# Development of a Database of Plant Diseases in Japan and a System for Making Microorganism Genetic Resources and their DNA Sequence Data Available to the Research Community

**Masaru TAKEYA[1]\*, Fukuhiro YAMASAKI[1], Shihomi UZUHASHI[2], Midori KUMAGAI[1], Hiroyuki SAWADA[1], Toshirou NAGAI[1], Keisuke TOMIOKA[1], Toyozo SATO[1], Takayuki AOKI[1] and Makoto KAWASE[1]**

[1] Genetic Resources Center, National Institute of Agrobiological Sciences (Tsukuba, Ibaraki 305–8602, Japan)

[2] Saskatoon Research Centre, Agriculture and Agri-Food Canada (Saskatoon, SK. S7N 0X2, Canada)

## Abstract

A database of plant diseases reported in Japan and a Web-based data-retrieval system to search for information stored in the database have been developed.  The retrieval system is linked with the database on microorganism genetic resources registered in the Genebank Project of the National Institute of Agrobiological Sciences (the NIAS Genebank), which provides access to detailed genetic information on strains of these microorganisms.  The NIAS Genebank manages the registration of so-called "passport" data for these genetic resources and controls their storage to make samples available to researchers worldwide.  Data management software has been developed to register the receipt, storage, testing, and multiplication of these genetic resources.  Search software has also been developed to seek suitable genetic resources based on several search criteria to facilitate their distribution from the stocks stored in the NIAS Genebank.  DNA sequence data for the D1/D2, 18S, and internally transcribed spacer ITS (ITS1–5.8S–ITS2) regions of rDNA and for the coding region of the β-tubulin gene for many of the conserved microorganisms can be accessed from the NIAS Genebank Web site.  The database on Japanese plant diseases and the NIAS Genebank search system are expected to become increasingly useful tools for research and education related to microorganisms and their genetic data, as well as for improving plant protection and food security.

**Discipline:** Information technology
**Additional key words:**  data management, genetic resources, Web-based retrieval system

## Introduction

Plant diseases can be defined based on the relationships between the host plants and the pathogens responsible for the disease symptoms.  Information on plant diseases such as disease names, pathogens, and related references is essential for their diagnosis, control, and treatment in order to sustain plant growth and development.

Massive quantities of data on plant diseases can now be distributed to users worldwide via various Internet technologies, which offer possibilities (e.g. search tools and database links) not otherwise available with printed publications.  Such Web pages have been constructed to record the common names of plant diseases, as in the case of the American Phytopathological Society (APS) Web site[1].  The APS Web site provides a hierarchical menu that links plant names with associated diseases.  The United States Department of Agriculture (USDA) has also developed a retrieval system to search fungal databases for information such as the locations of specimens, literature on hosts and fungi, and other relevant information[2].

The plant diseases reported in Japan have been compiled by the Phytopathological Society of Japan (PSJ) and

published in the book *Common Names of Plant Diseases in Japan*[7]. This book contains a host plant index and the names of diseases that have been reported on the hosts, together with detailed information including the causal pathogens and references to related information. The Genetic Resources Center of the National Institute of Agrobiological Sciences (NIAS) has recently developed a database of plant diseases in Japan[9] based on the book. The database retrieval system offers users a flexible search tool that provides access to the host plants, disease names, and pathogens. Many pathogens registered in the database are also stored in the NIAS Genebank[5] and available for distribution to researchers worldwide.

The NIAS Genebank engages in the exploration, collection, characterization, preservation, and distribution of microorganism genetic resources in collaboration with several sub-banks located at agricultural research institutes and NIAS, as the central bank. These genetic resources cover the following organisms, with a focus on agricultural and food research: filamentous fungi (including mushrooms), yeasts, bacteria, actinomycetes, plant viruses, animal viruses, bacteriophages, protozoa and nematodes. The resources curated by the NIAS Genebank can be used for research and educational purposes, including classification and identification studies, analyses of genetic diversity, and elucidation of their interactions with plants or animals. The genetic resources available for distribution to researchers include strains of *Xanthomonas oryzae* pv. *oryzae*[4] and *Mesorhizobium loti*[3], the entire genome sequences of which have been determined.

