

# DNA Chip Database

## for the Korean Functional Genomics Project

Sangsoo Kim, Ph.D.

Korea Research Institute of Bioscience & Biotechnology

E-mail: sskimb@mail.kribb.re.kr

### **Abstract**

The Korean Functional Genomics Project focuses on stomach and liver cancers. Specimens collected by six hospital teams are used in DNA microarray experiments. Experimental conditions, spot measurement data, and the associated clinical information are stored in a relational database. Microarray database schema was developed based on EBI's ArrayExpress. A diagrammatic representation of the schema is used to help navigate over many tables in the database. Field description, table-to-table relationship, and other database features are also stored in the database and these are used by a PERL interface program to generate web-based input forms on the fly. As such, it is rather simple to modify the database definition and implement controlled vocabularies. This PERL program is a general-purpose utility which can be used for inputting and updating data in relational databases. It supports file upload and user-supplied filters of uploaded data. Joining related tables is implemented using JavaScripts, allowing this step to be deferred to a later stage. This feature alleviates the pain of inputting data into a multi-table database and promotes collaborative data input among several teams. Pathological finding, clinical laboratory parameters, demographical information, and environmental factors are also collected and stored in a separate database. The same PERL program facilitated developing this database and its user-interface.

## CV

- 2000년-현재: 생명공학연구소 유전체연구센터 책임연구원, 과학기술부 21세기 프론티어사업 “인간유전체기능연구사업단” Bioinformatics 책임자
- 1999년-2000년: (주)LG화학 책임연구원, Bioinformatics 팀장
- 1995년-1998년: (주)LG화학 책임연구원, Biopharmaceutical Design 팀장
- 1988년-1995년: (주)LG화학 선임연구원, Biopharmaceutical Design 팀장
- 1986년-1988년: Purdue University, Postdoctoral Res. Ass., Dept. of Bio. Sci.
- 1983년-1986년: Iowa State University, Ph.D. in Physical Chemistry
- 1981년-1983년: Seoul National University, M.S. in Physical Chemistry
- 1977년-1981년: Seoul National University, B.S. in Chemistry

# **cDNA Microarray Database for Korean Functional Genomics Project**

Oct. 2001

Korea Research Institute of  
Bioscience & Biotechnology  
Sangsoo Kim

## **Scope and Objectives of Korean Function Genomics**

- **Establishment of basic infrastructure for Functional Genomics**
  - HT sequencing, DNA chip, Proteomics, bfx
- **Discovery of gene markers associated with diseases common in Korean population**
  - Diagnostic and prognostic markers of stomach or liver cancer

## Phase I: Infrastructure

### 1st year

- Tissue collection
- **HT sequencing**
- Korean-type SNP
- **DNA microarray**
- Proteomics
- **Bioinformatics**

### 2nd year

- Cell-based assay
- Model organisms
- KO/tg mice
- Structural genomics
- Chemical genomics
- Protein interaction

## Bioinformatics Mission

### v Data Archiving

- Scalable, Secure, Flexible, Affordable

### v Analysis Support

- State-of-the-art algorithms, Web-GUI

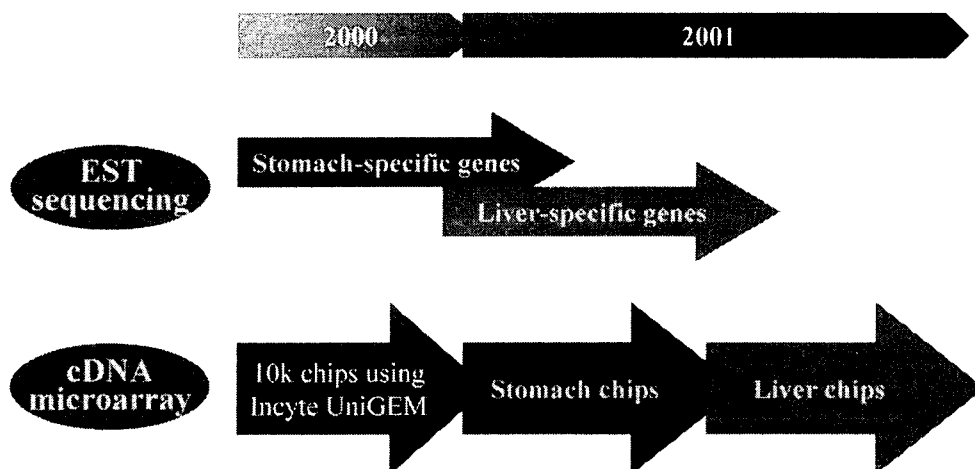
### • Sharing Information

- Data integration, Performance

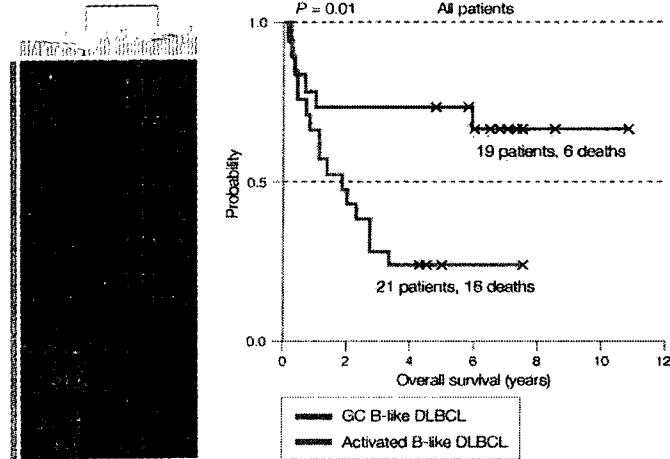
## Data Archiving

- v HT sequencing (1 team)
  - EST sequencing of cancer cDNA libraries
- v DNA microarray (2 teams)
  - 2.5k chips (now), 10k (end of this year) chips
- v Tissue collection (6 teams)
  - 3 teams for stomach and liver cancers, respectively
- Korean-type SNPs (5 teams)
  - Scoring known markers, Discovering cSNPs
- Proteomics (2 teams)
  - 2-D Gel images and processed data

