

1 **GFF3sort: a novel tool to sort GFF3 files for tabix**

2 **indexing**

3 Tao Zhu, Chengzhen Liang, Zhigang Meng, Sandui Guo<sup>\*</sup>, and Rui Zhang<sup>\*</sup>

4

5 Biotechnology Research Institute, Chinese Academy of Agricultural Sciences, 100081,

6 Beijing, China

7

8 Correspondence<sup>\*</sup>

9 Rui Zhang

10 [zhangrui@caas.cn](mailto:zhangrui@caas.cn)

11 Sandui Guo

12 [guosandui@caas.cn](mailto:guosandui@caas.cn)

13

## 1 **Abstract**

2 **Background:** The traditional method of visualizing gene annotation data in JBrowse  
3 is converting GFF3 files to JSON format, which is time-consuming. The latest version  
4 of JBrowse supports rendering sorted GFF3 files indexed by tabix, a novel strategy  
5 that is more convenient than the original conversion process. However, current tools  
6 available for GFF3 file sorting have some limitations and their sorting results would  
7 lead to erroneous rendering in JBrowse.

8 **Results:** We developed GFF3sort, a script to sort GFF3 files for tabix indexing.  
9 Specifically designed for JBrowse rendering, GFF3sort can properly deal with the  
10 order of features that have the same chromosome and start position, either by  
11 remembering their original orders or by conducting parent-child topology sorting.  
12 Based on our test datasets from seven species, GFF3sort produced accurate sorting  
13 results with acceptable efficiency compared with currently available tools.

14 **Conclusions:** GFF3sort is a novel tool to sort GFF3 files for tabix indexing. We  
15 anticipate that GFF3sort will be useful to help with genome annotation data  
16 processing and visualization.

17

18 **Keywords:** GFF3, JBrowse, Visualization, Tabix

19

20

## 1 **Background**

2 As a powerful genome browser based on HTML5 and JavaScript, JBrowse has been  
3 widely used since released in 2009[1, 2]. According to its configuration document[3],  
4 it works by first converting genome annotation data in GFF3 file formats to JSON  
5 files by a built-in script “flatfile-to-json.pl”, and then rendering visualized element  
6 models such as genes, transcripts, repeat elements, etc. The main problem, however, is  
7 that this step is extremely time-consuming. The time is proportional to the number of  
8 feature elements in GFF3 files (Additional file 1). Even for small genomes like yeast  
9 (*Saccharomyces cerevisiae*), it takes ~10 seconds to finish the conversion. For large  
10 and deeply annotated genomes such as that of humans, the time increases to more  
11 than 15 minutes. In addition, through the conversion process, a single GFF3 file is  
12 converted to thousands of piecemeal JSON files, thus putting a heavy burden on the  
13 ability to back up and store data.

14 In the recently released JBrowse version (v1.12.3), support for indexed GFF3  
15 files has been added[4]. In this strategy, the GFF3 file is compressed with bgzip and  
16 indexed with tabix[5], which generates only two data files: a compressed file (.gz) and  
17 an index file (.tbi). Compared with the traditional processing protocol, the whole  
18 compression and index process could be finished within a few seconds even for large  
19 datasets such as the human genome annotation data (Additional file 1). The tabix tool  
20 requires GFF3 files to be sorted by chromosomes and positions, which could be  
21 performed in the GNU sort program or the GenomeTools[6] package (see [7]). When  
22 dealing with feature lines in the same chromosome and position, both of the tools  
23 would sort them in an ambiguous way that usually results in parent features being  
24 placed behind their children (Figure 1A). Although this is still valid in tabix indexing,  
25 it would causing erroneous rendering in JBrowse[8] (Figure 1A). Currently there is no  
26 additional options or arguments for current tools to break such tied features by  
27 parent-child relationship. In the absence of a suitable bug fix to JBrowse, an

1 alternative sorting tool is needed to resolve this problem.

2 Here, we present GFF3sort, a novel tool to sort GFF3 files for tabix indexing.  
3 Compared with GNU sort and GenomeTools, GFF3sort produces sorting results that  
4 could be correctly rendered by JBrowse while still keeps enough efficiency. We  
5 anticipate that GFF3sort will be a useful tool to help with processing and visualizing  
6 genome annotation data.

