# D4

# DIRECT SUBMISSION SYSTEM AND LITERATURE ANNOTATION OF RICE GENES IN ORYZABASE

Yukiko Yamazaki[1], Rie Tsuchiya[1], Takao Asanuma[2], Yo Shidahara[2], Shingo Sakaniwa[1]

(1. National Institute of Genetics, Mishima, Japan, 2. NalaPro Technologies, Tokyo, Japan)

## Abstract

Oryzabase (http://www.shigen.nig.ac.jp/rice/oryzabase/) is a comprehensive rice science database [1]. It houses a variety of genetic resources, relevant literatures, gene dictionary, DNA sequences, and basic information such as developmental biology and anatomy. In order to keep the gene dictionary up-to-date, literature annotation has been conducted manually since 1995. However as the publication of journal articles increases year by year after genomic sequences were released, it became more difficult to update the dictionary timely and in high quality without sufficient annotators. To overcome this difficulty, we applied machine learning and text-mining to extract known and unknown genes from journals. The machine extraction followed by manual annotation achieved promising results and increased efficiency in manual annotation.

Furthermore a direct submission system where rice researchers can deposit new genes according to the standardized nomenclature [2] became operational in 2008. Recent advances will be introduced.

1. Kurata, N. and Y. Yamazaki., Oryzabase, An Integrated Biological and Genome Information Database for Rice. *Plant Physiology* (2006) 140, 12-17
2. Susan R. McCouch, Gene Nomenclature System for Rice, *Rice* (2008) 1:72-84

## Future plan

Make the gene dictionary more up-to-date, accurate, and comprehensive.

- improve precision of machine extraction
- collaboration with RAP group on gene annotation
- encourage researchers to submit new genes before publication
- encourage researchers to give feedback on Oryzabase genes

