# An Automated System for EST Functional Analysis and Its Application in Rice Genes

**Hsiu-Ying Lu (呂秀英)[1][3], Cheng-Tao Chen (陳政道)[1], Yi-Chia Chiu (邱怡嘉)[2], Chun-Tang Lu (呂椿棠)[1], Meng-Li Wei (魏夢麗)[1]**

[1]Agricultural Research Institute, Taiwan, ROC.
[2]Graduate Institute of Biotechnology, Chaoyang University of Technology, Taiwan, ROC.
[3]Correspondence author (iying@wufeng.tari.gov.tw)

## Motivation

Expressed sequence tag (EST), the partial fragment of cDNA sequence, represents the most extensive available survey of the transcribed portion of the genome. Rice (*Oryza sativa* L.) is the world's most important food crop; it is also a good model for studies of monocot plants. A vast number of rice EST sequences in the public databases provide an important resource for functional annotation of rice genome. However, ESTs are generated and deposited in the public domain, as redundant, un-annotated, single-pass reactions, with virtually no biological content.  While efforts to assemble, organize and annotate raw EST sequence data have been developed employing diverse strategies, they typically read single input files, produce single output files and require extensive manual intervention by the user. Thus, there is a need for developing an integrated and fully automated program to perform the assembly, annotation, functional classification and statistical analysis of large quantities of EST sequences.

## Material and Methods

- The complete procedure for EST functional analysis was established as (Figure 1): generation of a non-redundant data set (unique sequences) after clustering and assembling the raw ESTs (by CAP3 software), functional annotation of unique sequences (by BLAST similarity search to public annotated TC database in TIGR), functional classification (by linking to GO and MIPS catalog systems from tentative annotation of the selected TC sequences), and statistical analysis of gene functional representation.

- We firstly downloaded all the sequences of TC database, CAP3 and BLAST softwares, as well as GO and MIPS catalogs. Then, TC-sequence to function-catalog relations were established and a MySQL database was created. All programs of this system were written in Perl script and developed on the Linux platform.  To store and manage sequence data and analytical results, PhpMyAdmin package was used to connect the user's web browser and MySQL database.

- Brown planthopper (*Nilaparvata lugens* (Stal.)) is the major pest in rice so is the rice blast disease (infected by *Magnaporthe grisea*) the most serious problem. The automated system was implemented in the rice ESTs induced by *N. lugens* (with the total amount as 188) and *M. grisea* (with the total amount as 84,705) from NCBI-dbEST database. We spent 3 months finishing the functional analysis of 188 ESTs induced by *N. lugens* using manual intervention. However, it's impossible to accomplish the entire process of the huge amount of ESTs induced by *M. grisea* within one year without automated system.
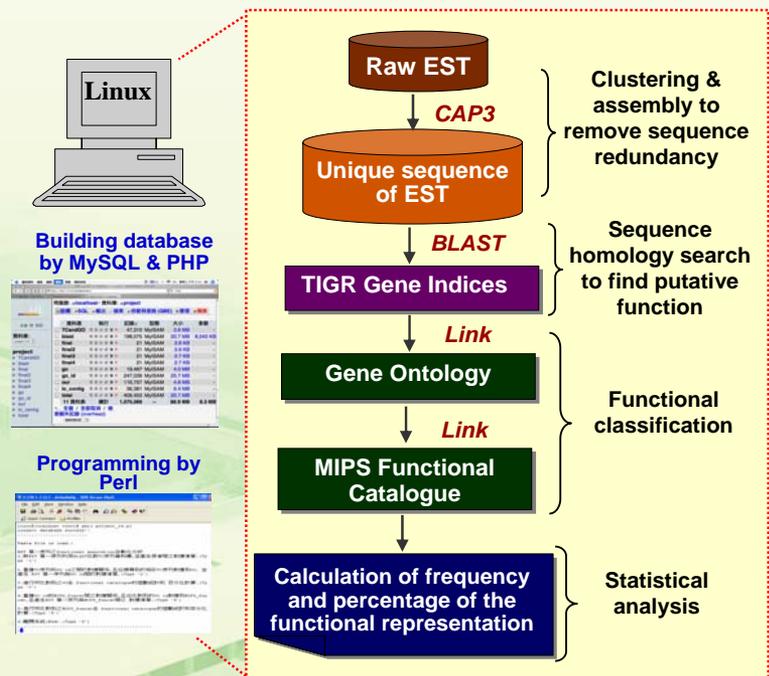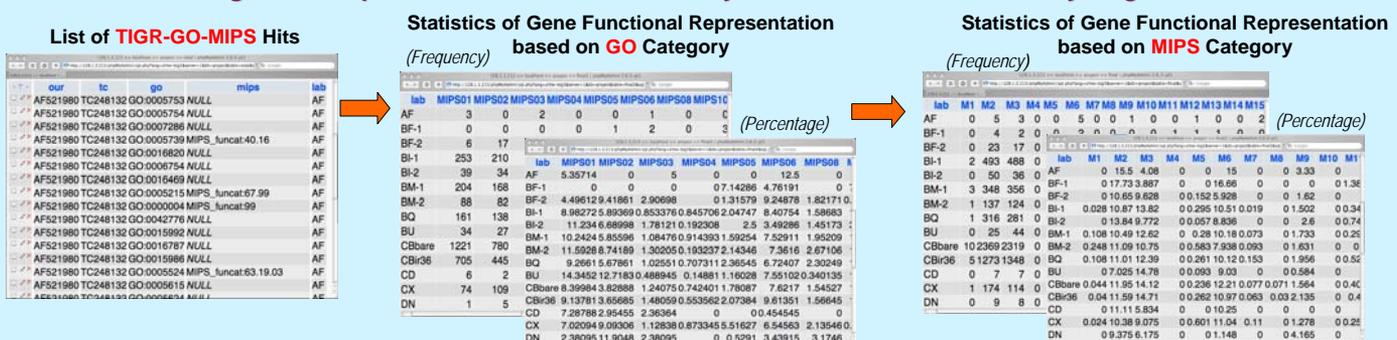


Figure 1. System Structure and Pipeline Strategy of automated system

## Results and Discussion

- This automated system considerably reduced analysis time from 3 months to 2 hours for the 188 rice ESTs induced by *N. lugens*. It also allows us to finish the fully functional analysis of 84,705 rice ESTs induced by *M. grisea* within 7 days.
- The output results contain 5 kinds of records and can be saved to an EXCEL file (an example shown in Figure 2).
- The findings of rice EST annotation help clarify the functional mechanisms of rice genes resistant to pest and disease.

### Figure 2. Output Results of Functional Analysis for Rice ESTs induced by *M. grisea*

行政院農業委員會農業試驗所
**Agricultural Research Institute, COA**