Genome Visualization and functional genomics Analysis

**STEP INTO GENOVA**

# Part 1. Views Introduction

GENOVA provides three views for genome visualization and modification. They are the feature editor (WorkBench), the marker editor (Marker) and the visualization canvas (Draw). The following figure is a typical view of the feature editor, which is initiated automatically once the genome file (GenBank format) has been read. In this example, the *Staphylococcus aureus* RN1 genome sequence is studied in the table. Users could insert or delete nucleotide bases, genomic elements (also in batch) and feature entries of interest into the genome. Note there are two more features than in the common GenBank content, i.e., '***category***' and '***term***', where users can specify the desired tag and favorite color. Marker view is a secondary table which is designed to note the genomic rearrangement both in table and in the canvas. The last view is the canvas, users can invoke it using the "canvas" button, there GENOVA provides abundant options for obtaining a optimal visualization, zooming, range, changing terms, etc.

# 1. Feature Editor



**Fig. T1 Feature Editor – WorkBench**. Features can be managed in the table, sorting and changing the values. All the cells except the first column '#' can be modified. The leading novel genome is generated and synchronized in the canvas. Operations can be carried using the buttons in the toolbar, the function are listed as below:

- Major buttons of the horizontal toolbar

▱ Load a sequence file of GenBank format

▱ Save as a sequence file of GenBank

▱ Save the genomic map with notes and markers into a picture file

▱ Plot a genomic map

▱ Content Analysis

ⓘ License information

▱ Closing the workspace


- Editor buttons in the vertical bar


▱ Column configuration of the feature table

▱ Categorization palette

▱ Search the keywords in the specified column

▱ Set the values in a certain column of selected rows

✚ Add a feature entry into the table

▬ Remove selected feature entries from the table, the corresponding nucleotide sequence can be also removed out of the genome after confirmation

▱ Insert nucleotide bases into the genome

✪ Delete nucleotide bases of desired range from the genome


Moreover, there is a "Filter" row below the editor, which enables the user to pick and show only feature entries according to keywords, e.g. "CDS" as filter will hide all entries which contain no "CDS". Users can easily type the category name in the "category" column. The system will give a random color, which always can be changed in the palette. All entries assigned in a category will take the color specified previously. Genes involved in the same island or in the same operon are given with preference (default setting) the same or similar color in order to distinguish these genes from others, e.g., in the table, light purple denotes all genes belong to the isoleucine-valine operon.
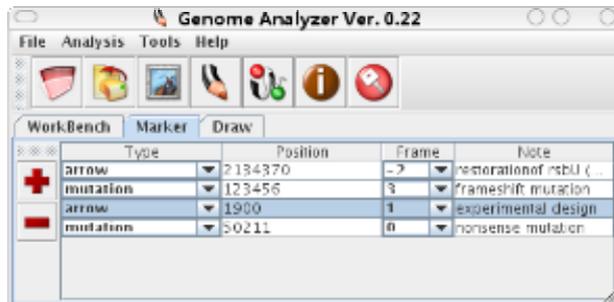
## 2. Marker Editor (Marker)



**Fig. T2 Marker Editor.** Markers can be inserted, edited and removed in this view. In the "Type" column, users can specify the favorite symbol (mutation is the "X" symbol), the corresponding note information will be painted as well in the canvas. Note that the reading frame (+3,+2,+1,0,-1,-2,-3) should also be assigned according to the place users would like to insert or plot additional markers. The left control panel enables to insert the marker or remove the selected one.
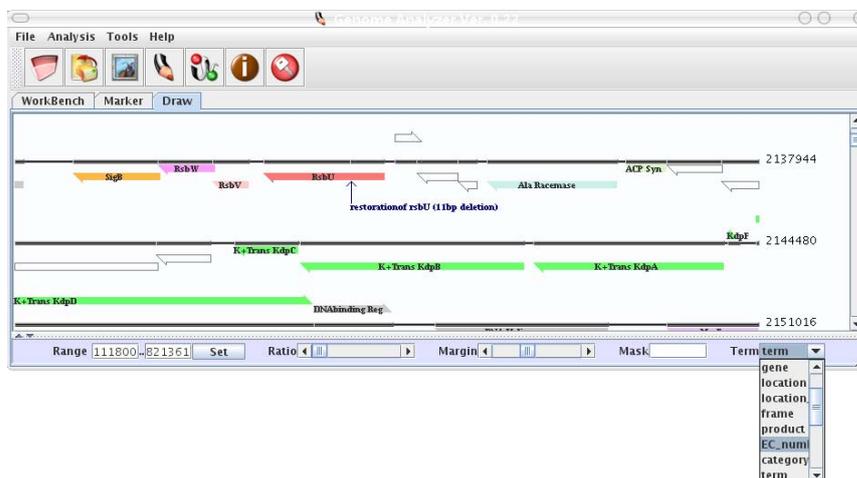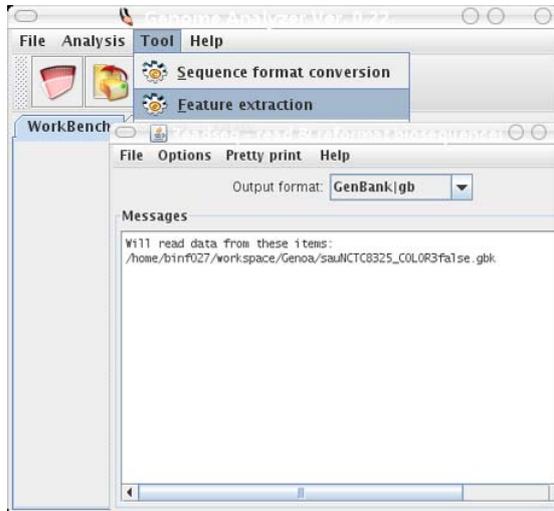
## 3. Visualization Canvas (Draw)