## 7 **Implementation**

8 GFF3sort is a script written in Perl. It uses a hash table to store the input GFF3  
9 annotation data (Figure 1B). For each feature, the chromosome ID and the start  
10 position are stored in the primary and secondary key, respectively. Features with the  
11 same chromosome and start position are grouped in an array in the same order of their  
12 appearance in the original GFF3 data. After sorting the hash table by chromosome IDs  
13 and start positions, GFF3sort implemented two modes to sort features within the array:  
14 the default mode and the precise mode (Figure 1B). In most situations, the original  
15 GFF3 annotations produced by genome annotation projects have already placed  
16 parent features before their children. Therefore, GFF3sort returns the feature lines in  
17 their original order, which is the default behavior. In some situations where orders in  
18 the input file has not obeyed the parent-child relationship, GFF3sort would sort them  
19 according to the parent-child topology using the sorting algorithm of directed acyclic  
20 graph[9], which is the most precise behavior but costs a little more computational  
21 source.

22 In order to test the performance of GFF3sort, the GFF3 annotation files of seven  
23 species, *Saccharomyces cerevisiae* (R64-1-1), *Aspergillus nidulans* (ASM1142v1),  
24 *Chlamydomonas reinhardtii* (INSDC v3.1), *Drosophila melanogaster* (BDGP6),  
25 *Arabidopsis thaliana* (Araport11), *Rattus norvegicus* (Rnor\_6.0), and *Homo sapiens*  
26 (GRCh38), were downloaded from the ENSEMBL database [10]. All the tests were

1 conducted on a SuperMicro® server equipped with 80 Intel® Xeon® CPUs  
2 (2.40GHz), 128 GB RAM, and running the CentOS 6.9 system. By default, CentOS  
3 6.9 carries GNU sort v8.4, a relatively old version released in 2010. Therefore, we  
4 downloaded and installed a new version (v8.28) from the official repository of GNU  
5 Coreutils[11]. Both the old and the new version of GNU sort would be used in  
6 performance test.

## 7 **Results and Discussion**

8 GFF3sort takes a GFF3 file as its input data and returns a sorted GFF3 file as output.  
9 Several optional parameters are provided such as turning on the precise mode, sorting  
10 chromosomes in different ways and properly dealing with inline FASTA sequences.  
11 Element models sorted by GFF3sort can be correctly rendered by JBrowse (Figure  
12 1C).

13 Besides the fixation of JBrowse rendering, GFF3sort has also other advantages over  
14 traditional tools. Compared with the GNU sort program, GFF3sort can properly deal  
15 with GFF3-specific lines or directives that are preceded by the '##' symbol, such as  
16 the topmost GFF version line and the heading sequence-region line. Compared with  
17 the GenomeTools, GFF3sort runs significantly faster (Additional file 1). In the default  
18 mode, GFF3sort saves ~70% running time in our seven test datasets. The precise  
19 mode takes longer time but still runs faster than GenomeTools, especially for large  
20 annotation data such as human. While keeping a high running speed, the memory  
21 consumption is still acceptable (Additional file 1). For the largest annotation dataset  
22 (the GRCh38 annotation version of human) with a ~400MB GFF3 file, the memory  
23 usage of GFF3sort is ~758MB, ~40% less than GenomeTools.

## 1 **Conclusions**

2 In conclusion, GFF3sort is a novel tool to sort GFF3 files for tabix indexing and  
3 therefore can be used to visualize annotation data in JBrowse appropriately. It has a  
4 fast running speed compared with similar, existing tools. We anticipate that GFF3sort  
5 will be a useful tool to simplify data processing and visualization.

## 6 **Figure Legends**

7 **Figure 1. The motivation for, outlines of, and action effects of GFF3sort.** A) An  
8 example of incorrectly sorted GFF3 data and its snapshots in JBrowse. Blocks with  
9 the same start position are marked in blue-yellow stripes. The two lines (mRNA)  
10 marked in red were placed after their sub-features (exon or UTR). Such incorrect  
11 placement leads to losing the first exon in JBrowse rendering results. See Additional  
12 file 2 for the full annotation lines. B) Outlines of GFF3sort. C) An example of  
13 correctly sorted data by GFF3sort and its snapshots in JBrowse. In this example, the  
14 two lines (mRNA) marked in red were correctly placed before their sub-features,  
15 allowing JBrowse to render them properly.