## Tissue-specific UniGene Chips



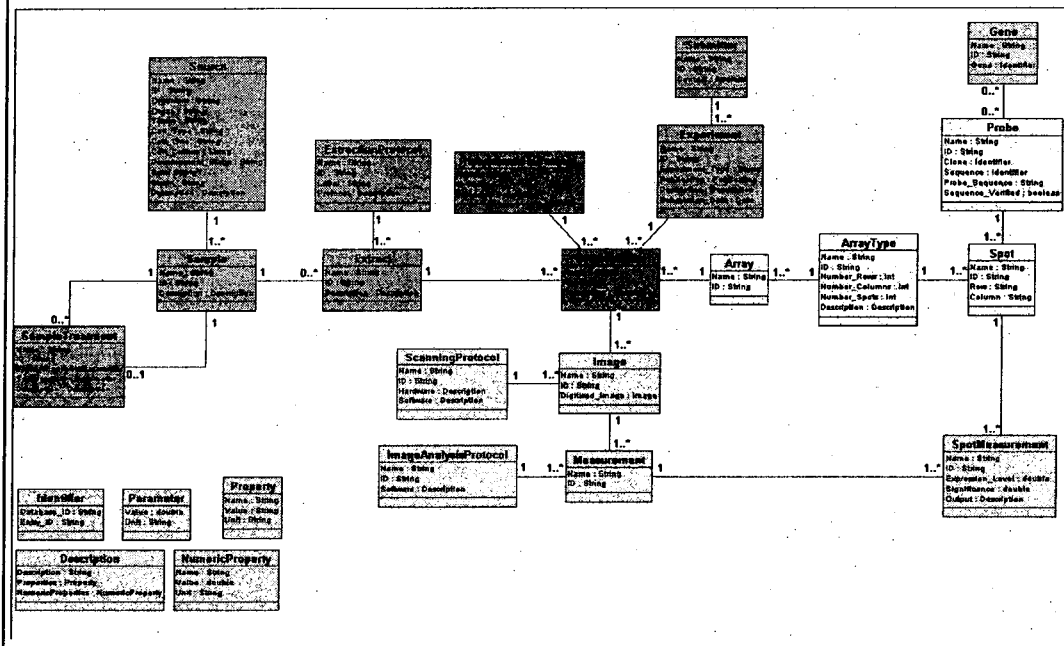
## Discovery of diffuse large B-cell lymphoma (DLBCL) clinical subtypes by DNA array profiling



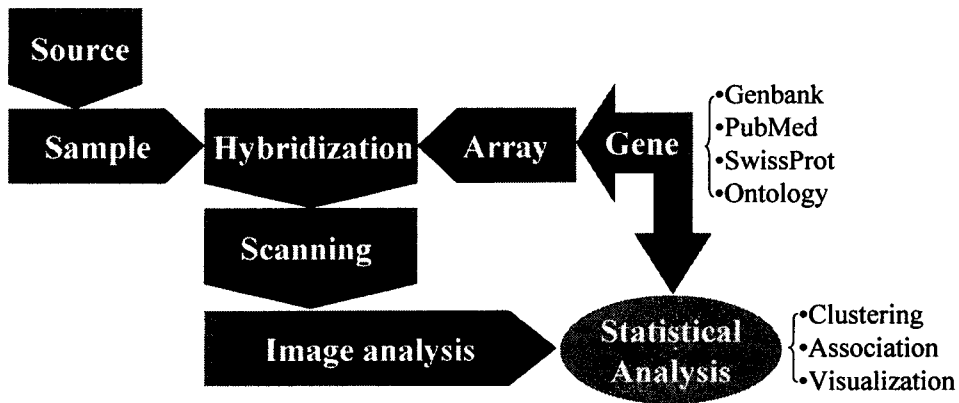
Nature Reviews Genetics 1: 48-56 (2000)  
MOLECULAR PROFILING OF HUMAN CANCER

Nature Reviews | Genetics

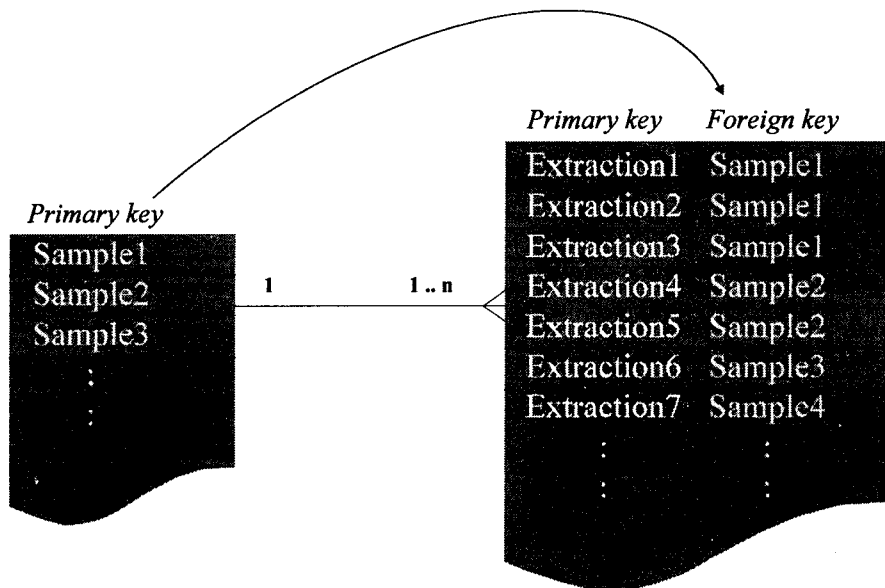
## EBI Microarray Database Schema (www.ebi.ac.uk/arrayexpress)