**Fig. T3 Visualization canvas.** The control panel below the canvas provides the possibility to set the visualization range ("Set" button), modify the zoom factor ("Ratio" slider), increase or decrease the line margin ("Margin" slider) , specify the term visualized and used for gene labeling ("Term"), e.g., *locus_tag, gene, protein_id, EC_number*, etc. There is another function named "Mask", which plays a crucial function to hide redundant symbols in the terms, in particular when the "locus_tag" is chosen as term: Similar leading words, e.g., the prefix of "SAOUHSC_" can be hidden to prevent cluttering of the figure using the mask.
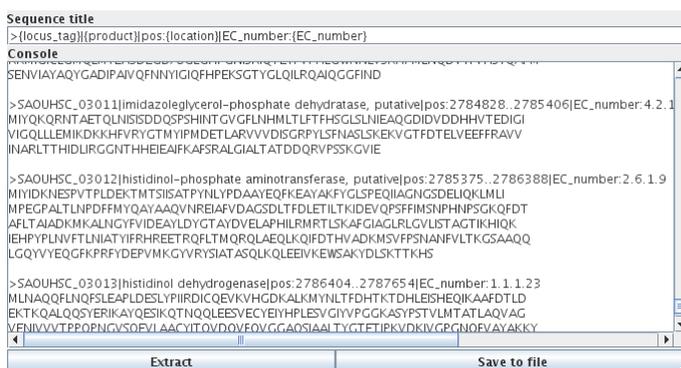
# 4. Tools

1)



2)



3)



**Fig. T4 Tools incorporated into the GENOVA.** A tool-kit allows multiple genome sequence format conversion (1; readseq [suppl. 8]) and comprehensive feature extraction (2), as well as construction of a proteome file from GenBank format genome data (3; as long as the input file contains translated protein entries). All this can be rapidly managed in seconds.

## Part 2. Typical steps (suggested tutorial tour)

1. Import the sequence file with the annotation into the GENOVA software using the "***open***"  button. (sample sequence in the sample sub-directory)



**Fig. T5 Open a GenBank file located in the sample directory.**

2. Edit the features and sequences in the "feature editor" (***workbench*** tab). Firstly please type CDS in the bottom filter text-field, click the "***Exec***" button or simply press return key, so that only CDS entries are listed in the editor. Scroll them to the desired region, e.g., SAOUHSC_02300 in the following figure.



**Fig. T6 Overview of operations in the genome feature editor.**

3. Categorize them by writing directly into cells of the "***category***" column, e.g., *rsb_sigB*, or any other desired name (ilV, RsbV, RsbU in the sample file), system will automatically detect the repeats and assign a random color. Of course you are welcome to specify a favourite color using the "***palette***"  button directly. This table is actually serving as a collector of function groups.