## 16 **Additional files**

17 Additional file 1: Benchmark data. This file displays: 1) the detailed running time of  
18 GFF3-to-JSON conversion and the bgzip-tabix process on our test datasets; 2) the  
19 detailed running time and 3) memory usage of GFF3sort, GNU sort (v8.4 and v8.28),  
20 and GenomeTools on our test datasets. (PDF)

21 Additional file 2: The full GFF3 annotation lines used in Figure 1A and C. It is the  
22 gene AT1G01110 extracted from the *Arabidopsis thaliana* (Araport11) annotation  
23 files. It includes three plain-text files: raw.gff3, GNUsort.gff3 (Figure 1A),

1 and GFF3sort.gff3 (Figure 1C). (ZIP)

## 2 **List of abbreviations**

3 JBrowse: JavaScript-based genome browser

4 GFF3: General Feature Format, version 3

5 JSON: JavaScript Object Notation

6 HTML5: HyperText Markup Language, version 5

## 7 **Declarations**

### 8 **Ethics approval and consent to participate**

9 Not applicable.

### 10 **Consent for publication**

11 Not applicable.

### 12 **Availability of data and material**

13 Project name: GFF3sort

14 Project home page: <https://github.com/billzt/gff3sort>

15 Operating system(s): Linux

16 Programming language: Perl

17 Other requirements: No

18 License: No restrictions for academic users.

19 Any restrictions to use by non-academics: license needed

### 20 **Competing interests**

21 The authors declare that they have no competing interests.

## 1 **Funding**

2 This work is supported by grants from the National Science and Foundation of China  
3 (Grant No. 31771850) and the Ministry of Agriculture of China (Grant No.  
4 2016ZX08005004).

## 5 **Authors' contributions**

6 SG, RZ, and TZ initiated the idea of the tool and conceived the project. TZ designed  
7 the tool and analyzed the data. CL and ZM helped to test the tool. TZ wrote the paper.  
8 All authors read and approved the final manuscript.

## 9 **Acknowledgements**

10 We thank Dr. Miklos Csuros and other anonymous reviewers for their helpful  
11 comments.

## 12 **References**

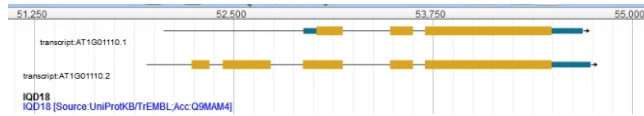
- 13 1. Skinner ME, Uzilov AV, Stein LD, Mungall CJ, Holmes IH: **JBrowse: A**  
14 **next-generation genome browser**. *Genome Res* 2009, **19**(9):1630-1638.
- 15 2. Buels R, Yao E, Diesh CM, Hayes RD, Munoz-Torres M, Helt G, Goodstein  
16 DM, Elisk CG, Lewis SE, Stein L *et al*: **JBrowse: a dynamic web platform**  
17 **for genome visualization and analysis**. *Genome Biol* 2016, **17**(1):66.
- 18 3. **JBrowse Configuration Guide**  
19 [[http://gmod.org/wiki/JBrowse\\_Configuration\\_Guide](http://gmod.org/wiki/JBrowse_Configuration_Guide)] Accessed 26 May 2017
- 20 4. **JBrowse-1.12.3: Maintenance Release** [<http://jbrowse.org/jbrowse-1-12-3/>]  
21 Accessed 26 May 2017
- 22 5. Li H: **Tabix: fast retrieval of sequence features from generic**  
23 **TAB-delimited files**. *Bioinformatics* 2011, **27**(5):718-719.
- 24 6. Gremme G, Steinbiss S, Kurtz S: **GenomeTools: a comprehensive software**  
25 **library for efficient processing of structured genome annotations**.  
26 *IEEE/ACM Trans Comput Biol Bioinformatics* 2013, **10**(3):645-656.
- 27 7. **JBrowse FAQ** [[http://gmod.org/wiki/JBrowse\\_FAQ](http://gmod.org/wiki/JBrowse_FAQ)] Accessed 26 May 2017
- 28 8. **Potential GFF3Tabix issues** [<https://github.com/GMOD/jbrowse/issues/780>]  
29 Accessed 26 May 2017
- 30 9. **Sort::Topological - Topological Sort - metacpan.org**



- 1            [\[https://metacpan.org/pod/Sort::Topological\]](https://metacpan.org/pod/Sort::Topological) Accessed 4 June 2017
- 2    10.    Aken BL, Achuthan P, Akanni W, Amode MR, Bernsdorff F, Bhai J, Billis K,
- 3            Carvalho-Silva D, Cummins C, Clapham P *et al*: **Ensembl 2017**. *Nucleic*
- 4            *Acids Res* 2017, **45**(D1):D635-D642.
- 5    11.    **Coreutils - GNU core utilities**
- 6            [\[https://www.gnu.org/software/coreutils/coreutils.html\]](https://www.gnu.org/software/coreutils/coreutils.html) Accessed 15 Sept 2017
- 7