# Microarray (Processes)



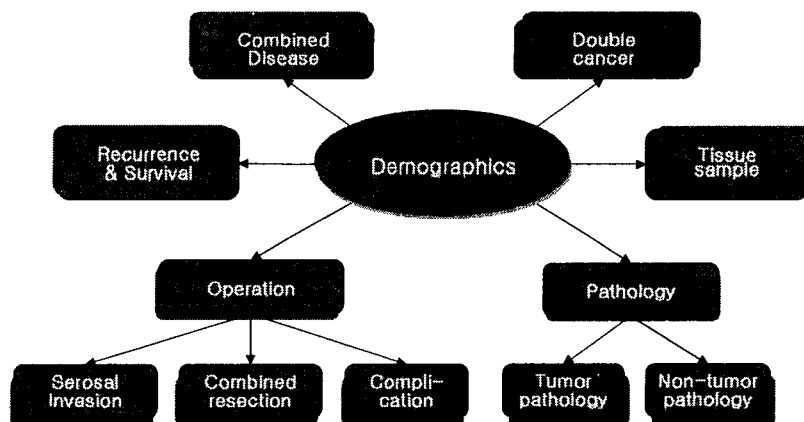
One "Sample" can be used several times for "Extraction"



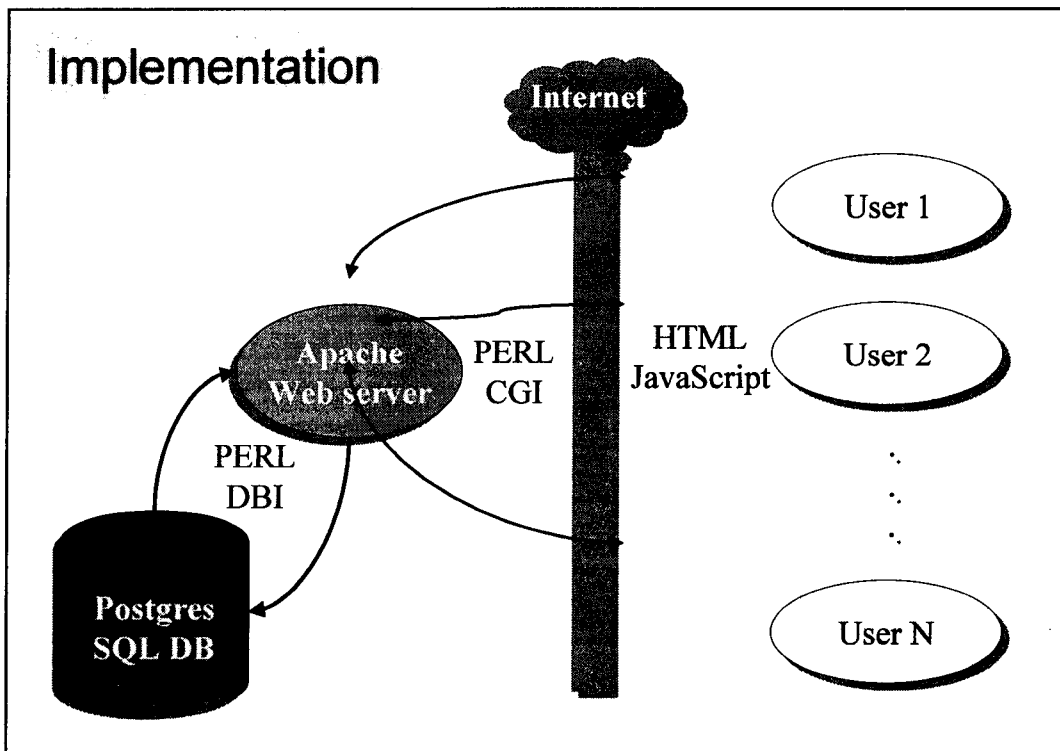
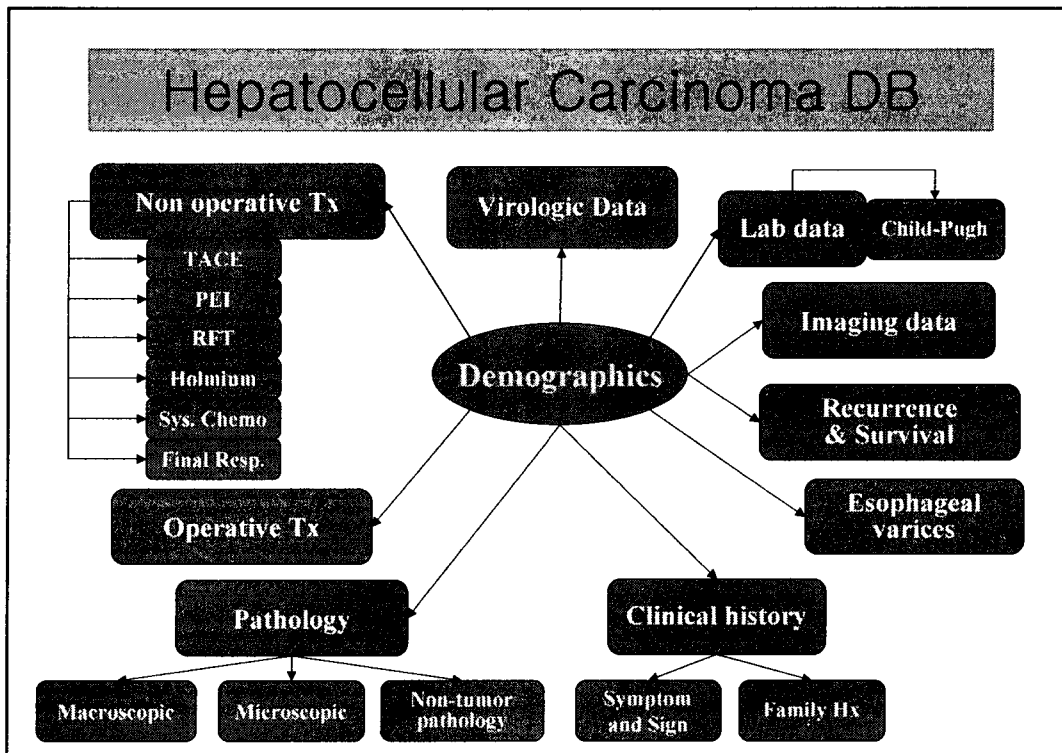
# Clinical Tissue Database