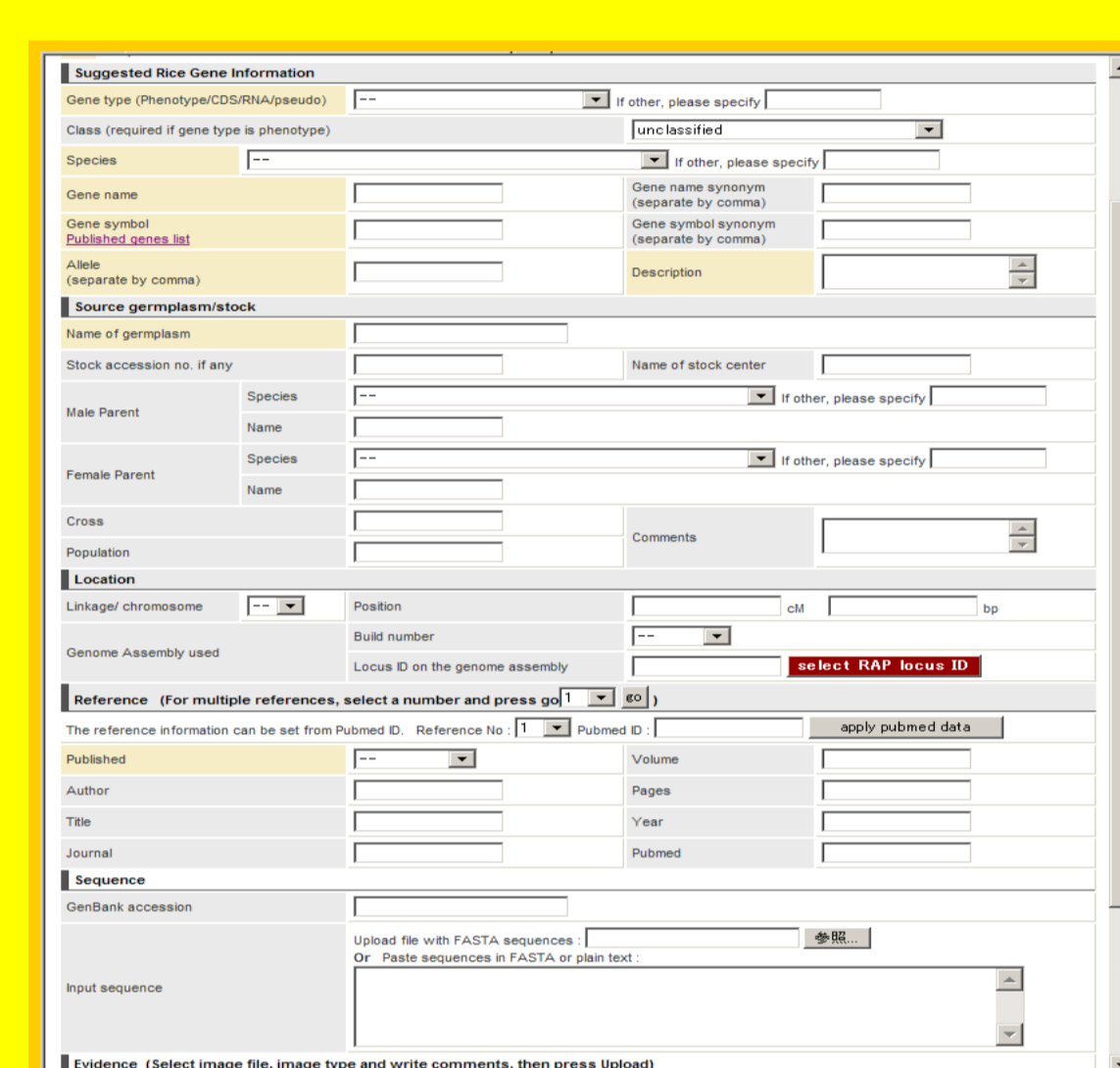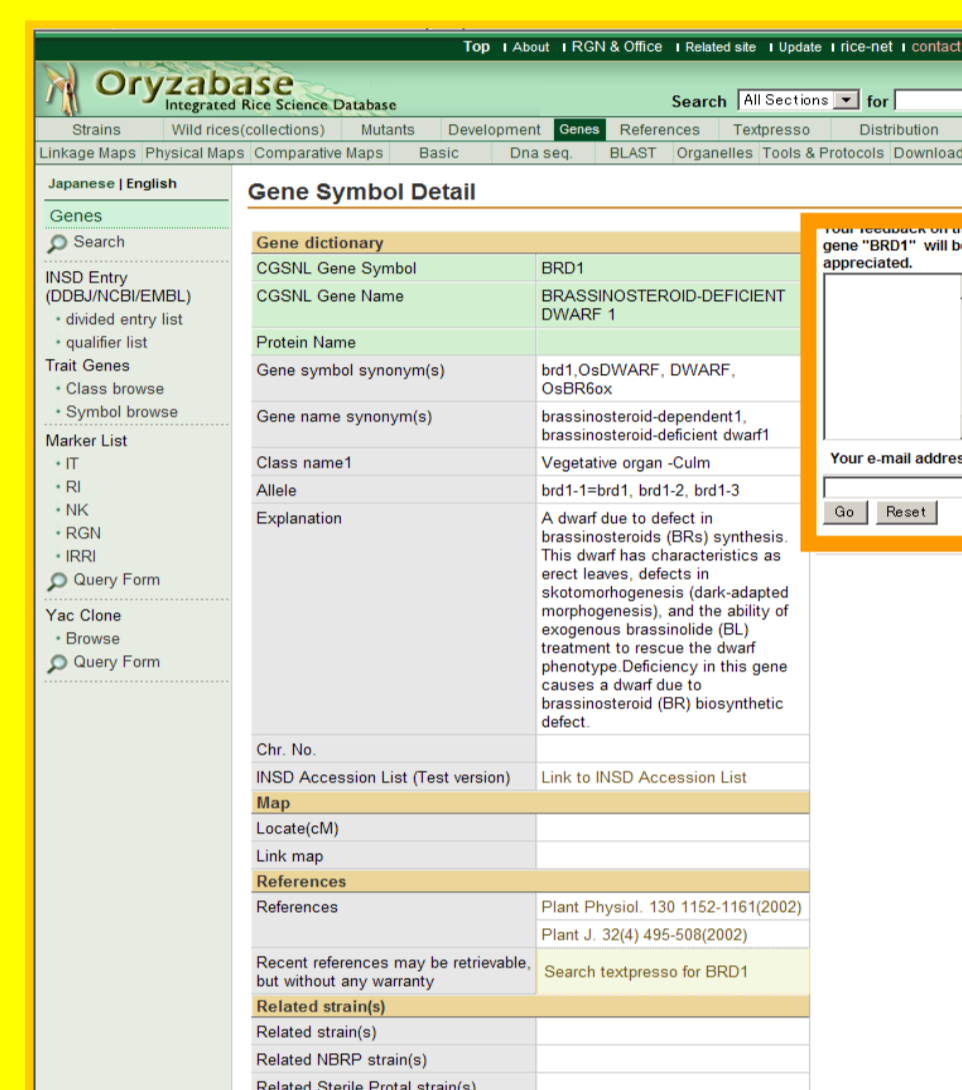