Data management software has also been developed to enable researchers to register records of new microorganisms and details of the associated research in the genetic resources database. A search system has also been developed to help researchers find genetic resources stored in the NIAS Genebank by using several criteria.

In this paper, we describe the contents of the database of plant diseases in Japan and give details of its operation. We subsequently introduce the data management software used to control the storage of these genetic resources and find relevant data on DNA sequences. The result is a tool that will provide increasingly strong support for research into microorganisms of interest in agriculture, food management, and related fields of study.

## Materials and methods

### 1. Development of the database

The database of plant diseases in Japan was based on the information contained in *Common Names of Plant Diseases in Japan*[7]. This book may contain details of one or more diseases for a single host plant, while a single disease may also have one or more pathogens. We examined the information in this book to detect hierarchical relationships among the common names of plant diseases. We then designed a relational database schema that accounted for these relationships and that expressed the data structure for the plant diseases appropriately. The database consists of several tables, including those for host plants, plant pathogens, plant disease names, and so on. The table of host plants includes the following columns: host plant code, Japanese name, and English name. The table of plant pathogens has the following columns: pathogen code and scientific name (Latin binomials). The table of plant disease names has the following columns: disease name code, Japanese common name, English name, and synonym. Subsequently, a linkage of the tables of host plants, plant pathogens, plant disease names, and so on were created to establish the basic structure of the plant disease database, by applying the codes designated in the connected tables. Data on their references and related information are also included in the database.

The database was initially stocked with data from the book on plant disease names[7]. A series of addenda (http://www.ppsj.org/mokuroku.html [in Japanese]) published after the initial publication of the book in 2000 has also been added to the database. The database now includes 1,948 host plants and 11,440 common names of associated plant diseases. Eight host plants (or plant groups) (rice, citrus, apple, grape, tomato, pear, Japanese cedar, and tea) have more than 70 common disease names associated with them (Table 1); other species in the database have smaller numbers of entries. The pathogens recorded in the database have been classified into nine groups (Table 2), which are dominated by fungi (including mushrooms and yeasts) (73.4% of the total), nematodes (11.3%), bacteria (including actinomycetes) (6.6%), and plant viruses (5.9%).

### 2. Web-based data retrieval system for the database

We have constructed a Web-based data retrieval system to facilitate searching of the database. The search interface is now available at the NIAS Genebank Web site (the "Database of Plant Diseases in Japan", http://www.gene.affrc.go.jp/databases-micro_pl_diseases_en.php). The retrieval system offers several user-friendly functions. The following examples of using the system are based on version 8 of Internet Explorer; however, the database is also accessible via Firefox 3 (Windows, Macintosh, Linux) and Safari 4 (Windows and Macintosh). To use the database, JavaScript and support for cookies must be enabled.

**Table 1. List of host plants for which more than 70 common names of plant diseases have been proposed**

| Host plant | Number of proposed common names of plant diseases |
|---|---|
| Rice (*Oryza sativa* L.) | 127 |
| Citrus (Citrinae including *Citrus* spp.) | 113 |
| Apple (*Malus pumila* Miller var. *domestica* Schneider) | 102 |
| Grape (*Vitis* spp.) | 86 |
| Tomato (*Lycopersicon esculentum* Mill.) | 79 |
| Pear (*Pyrus* spp.) | 73 |
| Japanese cedar (*Cryptomeria japonica* (Linn. fil.) D. Don) | 73 |
| Tea (*Camellia sinensis* (L.) Kuntze) | 72 |

**Table 2. Proportions of the total number of pathogens in the database in each of the nine categories**

| Pathogen category | Ratio (% of total entries) |
|---|---|
| Fungi (including mushrooms and yeasts) | 73.4 |
| Nematodes | 11.3 |
| Bacteria (including actinomycetes) | 6.6 |
| Plant viruses | 5.9 |
| Phytoplasmas | 0.5 |
| Algae | 0.4 |
| Mites and Insects | 0.2 |
| Viroids | 0.2 |
| Others | 1.5 |