- Epidemiological parameters
  - sex, age, education, lifestyle, job, home
  - family history
  - prior diseases, hepatitis
  - dietary habit, smoking, alcohol, occ. hazard
  - physical exam., clinical lab.
  - lifespan
- Pathological parameters
  - surgery, tumor size, stage, type, invasion, lymphocyte infiltration

## Gastric Carcinoma DB







# Database Features

- User authentication
  - cookie-based
- Navigation via ‘map diagram’
  - Sequential input not necessary; cooperative input
- Table definitions in DB
  - Flexible design change; controlled vocabulary
- Table-table relationships in DB
  - Copy values from upstream to downstream tables
- File-upload for long tables with filters
- Link-out to external URLs
- General-purpose DB interface

## Table definitions stored in DB – easy design change, controlled

TABLE	FIELD	TYPE	DESCRIPTION
demograp	hospital_code	character varying(10)	Hospital Code
demograp	ch_c	character varying(10)	Ch
demograp	social_no	character varying(14)	Social No
demograp	name_c	character varying(10)	Name
demograp	age_c	int2	Age
demograp	sex_c	character	Sex { M, F }
demograp	date_c	date	Enr_b
demograp	staff_c	int2	user ID
symp	name	character varying(10)	patient name
symp	date	date	Enr_b
symp	staff	int2	user ID
symp	abdominal_pain	int2	Abdominal pain { 1 (Yes), 2 (No) }
symp	abdominal_distension	int2	Abdominal distension { 1 (Yes), 2 (No) }
symp	weight_loss	int2	weight loss { 1 (Yes), 2 (No) }
symp	weakness	int2	weakness { 1 (Yes), 2 (No) }
symp	jaundice	int2	Jaundice { 1 (Yes), 2 (No) }
symp	vomiting	int2	Vomiting { 1 (Yes), 2 (No) }
symp	palpable_mass	int2	Palpable mass { 1 (Yes), 2 (No) }
symp	fever	int2	Fever { 1 (Yes), 2 (No) }
symp	other	character varying(32)	other
symp	no_symptom	int2	No symptom { 1 (Yes), 2 (No) }
family_hx	name	character varying(10)	patient name
family_hx	date	date	Enr_b
family_hx	staff	int2	user ID
family_hx	hx_of_chronic_liver_diseases	int2	Hx of chronic liver diseases { 1 (Yes), 2 (No) }
family_hx	hcc_hx	int2	HCC Hx { 1 (parents), 2 (sibling), 3 (offspring) }
family_hx	hbv_hx	int2	HBV Hx { 1 (parents), 2 (sibling), 3 (offspring) }
family_hx	hcv_hx	int2	HCV Hx { 1 (parents), 2 (sibling), 3 (offspring) }
clinical_history	name	character varying(10)	patient name
clinical_history	date	date	Enr_b
clinical_history	staff	int2	user ID
clinical_history	duration_of_hbv_hx	float4	Duration of HBV hx (years), 0 for none
clinical_history	duration_of_hcv_hx	float4	Duration of HCV hx (years), 0 for none
clinical_history	duration_of_cirrhosis	float4	Duration of cirrhosis (years), 0 for none
clinical_history	transfusion_hx	float4	Transfusion Hx (years), 0 for none

