

A New File Format and Tools for the Large-Scale Data Submission to DNA Data Bank of Japan (DDBJ)

Satoru Miyazaki¹ Hiroyuki Hashimoto² Akemi Shimada¹
smiyazak@genes.nig.ac.jp hirohash@genes.nig.ac.jp akshimad@genes.nig.ac.jp
Yoshio Tateno¹ Hideaki Sugawara¹
ytateno@genes.nig.ac.jp hsugawar@genes.nig.ac.jp

¹ Center for Information Biology, National Institute of Genetics,
1111 Yata, Mishima, Shizuoka 411-8540, Japan

² Hitachi Software Engineering Co., Ltd. Yokohama 231-0015, Japan

Keywords: DNA sequence, genome, large-scale submission, JAVA, object-oriented

1 Introduction

Thanks to the genome projects and the rapid development of high-throughput technique to determine DNA sequences, two types of very large-scale data have been submitted to DDBJ/EMBL/Genbank International Nucleotide Sequence Databases (INSD). One is a whole genome sequence of an organism. Another is a set of sequences for the phylogenetic study based on a ubiquitous gene among hundreds of organisms. It is very important for researchers to get accession numbers of their sequences in a timely fashion. INSD has to meet their needs even if the data is massive.

To make submission of data easier, INSD has developed and provided such tools for stand-alone usage and Web submission as Authorin, Sequin, SAKURA, Webin, BankIt. These submission tools, however, target relatively small scale data and are not suitable for massive data.

Therefore we at DDBJ developed a new data file format and off-line tools to support the submission of the large-scale sequence data.

2 Method and Results

2.1 A new file format

Authorin and Sequin use the transaction file format and ASN.1 format respectively. EMBL flat file format is required to send your data to EMBL database. These formats, however, have some points we have to consider. For instance, files described by ASN.1 format are difficult to edit directly, but easy to find some syntax error by computers. EMBL format is easy to edit by yourself, but difficult to describe many kinds of combination of feature keys and qualifiers, that are two level data descriptors, without syntax errors. To realize readability and easiness to parse, we designed a new file format for the description of biological annotation to sequences. It is a tab-delimited text format with five columns (Entry, Feature, Location, Qualifier and Value) (See Fig. 1).

2.2 Mass Submission Tool (MST)

MST is a stand-alone application program and is able to handle the annotation data consisting of over 100 features with various qualifiers for more than 10,000 sequences. This application is consist of three modules: annotation editor, sequence editor and parser tool. MST is developed based on the object-oriented concept by use of a library implemented by JAVA language. The features from the viewpoints of programming are:

1. executable on all-most all operating system
2. easy to modify user interface for the addition and the deletion of feature keys and qualifiers
3. each module is also executable as an application by minor modification of source code

Fig. 2 is sample windows of annotation editor on MST.

| Line | Feature Key | Qualifier | Value |
|------|-------------|-----------|---------------|
| 1 | ORIGIN | | EXCHANGIA (1) |
| 2 | ORIGIN | | EXCHANGIA (1) |
| 3 | ORIGIN | | EXCHANGIA (1) |
| 4 | ORIGIN | | EXCHANGIA (1) |
| 5 | ORIGIN | | EXCHANGIA (1) |
| 6 | ORIGIN | | EXCHANGIA (1) |
| 7 | ORIGIN | | EXCHANGIA (1) |
| 8 | ORIGIN | | EXCHANGIA (1) |
| 9 | ORIGIN | | EXCHANGIA (1) |
| 10 | ORIGIN | | EXCHANGIA (1) |
| 11 | ORIGIN | | EXCHANGIA (1) |
| 12 | ORIGIN | | EXCHANGIA (1) |
| 13 | ORIGIN | | EXCHANGIA (1) |
| 14 | ORIGIN | | EXCHANGIA (1) |
| 15 | ORIGIN | | EXCHANGIA (1) |
| 16 | ORIGIN | | EXCHANGIA (1) |
| 17 | ORIGIN | | EXCHANGIA (1) |
| 18 | ORIGIN | | EXCHANGIA (1) |
| 19 | ORIGIN | | EXCHANGIA (1) |
| 20 | ORIGIN | | EXCHANGIA (1) |
| 21 | ORIGIN | | EXCHANGIA (1) |
| 22 | ORIGIN | | EXCHANGIA (1) |
| 23 | ORIGIN | | EXCHANGIA (1) |
| 24 | ORIGIN | | EXCHANGIA (1) |
| 25 | ORIGIN | | EXCHANGIA (1) |
| 26 | ORIGIN | | EXCHANGIA (1) |
| 27 | ORIGIN | | EXCHANGIA (1) |
| 28 | ORIGIN | | EXCHANGIA (1) |
| 29 | ORIGIN | | EXCHANGIA (1) |
| 30 | ORIGIN | | EXCHANGIA (1) |
| 31 | ORIGIN | | EXCHANGIA (1) |
| 32 | ORIGIN | | EXCHANGIA (1) |
| 33 | ORIGIN | | EXCHANGIA (1) |
| 34 | ORIGIN | | EXCHANGIA (1) |

Figure 1: The format of Annotation file

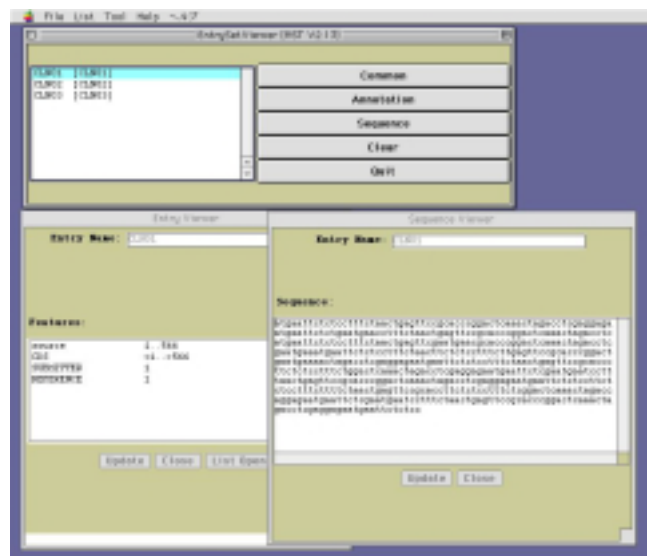


Figure 2: Sample windows on MST

References

- [1] Sugawara, H., Miyazaki, S., Gojobori, T. and Tateno, Y., DNA Data Bank of Japan dealing with large-scale data submission, *Nucleic Acids Res.*, 27(1):25–28, 1999.
- [2] Tateno, Y., Fukami-Kobayashi, K., Miyazaki, S., Sugawara, H. and Gojobori, T., DNA Data Bank of Japan at work on genome sequence data, *Nucleic Acids Res.*, 26(1):16–20, 1998.
- [3] Tateno, Y., Miyazaki, S., Ota, M., Sugawara, H. and Gojobori, T., DNA Data Bank of Japan (DDBJ) in collaboration with mass sequencing teams [In Process Citation], *Nucleic Acids Res.*, 28(1):24–26, 2000.