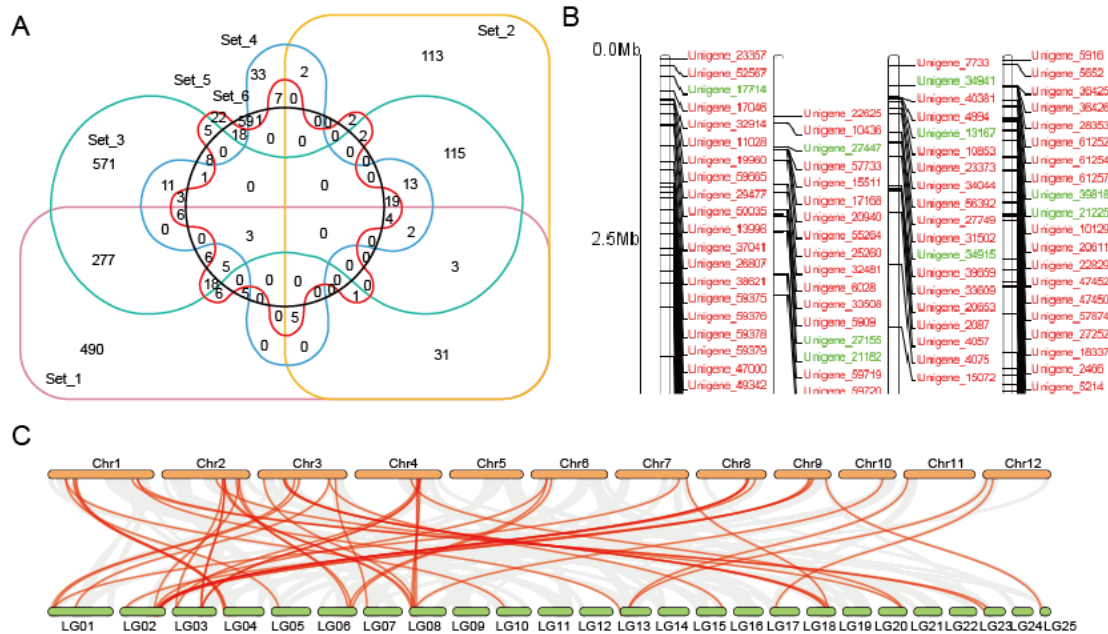112 **Map Gene on Genome** produces a graph showing positions of genes on chromosomes based on
113 input sequences file or position information (Figure 4B).
114 **Dual Synteny Plotter** visualizes results from MCScanX[2] in an interactive mode (Figure 4C).
115 **Domain/Motif Pattern Drawer** visualizes motif/domain information predicted by the MEME suite,
116 Pfam and NCBI CD-search database. Besides, we also integrate a function for gene structure
117 (exon/intron) presentation(Figure 5), which accepts a gtf/gff file and a list of gene IDs. It has no

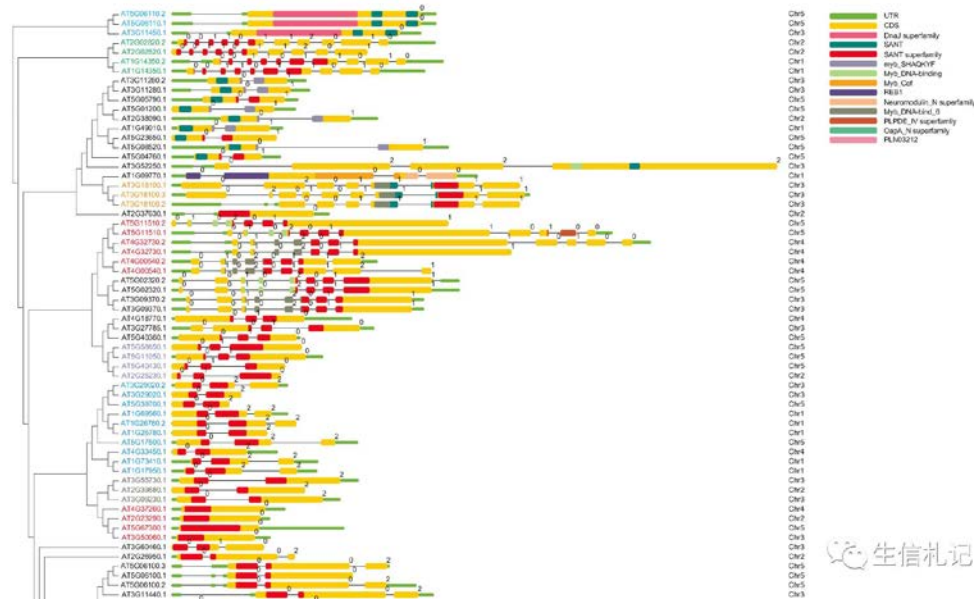118   limit to the size of input file and is easier to use than GSDS[3].

119



120

121   **Figure 4. Example graph of wonderful venn, map gene on genome and dual synteny plotter.**

122

## 3. Discussion and conclusion

124   Recent year, bioinformatics develops fast and has been emerging as a critical discipline for

125   biology study. Major efforts have been put into the development of pipelines, which emphasize

126   on the performance on 'big' bio-data analyses such as genome assembly, annotation, expression

127   analysis, easy-to-use tools with user-friendly interface is much less available. We present TBtools,

128   a toolset to meet the needs of wet-lab biologists for routine data work, who are unfamiliar to

129   programing or command-line environments. TBtools integrates most frequently used functions

130   into one handy executable jar file. Compare to existing softwares or pipelines, such as venny,

131   kalama and comprehesive R packages, TBtools is lighter and easy-to-use, with no requirement of

132   data upload, network connection and special computational skills. On the other hand, TBtools

133   contains a homemade graphic module 'JJplot2' (now JIGplot) which is a java implementation of

134   'Grammer of Graphic', making it very suitable for the production of publishable graphs. We

135   believe that TBtools will be a handy and useful toolset to facilitate downstream utilization of HTS

136   data, especially for web-lab biologists.

137

138

**Figure 5. Examples graph of Motif/Doman Patter Drawer**

140

## 4. Availability and requirements

142 TBtools is available as a jar file on https://github.com/CJ-Chen/TBtools. JRE verison not lower
143 than 1.6 is needed, which is pre-installed in Linux or Mac. Window users may have installed it
144 before using TBtools. Blast+ needs to be added into environment path if user want to use the
145 **Blast wrapper** application.

## Competing interests

147 No competing interest is delared.

148

## Authors' contributions

150 CC and HC designed the software. CC did the coding jobs. YH and RX initiated and supervised the
151 project.

## Acknowledgments

155 1. Neumann, R., Kumar, S., Haverkamp, T. H. & Shalchian-Tabrizi, K. BLASTGrabber: a
156 bioinformatic tool for visualization, analysis and sequence selection of massive BLAST data.
157 *BMC Bioinformatics* **15,** 128 (2014).
158 2. Wang, Y. *et al.* MCScanX: A toolkit for detection and evolutionary analysis of gene synteny and
159 collinearity. *Nucleic Acids Res.* **40,** 1–14 (2012).
160 3. Guo, A. Y., Zhu, Q. H., Chen, X. & Luo, J. C. [GSDS: a gene structure display server]. *Yi Chuan* **29,**
161 1023–1026 (2007).

162