## Navigation via 'map diagram'

- Sequential input not necessary; collaborative input

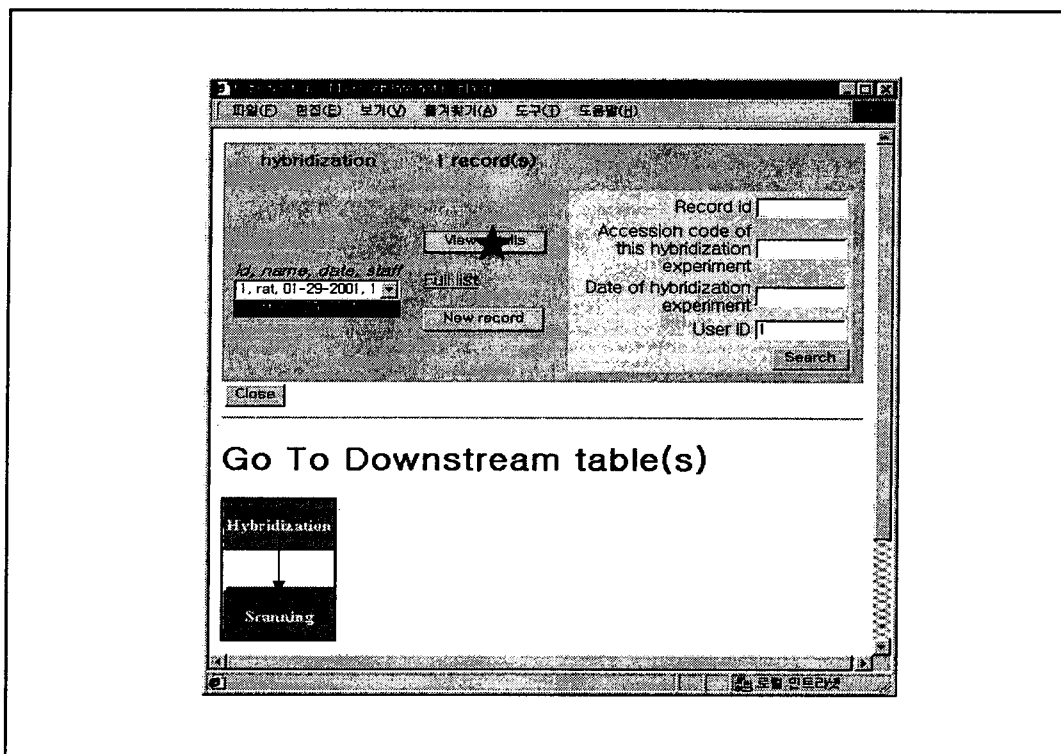
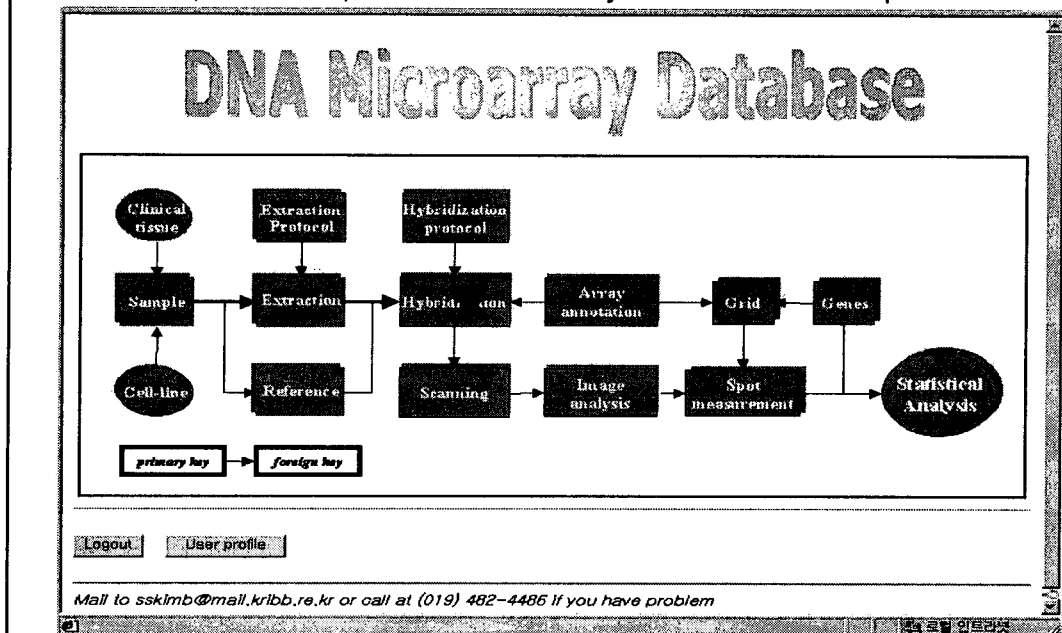


TABLE: hybridization, RECORD # 1

Field description	Field value
Accession code of this hybridization experiment	rat <input type="button" value="UNIQUE ?"/>
Array annotation table ID	1 <input type="button" value="Browse"/>
Array serial number in this batch	
RNA extraction table ID for the test sample	<input type="button" value="Browse"/>
Dye for the test sample	<input checked="" type="radio"/> Cy5 <input type="radio"/> Cy3
Reference sample table ID	<input type="button" value="Browse"/>
Dye for the reference sample	<input checked="" type="radio"/> Cy3 <input type="radio"/> Cy5
Hybridization protocol table ID	<input type="button" value="Browse"/>
Internal standard	<input checked="" type="radio"/> None <input type="radio"/> Included
Date of hybridization experiment	01-29-2001
User ID	
Outline and scope of this experiment	<input type="button" value="Browse"/>

array\_annotation: 2 record(s)

id, name, date, size	Record id	Accession code of this array record	Date of array printing	User ID
2, hyu_hum_01, 01-27-2001, 3				