**Fig. 1. Example of starting a search with the plant disease names database**

Data can be retrieved by host plant, disease name, pathogen, or a combination of these criteria. To reduce typing, an "auto-complete" function suggests possible matches as you begin typing in each field. Multiple selections can be specified by using commas, and four kinds of partial match options (i.e. perfect, prefix, suffix, or partial) are available. In an actual search for a plant disease, the host name may start with "rice" (for example) and the disease name may include "leaf" (Fig. 1).

The search result is displayed by clicking the "Search" button. The first row of the results contains the total number of datasets retrieved, an icon to download the search results in Microsoft Excel format (the .XLS format rather than the newer .XLSX format), and a tool to show additional pages when the search results exceed 25 diseases. Basic data on plant diseases are shown in tabular format, including the host plants, disease names, and pathogens. Detailed information is expressed by selecting an individual disease name (Fig. 2). The detailed information window provides links to related microorganism genetic resources stored in the NIAS Genebank.

## 3. Data management

The microorganism section of the NIAS Genebank conserves genetic resources for microorganisms associated with food and agriculture and manages data on their origins and histories (so-called "passport" data), as well as on their storage control. The collected resources are then identified and characterized by experts before becoming part of the NIAS Genebank. Culture and multiplication of these resources ensures their long-term preservation and maintenance under stable conditions, including L-drying (vacuum-drying directly from their liquid phase), lyophilization (vacuum-freeze-drying), or deep-freezing in liquid nitrogen vapor. The microorganism genetic resources are each assigned with unique descriptors called "MAFF numbers", which are the prefix MAFF followed by the strain number. Various types of information on the preserved genetic resources are added to the genetic resources database from time to time and then made publicly available online.
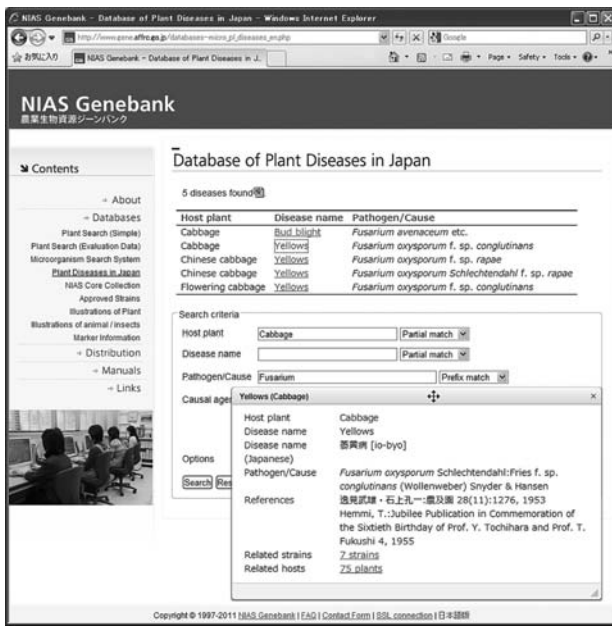
**Fig. 2.** Example of the kind of search result provided by the Web-based data retrieval system after a search of the plant disease database