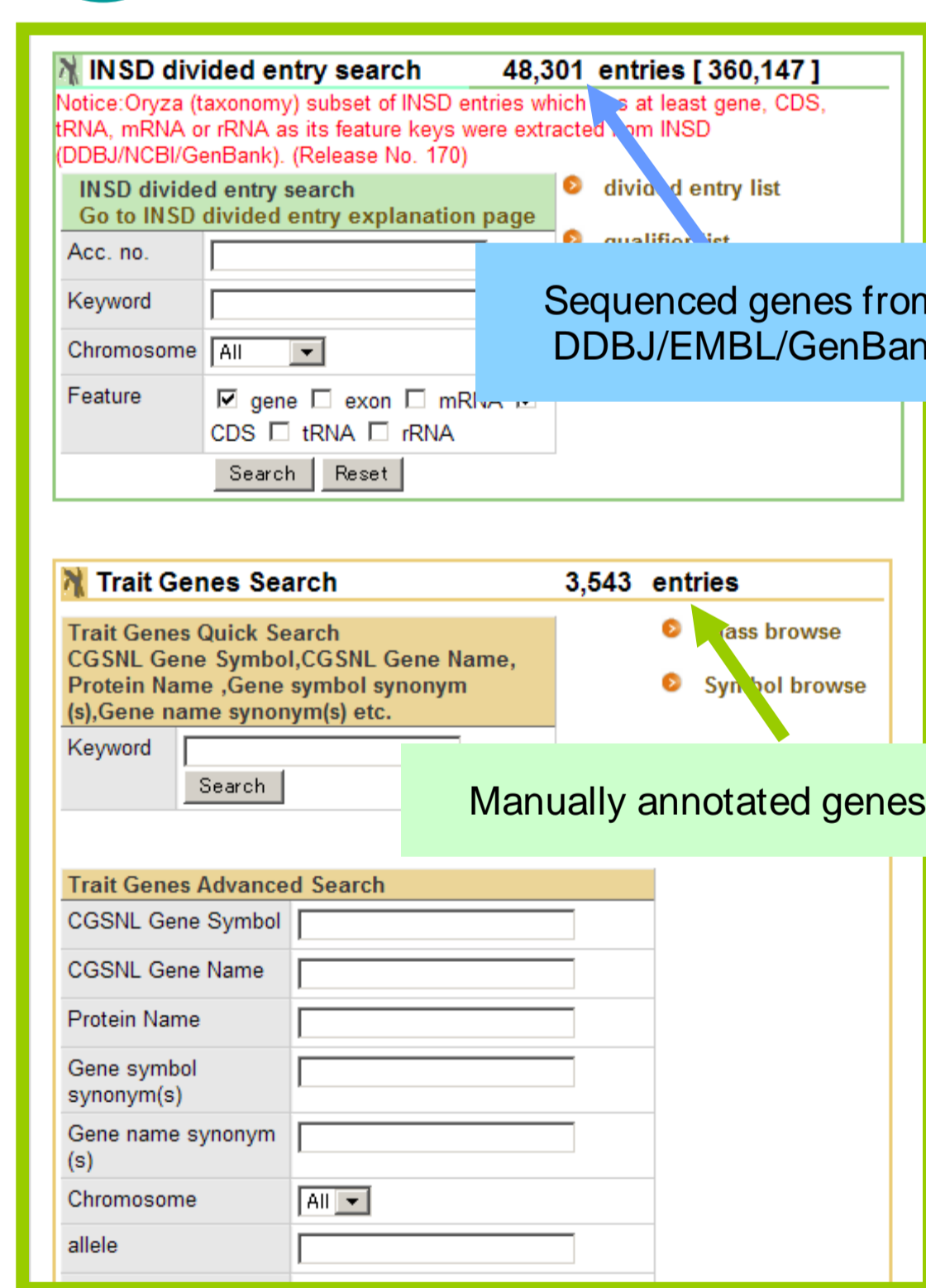## *Oryzabase* HP