Copy values in this Record to hybridization table

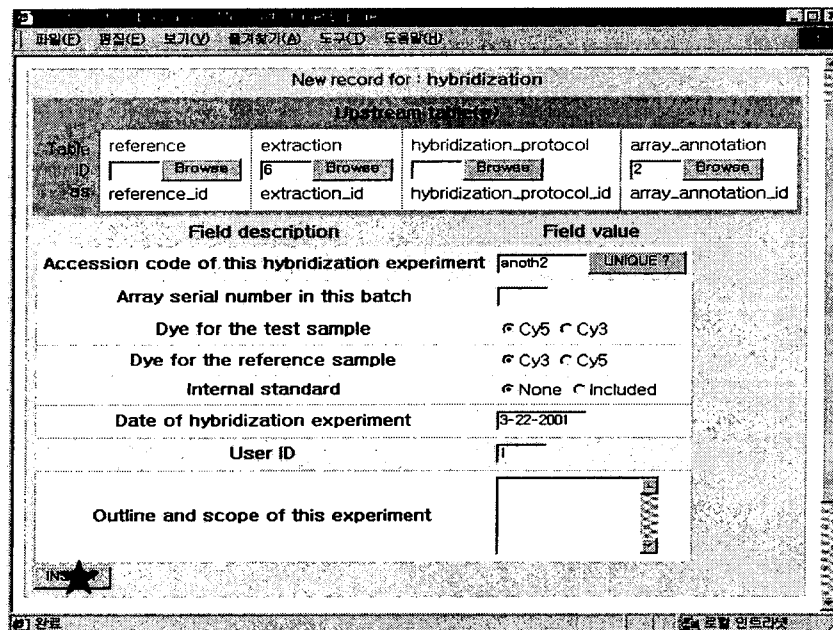
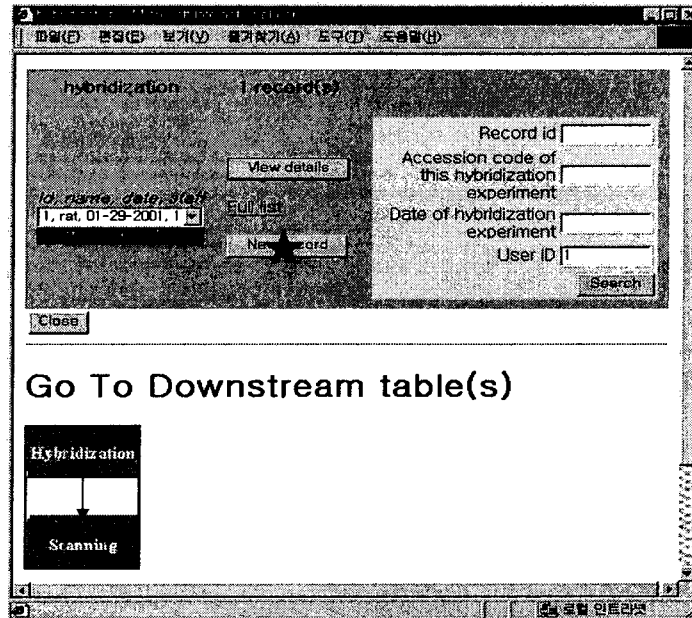
TABLE: array\_annotation, RECORD # 1