**Fig. T7 Categorization palette.**

4. Click the "***Plot***" button  to visualize the ORFs, the ***Term*** combo-box in the bottom control panel allows users to switch different tags appearing in the ORFs, e.g., *locus_tag, gene name, position, function, EC number* as well as the ***term***, which is in addition in the feature editor (column position see figure of step.2) allows any desired notes from users without changing the other original annotation. The following figure is an example when the user specify the *locus_tag* as the ORF names. Here users are suggested to remove the redundant prefix using the "***Mask***", please input SAOUHSC_0 in the text-field. Any regular expression can be used here as the mask filter. In the following steps, we suggest to choose the term in the term combo-box for the illustration. The other slider "***ratio***" is for the zoom factor, and "***margin***" is for increasing or decreasing the line margin.
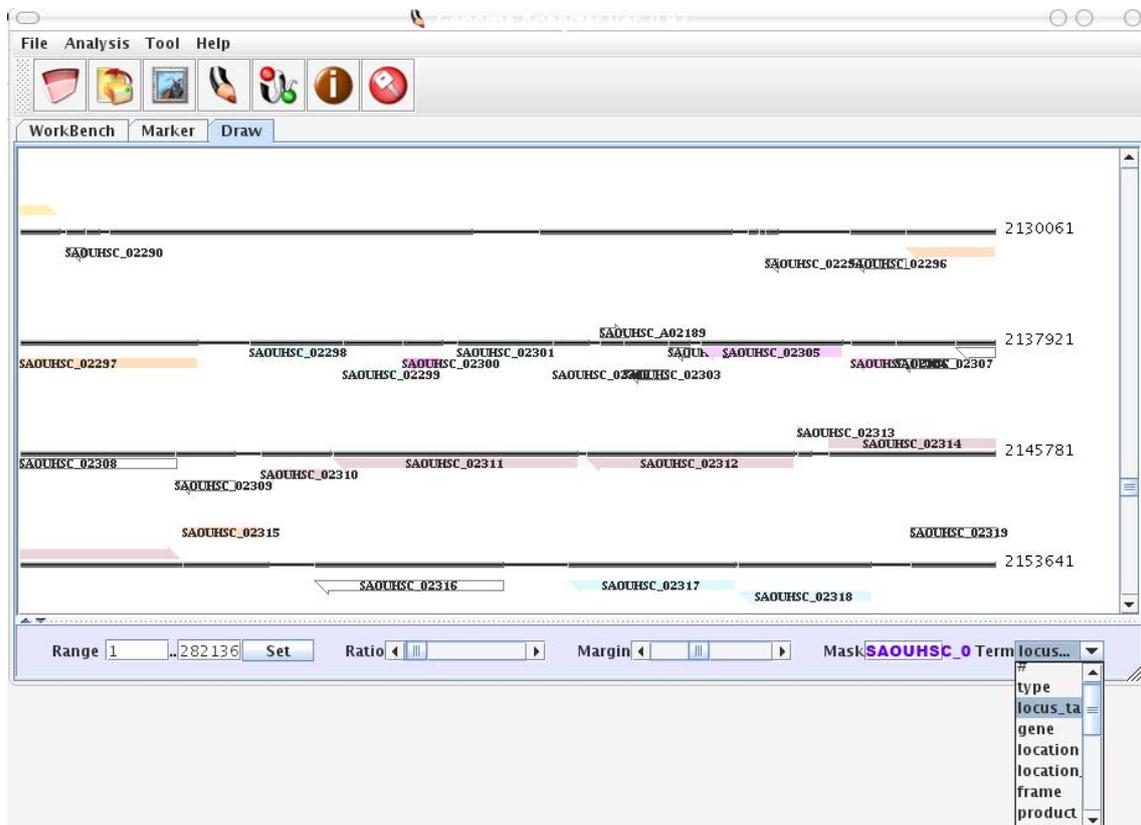


**Fig. T8 Canvas operation.**

5. Add some markers in the "***marker***" view (Fig. T2) to note down the genome region as well as the interesting rearrangement events. Frame is critical for acquiring an optimal figure, since the ORFs are drawn according to their reading frame numbers (3,2,1,0,-1,-2,-3). Here 0 suggests the arrow should point to the genome scale line directly.

6. Observe them again by selecting the "***draw***" tab in the index (see Fig.2 in the GENOVA manuscript)

7. Click the "***Statistics***" button to inspect the content analysis.
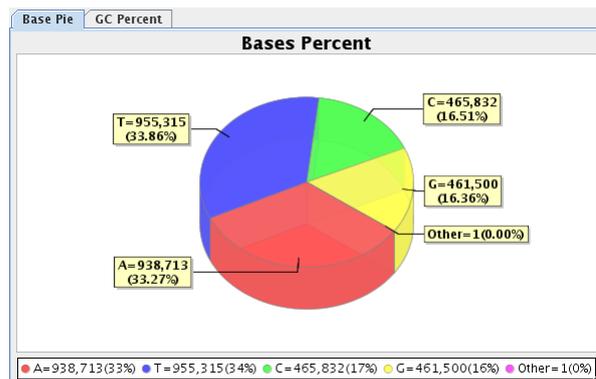


**Fig. T8 Content analysis.**

8. Optional: Edit the features, inserting and removing entries, adding and deleting nucleotide sequences (Fig. T1). This operation will change the genome length and all the ORFs, the software will return the result list of all the deleted entries, certain entries partially changed which required for a manual confirmation or modification.

9. Save the novel genome as a GenBank file using the "***save***" button. The categorization information will also be stored.

10. The resulting genome figure can be exported into a picture file using the "***save genome map***" button.

11. Optional: Toolkit to extract feature entries or convert the file into different formats or a proteome sequence-file (Fig. T4). The feature extraction tool allows users to select the feature type firstly in the upper list, specify the desired title-format and generate a sequence collection in fasta format. The "***save to file***" button enables to save as a permanent file in the disk.

12. Latest update, other tutorials or more sample files please visit our website located at http://genova.bioapps.biozentrum.uni-wuerzburg.de, thank you very much and we are looking forward to hearing from you.