# www.shigen.nig.ac.jp/rice/oryzabase/

### 3 Rice Textpresso
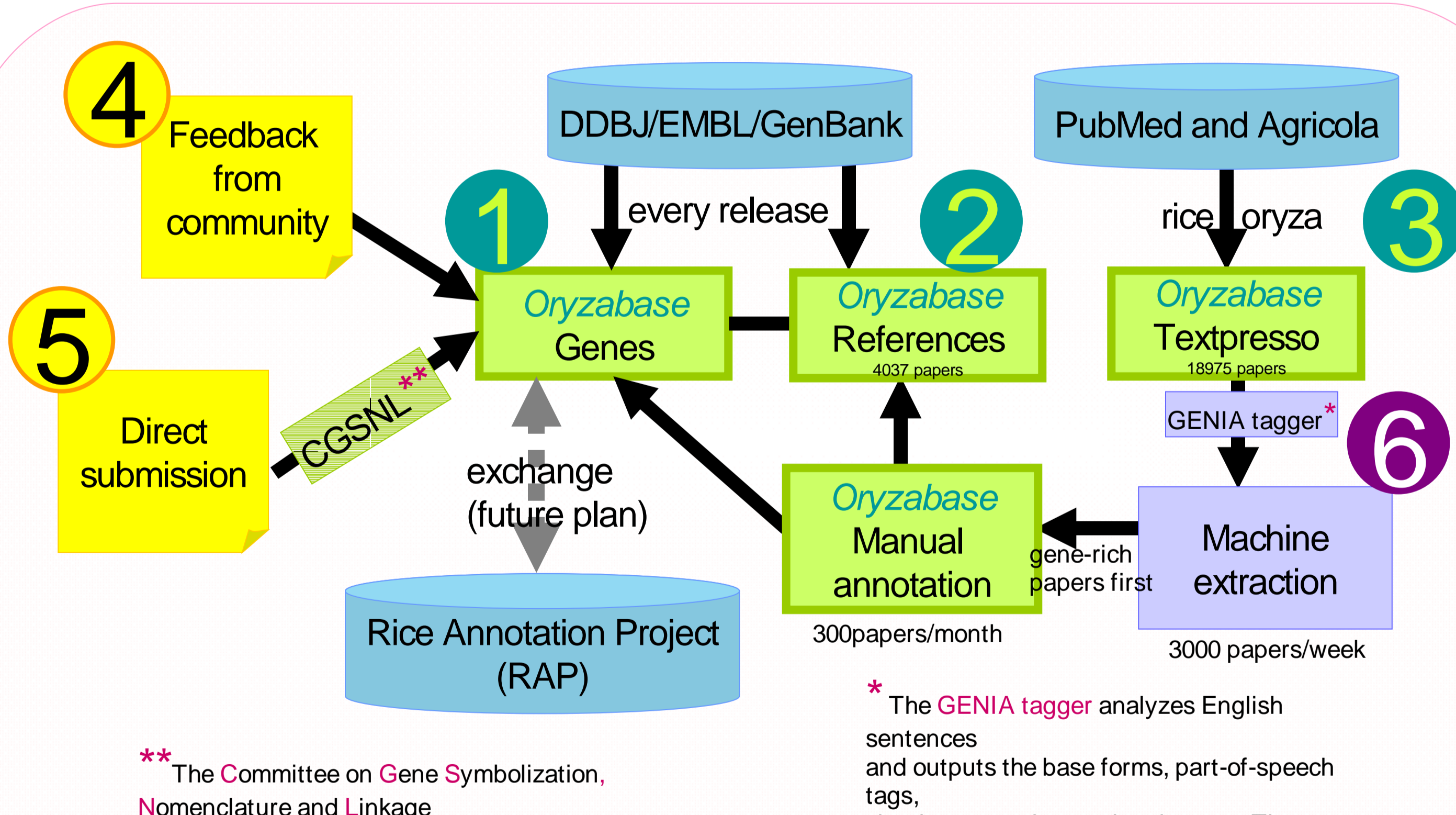
abstract (18889), title (18975)

### 4 Feedback to a gene

### 5 on-line gene submission system

## 1 *Oryzabase* genes

INSD divided entry search 48,301 entries [360,147]

Sequenced genes from DDBJ/EMBL/GenBank

Trait Genes Search 3,543 entries

Manually annotated genes

## *Oryzabase* gene annotation flow

4 Feedback from community

5 Direct submission

1 *Oryzabase* Genes

2 *Oryzabase* References 4037 papers

3 *Oryzabase* Textpresso 18975 papers

DDBJ/EMBL/GenBank — every release

PubMed and Agricola — rice, oryza

CGSNL **

exchange (future plan)

Rice Annotation Project (RAP)