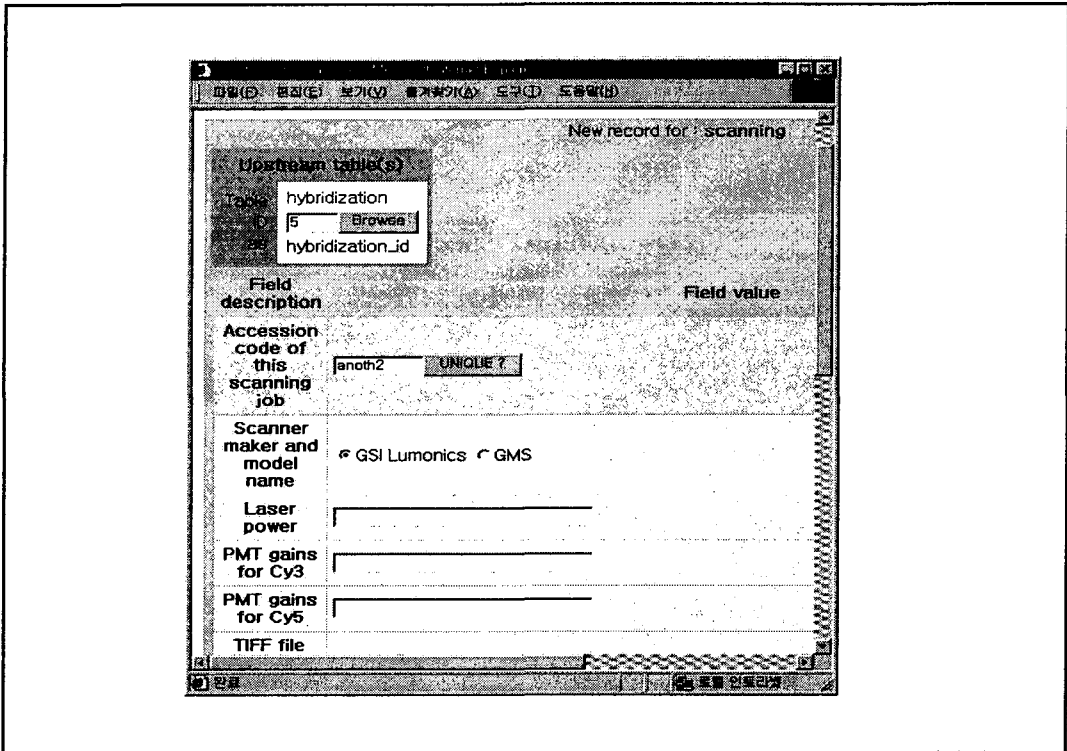
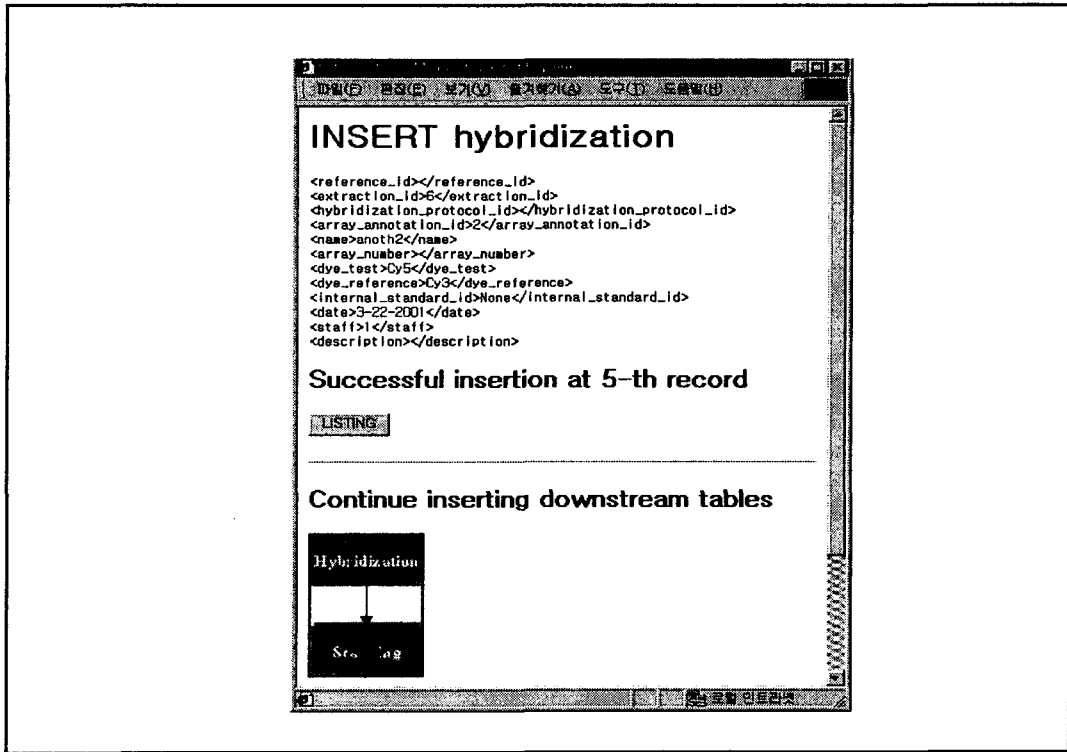
Field description	Field value
Accession code of this array record	hyu_ral.1. <input type="checkbox"/> UNIQUE ?
Number of block columns	4
Number of block rows	4
Total number of blocks	16
Number of columns in each block	16
Number of rows in each block	15
Total number of spots in each block	240
Total number of spots in an array	3840
Diameter of a spot in um	
Distance between centers of	

TABLE: hybridization, RECORD # 1

Field description	Field value
Accession code of this hybridization experiment	rat <input type="checkbox"/> UNIQUE ?
Array annotation table ID	<input type="text"/> <input type="button" value="Browse"/>
Array serial number in this batch	<input type="text"/>
RNA extraction table ID for the test sample	<input type="text"/> <input type="button" value="Browse"/>
Dye for the test sample	<input checked="" type="radio"/> Cy5 <input type="radio"/> Cy3
Reference sample table ID	<input type="text"/> <input type="button" value="Browse"/>
Dye for the reference sample	<input checked="" type="radio"/> Cy3 <input type="radio"/> Cy5
Hybridization protocol table ID	<input type="text"/> <input type="button" value="Browse"/>
Internal standard	<input checked="" type="radio"/> None <input type="radio"/> Included
Date of hybridization experiment	01-28-2001
User ID	<input type="text"/>
Outline and scope of this experiment	<input type="text"/>

UPDATE this record    INSERT this as a new record





Field description	Field value
Accession code of this image analysis job	7 <input type="button" value="Browse"/>
Name and version number of image analysis program	MAAS GenePix 3.0 Image Quantarray ScanAnalyzer
Signal range accepted (eg. 15-85%)	
JPEG file name of the synthetic image	
GIF file name of the log-scatter plot	
Excel file name of the data processing	
Spot measurement file name	ka_6h_-78-40_data-1.txt <input type="button" value="SPOTS"/>