The storage control processes used by the microorganism section of the NIAS Genebank are also recorded in the database. We have developed data management software to control the input, storage, and updating of data in the database and organisms in the bank. Figure 3 is a schematic diagram of the control process; rounded rectangles denote specific data management programs. The receipt program ("P. receipt") registers passport data such as MAFF numbers, the scientific names of the microorganisms, their prior designations, and their origins or sources when microorganisms are isolated from natural habitats and added to the collection of genetic resources. The storage program ("P. storage") registers storage control data that specifies the storage containers and locations and numbers of stored materials inside them. For quality control, the macro- and microscopic features of the individual preserved microorganism stocks are examined to determine their rates of degeneration and possible contamination, as well as to confirm their survival rates during preservation. The test program ("P. test") is applied to record these survival rates. When the quantity of
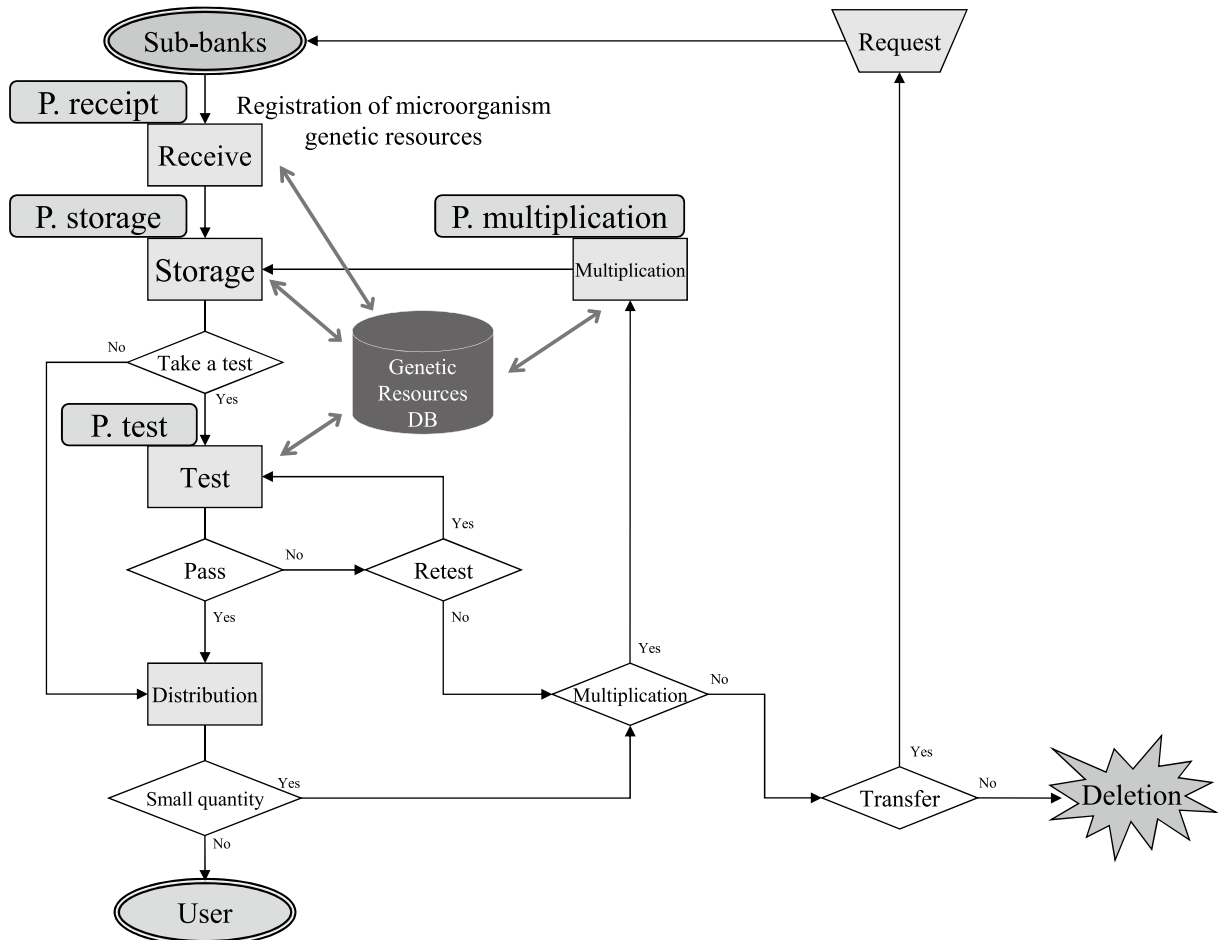


**Fig. 3.** Schematic diagram of the storage control system used to manage the microorganism genetic resources stored in the NIAS Genebank

a stock falls below a pre-determined minimum value, the survival rate is too low, or the stock quality is insufficient to permit distribution due to contamination or other problems, distribution of the resource is tentatively stopped. The multiplication program ("P. multiplication") is also used to register storage control information related to multiplication of the stocks. When the main stocks preserved in the central bank at NIAS are of insufficient quality for distribution, requests can be sent to the subbanks at other institutes participating in this program to transfer new samples of original stocks to the central bank at NIAS.