*Oryzabase* Manual annotation
300papers/month

gene-rich papers first

Machine extraction
3000 papers/week

GENIA tagger *

6

** The Committee on Gene Symbolization, Nomenclature and Linkage

* The GENIA tagger analyzes English sentences and outputs the base forms, part-of-speech tags, chunk tags, and named entity tags. The tagger http://www-tsujii.is.s.u-tokyo.ac.jp/GENIA/tagger/ specifically tuned for biomedical text such as MEDLINE abstracts.
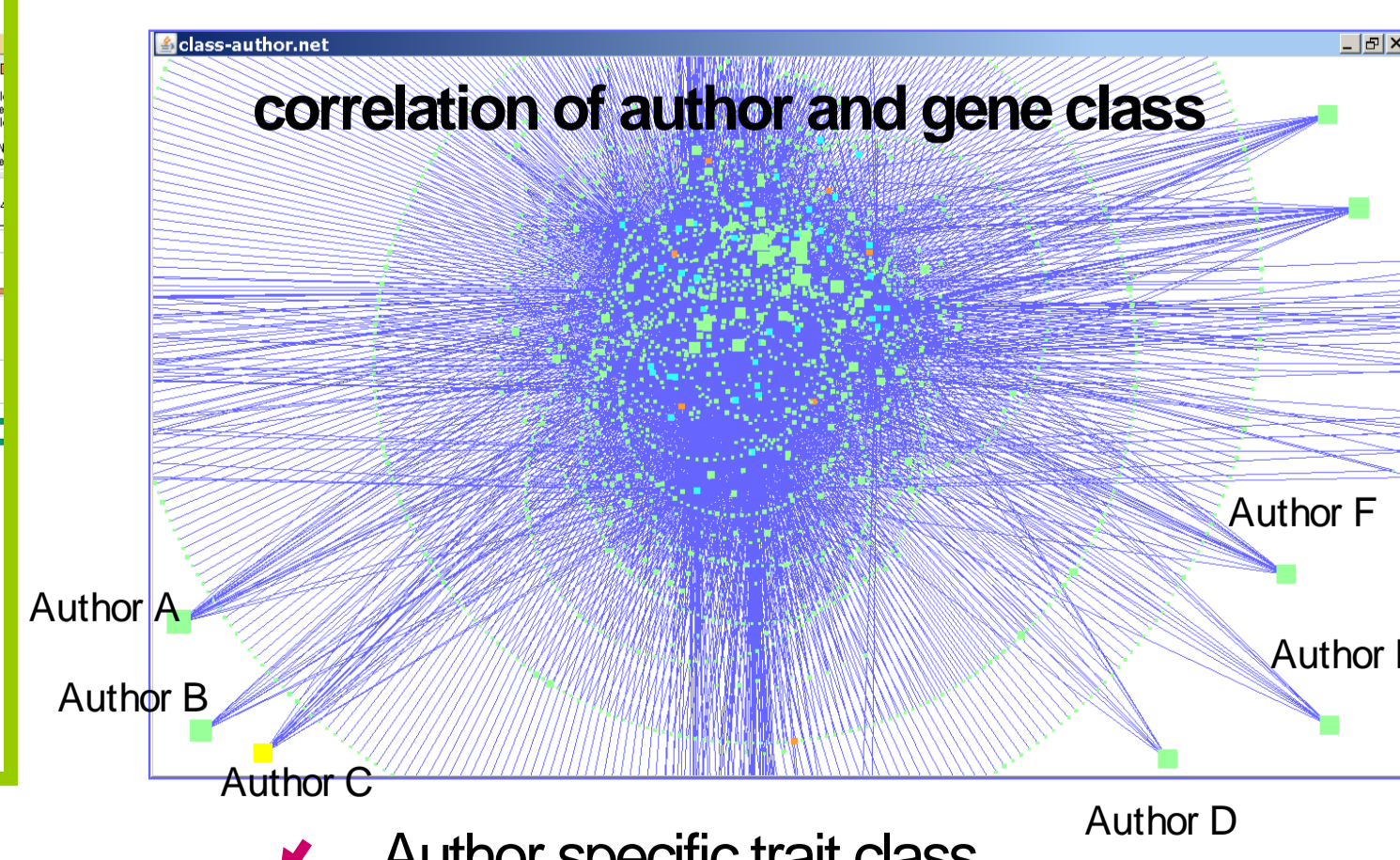
### 6-1 Materials used for machine extraction

Oryzabase Gene Dictionary (as referencedata)
Rice DNA marker name (from RGP)
WordNet Dictionary for general words