Controlled vocabulary

```

INSERT image_analysis

<scanning_id>7</scanning_id>
<name>anoth2</name>
<program>MAAS</program>
<signal_range></signal_range>
<image_path></image_path>
<scatter_path></scatter_path>
<excel_path></excel_path>
<spot_measurement>C:\3A\5Cribb\is\5Cnew\5CXC7X01\XBEXE7XB4\EBX5Cdma\5Cka_6h_-78-40_data-1.txt</spot_measurement>
<date>3-22-2001</date>
<staff>1</staff>
<normalization_program></normalization_program>
<normalization_algorithm1>All spots</normalization_algorithm1>
<normalization_algorithm2>No</normalization_algorithm2>
<normalization_algorithm3>No</normalization_algorithm3>
<description></description>

File uploading C:\Wkribbmis\Wnew\한양대\Wdma\Wka_6h_-78-40_data-1.txt

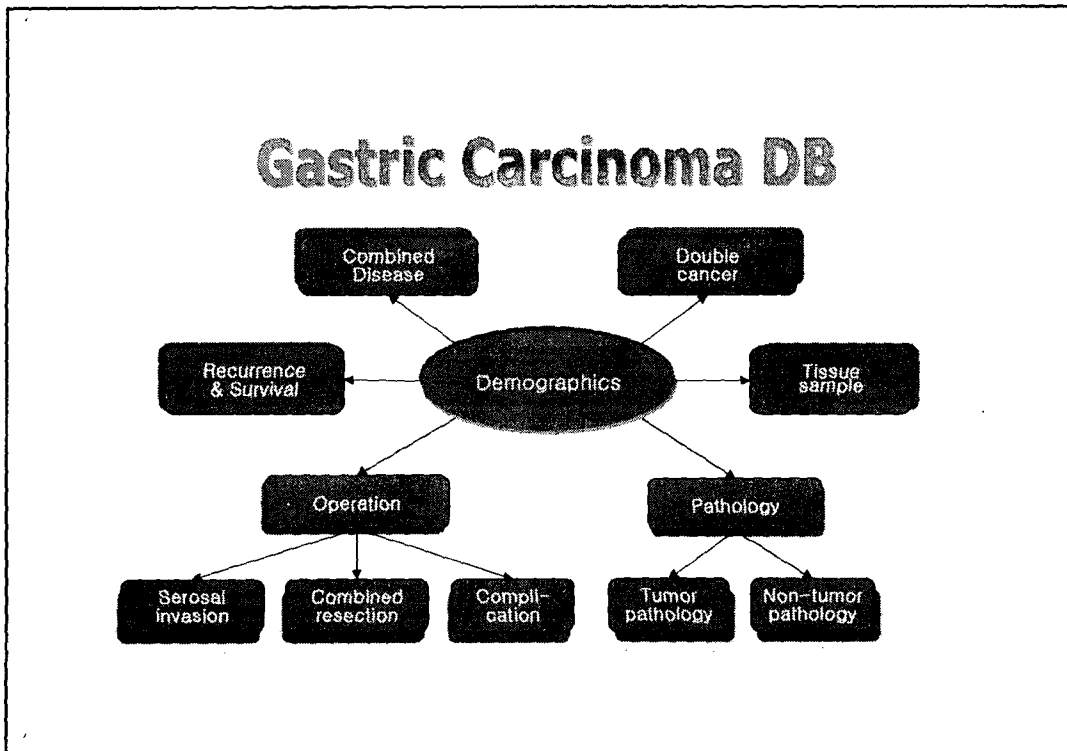
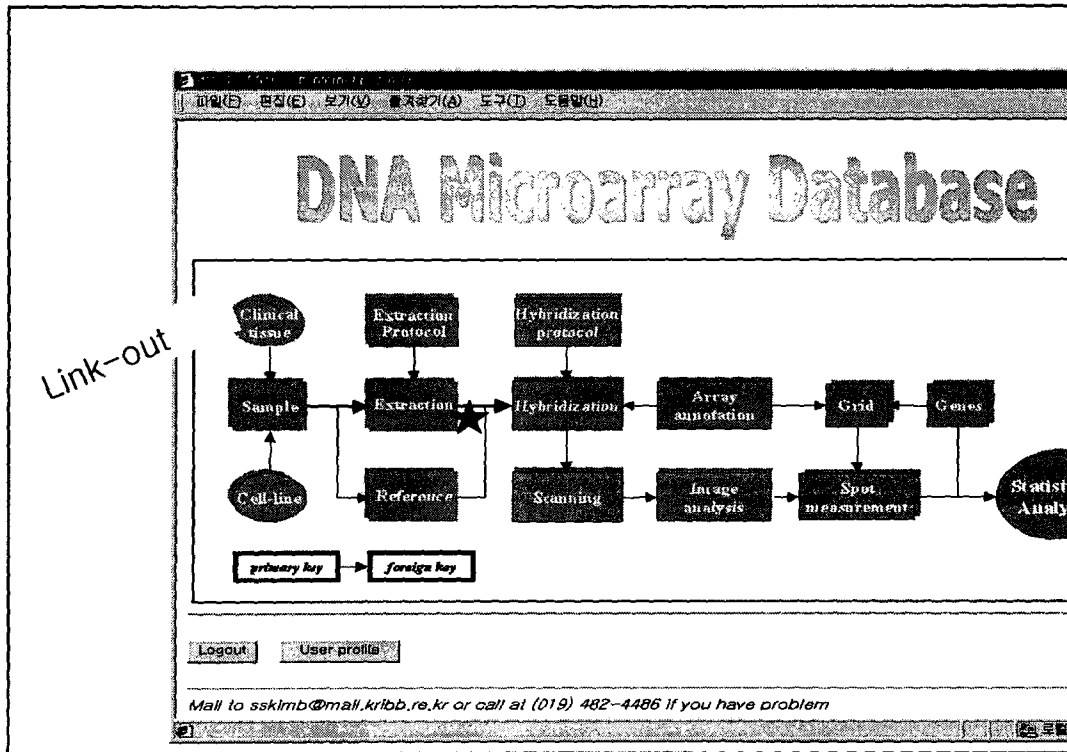
Successful insertion at 4-th record

LISTING

No further downstream tables

```





**Hierarchical Clustering of Stomach Cancer Samples**

- 20 normal+tumor pairs

- courtesy of Inchul Lee, Asan Medical Center

Normal    Cancer+*normal*