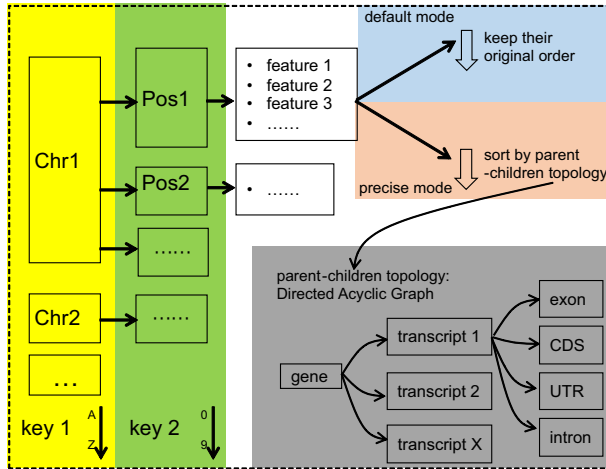
A

araport11	exon	51953	52346	+	Parent=AT1G01110. 2
araport11	five_prime_UTR	51953	52238	+	Parent=AT1G01110. 2
araport11	gene	51953	54737	+	ID=AT1G01110
araport11	mRNA	51953	54737	+	ID=AT1G01110. 2;Parent=AT1G01110
araport11	exon	52061	52730	+	Parent=AT1G01110. 1
araport11	five_prime_UTR	52061	52730	+	Parent=AT1G01110. 1
araport11	mRNA	52061	54689	+	ID=AT1G01110. 1;Parent=AT1G01110
araport11	CDS	52239	52346	+ 0	ID=CDS: AT1G01110. 2;Parent=AT1G01110. 2
araport11	CDS	52434	52730	+ 0	ID=CDS: AT1G01110. 2;Parent=AT1G01110. 2
araport11	exon	52434	52730	+	Parent=AT1G01110. 2;
araport11	CDS	52938	53183	+ 0	ID=CDS: AT1G01110. 2;Parent=AT1G01110. 2
araport11	exon	52938	53183	+	Parent=AT1G01110. 1;
araport11	exon	52938	53183	+	Parent=AT1G01110. 2;
araport11	five_prime_UTR	52938	53021	+	Parent=AT1G01110. 1
araport11	CDS	53022	53183	+ 0	ID=CDS: AT1G01110. 1;Parent=AT1G01110. 1
araport11	CDS	53484	53624	+ 0	ID=CDS: AT1G01110. 1;Parent=AT1G01110. 1
araport11	CDS	53484	53624	+ 0	ID=CDS: AT1G01110. 2;Parent=AT1G01110. 2
araport11	exon	53484	53624	+	Parent=AT1G01110. 1
araport11	exon	53484	53624	+	Parent=AT1G01110. 2



B

### Hash Table



C

araport11	gene	51953	54737	+	ID=AT1G01110
araport11	mRNA	51953	54737	+	ID=AT1G01110. 2;Parent=AT1G01110
araport11	five_prime_UTR	51953	52238	+	Parent=AT1G01110. 2
araport11	exon	51953	52346	+	Parent=AT1G01110. 2
araport11	mRNA	52061	54689	+	ID=AT1G01110. 1;Parent=AT1G01110
araport11	exon	52061	52730	+	Parent=AT1G01110. 1
araport11	five_prime_UTR	52061	52730	+	Parent=AT1G01110. 1
araport11	CDS	52239	52346	+ 0	ID=CDS: AT1G01110. 2;Parent=AT1G01110. 2
araport11	exon	52434	52730	+	Parent=AT1G01110. 2;
araport11	CDS	52434	52730	+ 0	ID=CDS: AT1G01110. 2;Parent=AT1G01110. 2
araport11	exon	52938	53183	+	Parent=AT1G01110. 2;
araport11	CDS	52938	53183	+ 0	ID=CDS: AT1G01110. 2;Parent=AT1G01110. 2
araport11	five_prime_UTR	52938	53021	+	Parent=AT1G01110. 1
araport11	exon	52938	53183	+	Parent=AT1G01110. 1;
araport11	CDS	53022	53183	+ 0	ID=CDS: AT1G01110. 1;Parent=AT1G01110. 1
araport11	exon	53484	53624	+	Parent=AT1G01110. 1
araport11	CDS	53484	53624	+ 0	ID=CDS: AT1G01110. 1;Parent=AT1G01110. 1
araport11	exon	53484	53624	+	Parent=AT1G01110. 2
araport11	CDS	53484	53624	+ 0	ID=CDS: AT1G01110. 2;Parent=AT1G01110. 2