GENIA Tagger ( Corpus) *
 Protein
 DNA
 RNA
 Cell Line
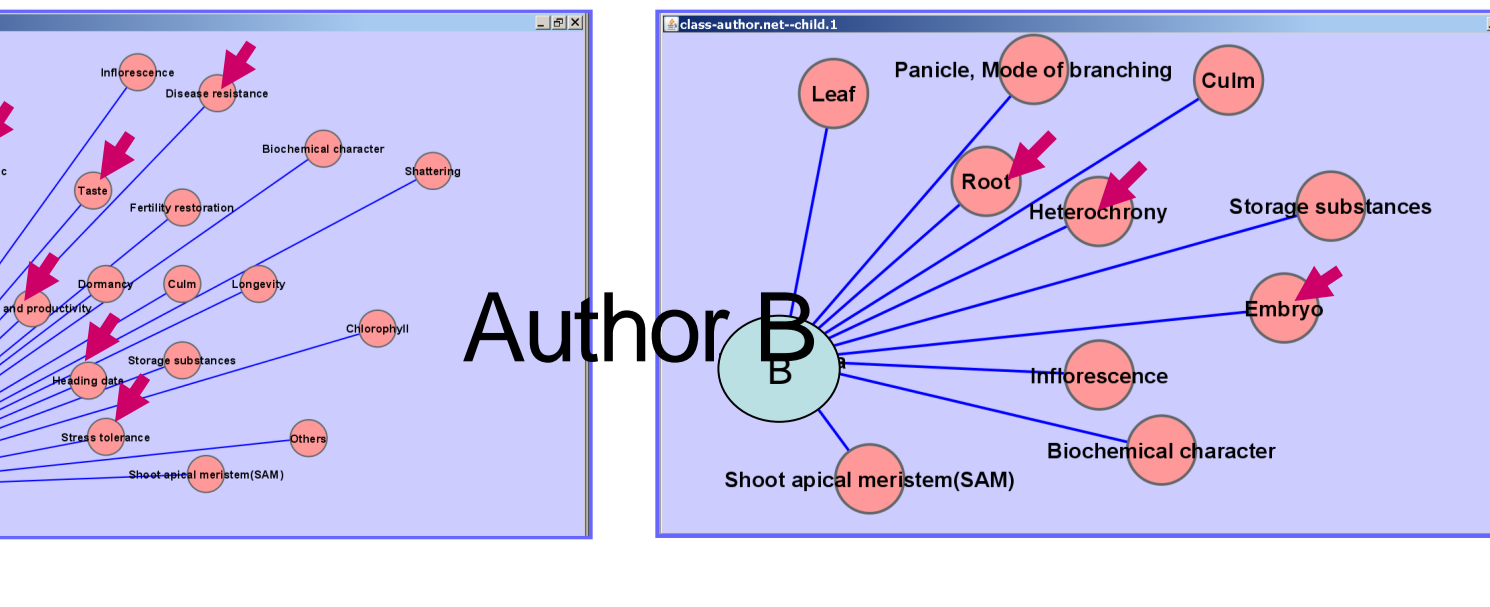