## 4. Data search system

The passport data for the microorganism genetic resources and related information are provided by the NIAS Genebank Web site (http://www.gene.affrc.go.jp/databases-micro_search_en.php). We have constructed a search system to seek genetic resources conserved in the NIAS Genebank by using the following search criteria:

- Category: Search by groups of microorganisms such as "fungi" or "bacteria". Multiple categories can be selected simultaneously.
- MAFF No.: Search by the 6-digit MAFF strain IDs. A range can be specified using commas (which mean "or") and hyphens (which mean "between"). Examples: "118089, 118147", "211368-211370".
- Scientific names: Search by the Latin binomials of the microorganisms. An auto-complete function suggests possible matches for entries as you begin typing in this field. Multiple selections are possible by using commas and a partial match option. Examples: Xantho [Prefix match], rhizobium [Partial match].
- Designation: Search by prior strain designations described in earlier research literature. Multiple selections are possible by using commas and a partial match option. Examples: S96, TNPG [Prefix match].
- Sources: Search by the origin or source of the isolate (scientific names and major common names). Multiple selections are possible by using commas and a partial match option. Examples: Aspergillus, malt [Partial match].
- Locations: Search by the collection sites of the resource. Multiple selections are available by clicking on additional location names while holding down the Ctrl key. Extension of a selection to include a range of additional locations is done by selecting the first location, then clicking on the final location name in the range while holding down the Shift key.
- Properties: Search by properties of the resource. Individual specifications are made by clicking while holding down the Ctrl key, and a range of entries is

selected by clicking while holding down the Shift key, as described above for locations. Optional switches can also be used to search for approved, (ex-) type, and reference strains, or strains for which DNA sequence data are available in the database.
- Search button: Clicking this button starts the database query. Each search criterion is connected by "AND" (which represents an intersection of the dataset for each criterion). Empty fields are ignored.

The first row of search results contains the total number of microorganism genetic resources that match the search criteria, together with an icon to download the results in Microsoft Excel .XLS format (up to 1000 items), and a tool for moving to the next page when the search result exceeds 25 items. The data are summarized in tabular form. The results can be sorted in ascending order by clicking the column header that specifies the sort parameter. Detailed information is displayed by clicking on the MAFF number for each result. The abbreviations in the remarks column in the summary table have the following meanings:

- A: approved strain
- R: reference strain
- T: (ex-) type strain
- L: literature information available
- S: sequence datum or data available.

For part of the microorganism genetic resources stored in the NIAS Genebank, DNA sequence data are available for the D1/D2, 18S, and internally transcribed spacer (ITS) regions of rDNA and for the gene that encodes β-tubulin. To facilitate use of the sequence data for taxonomy and identification[6, 8], the search system was also designed to supply the available sequence data in the Multi-FASTA format (Fig. 4). Single-FASTA files are also available from the individual information window (Fig. 5).