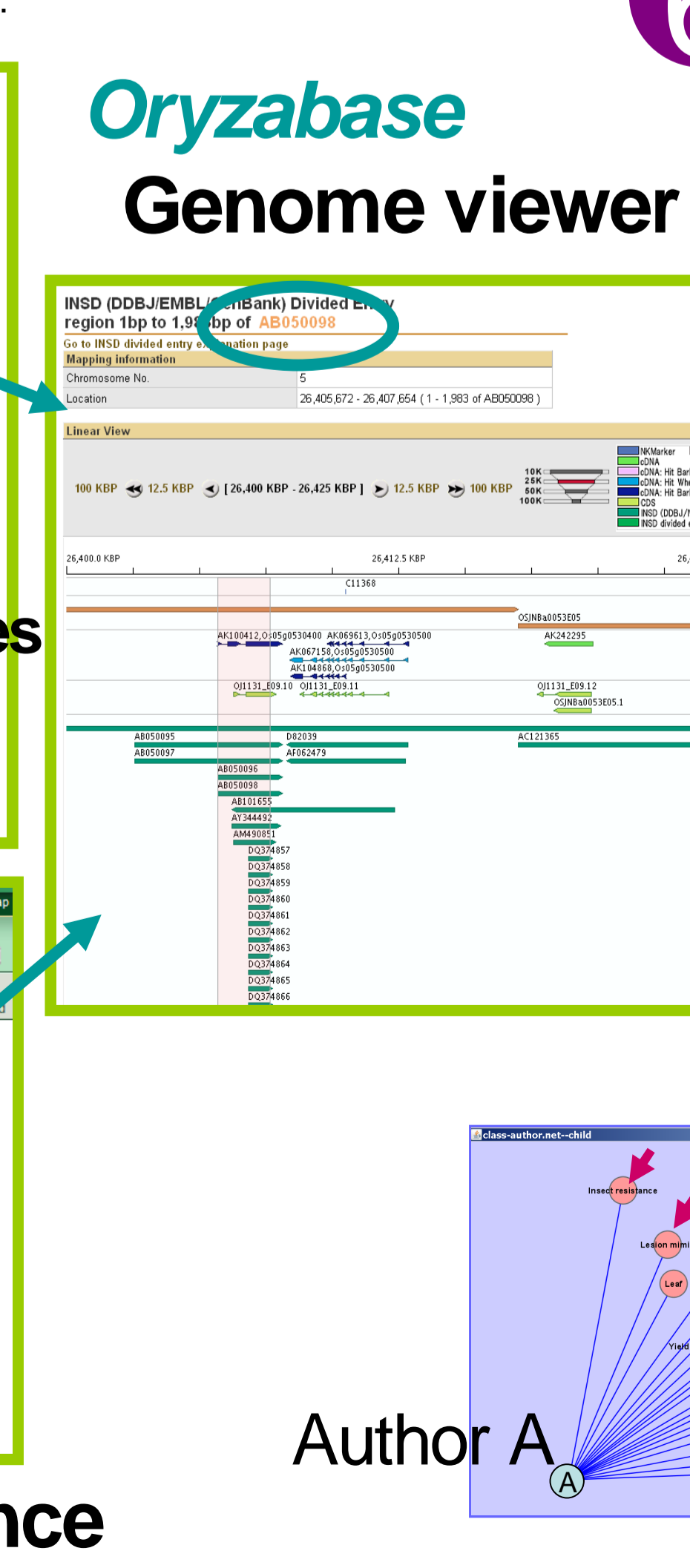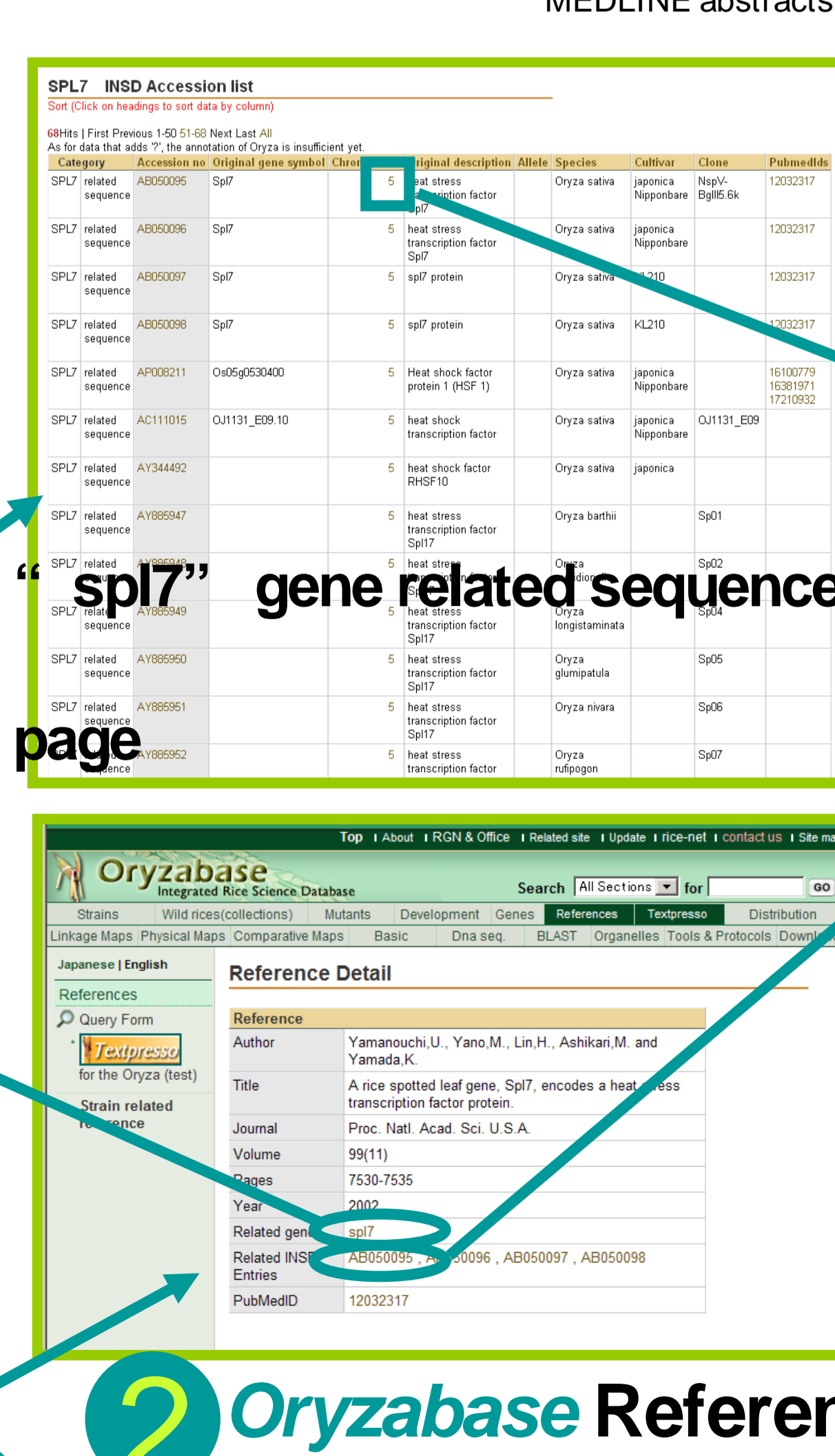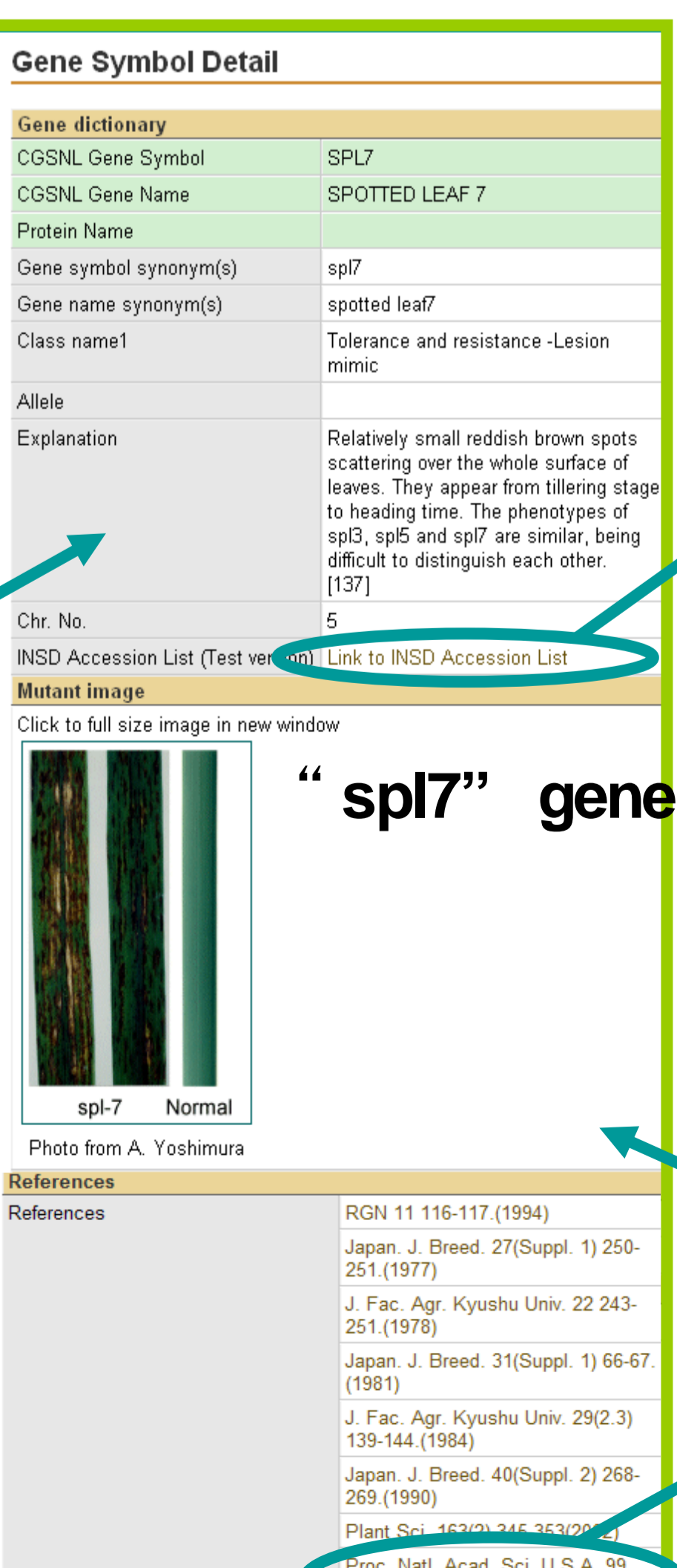 Cell Type

### 6-2 Result of natural language processing

| Journal | Number extracted genes (A) | Manually extracted genes (B) | Machine extracted genes (C) % | Match (D=C/B)% | Precision (D=C/B)% | Recall (E=C/A)% | F-value 2x(DxE)/(D+E) |
|---|---|---|---|---|---|---|---|
| Nature/Science | 127 | 85 | 206 | 67 | 32.5 | 78.8 | 46 |
| TAG* | 103 | 47 | 171 | 33 | 19.3 | 70.2 | 30.3 |
| Plant Cell etc. | 233 | 226 | 656 | 181 | 27.9 | 80.8 | 41.5 |
| Total | 463 | | | | | | |

### 6-3 To improve precision

-- Use correlation of author and gene class
-- Use correlation of author and co-author
-- Improve quality of dictionary
-- Adjust GENIA Corpus to rice related words

correlation of author and gene class

Author A
Author B
Author C
Author D
Author E
Author F

Author specific trait class

Author A
Author B

## 1 *Oryzabase* genes list

Gene Symbol Detail

CGSNL Gene Symbol SPL7
CGSNL Gene Name SPOTTED LEAF 7
Gene symbol synonym(s) spl7
Gene name synonym(s) spotted leaf7
Class name1 Tolerance and resistance -Lesion mimic
Explanation Relatively small reddish brown spots scattering over the whole surface of leaves. They appear from tillering stage to heading time. The phenotypes of spl3, spl5 and spl7 are similar, being difficult to distinguish each other. [137]
Chr. No. 5

"spl7" gene page

"spl7" gene related sequences

## *Oryzabase* Genome viewer

spl-7   Normal
Photo from A. Yoshimura

## 2 *Oryzabase* Reference

Reference Detail
Author Yamanouchi U., Yano.M., Lin.H., Ashikari.M. and Yamada.K.
Title A rice spotted leaf gene, Spl7, encodes a heat stress transcription factor protein

Reserved gene symbol
(before publication)