## Results and discussion

A database of plant disease information has been constructed for common plant diseases in Japan[7], and a Web-based data retrieval system has been developed to facilitate searches for data on the same. Unlike the hierarchical menu system that links plant names with associated diseases, this system supports the more flexible retrieval of data by supporting arbitrary query combinations that combine search criteria based on the names of the host plant, disease, and pathogen. Table 1 lists eight host plants (or plant groups) for which more than 70 plant diseases are listed; rice, one of the key food crops in Japan, appears at the top of the list, with 127 entries. Table 2 summarizes the pathogen groups with the greatest repre-
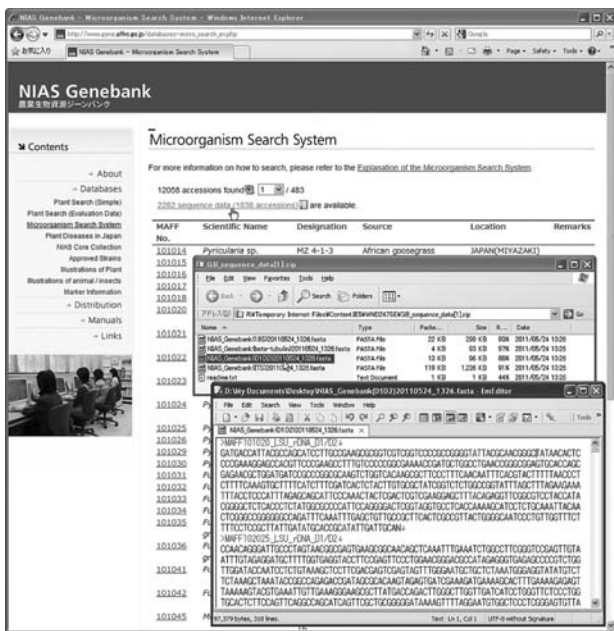
**Fig. 4.** DNA sequence data for the ITS, 18S, and D1/D2 regions of rDNA and for the gene that encodes β-tubulin can be downloaded from the search results in Multi-FASTA format



**Fig. 5.** Single-FASTA data for the ITS, 18S, and D1/D2 regions of rDNA and for the gene that encodes β-tubulin can be downloaded from the detailed information window in the search results

sentation in the database. Fungi, including mushrooms and yeasts, are the most frequently reported category of plant pathogens in Japan.

Although the current database and search tools represent a promising start, much room for future improvements to the system remains. Improvement of the system used to store the microorganism genetic resources may be requested by the microorganism section of the NIAS Genebank. A BLAST search tool for use with the stored DNA sequence data will be included in a future system update, while links between the plant disease database and actual research data concerning the pathogenicity of the microorganism genetic resources (e.g. based on data obtained from inoculation tests) may become available in the near future. We foresee that this system will become an increasingly useful tool as these features and new disease data are added.

## Acknowledgments

## References

1. American Phytopathological Society: Common Names of Plant diseases. http://www.apsnet.org/publications/commonnames/Pages/default.aspx.
2. Farr, D.F., & Rossman, A.Y. (2011) Fungal Databases, Systematic Mycology and Microbiology Laboratory, United States Department of Agriculture, Agricultural Research Service. http://nt.ars-grin.gov/fungaldatabases/index.cfm.
3. Kaneko, T. et al. (2000) Complete genome structure of the nitrogen-fixing symbiotic bacterium *Mesorhizobium loti*. *DNA Res.*, **7**, 331–338.
4. Ochiai H, et al. (2005) Genome sequence of *Xanthomonas oryzae* pv. *oryzae* suggests contribution of large numbers of effector genes and insertion sequences to its race diversity. *JARQ*, **39**, 275–287.
5. Okuno, K. et al. (2005) Plant genetic resources in Japan: platforms and destinations to conserve and utilize plant genetic diversity. *JARQ*, **39**, 231–237.
6. Peterson S.W. & Kurtzman C.P. (1991) Ribosomal RNA sequence divergence among sibling species of yeasts. *System. Appl. Microbiol.,* **14**, 124–129.
7. Phytopathological Society of Japan (2000) *Common Names of Plant Diseases in Japan: First edition*. Japan Plant Protection Association, Tokyo, Japan, pp.857 [In Japanese].
8. Sugita T. et al. (1999) Identification of medically relevant *Trichosporon* species based on sequences of internal transcribed spacer regions and construction of a database for *Trichosporon* identification. *J. Clin. Microbiol.,* **37**, 1985–1993.
9. Takeya, M. et al. (2011) NIASGBdb: NIAS Genebank databases for genetic resources and plant disease information. *Nucleic Acids Res.*, **39**, D1108–D1113.