# KGG: A systematic biological Knowledge-based mining system for Genome-wide Genetic studies (Version 3.5)

# User Manual

*Miao-Xin Li, Jiang Li*

Department of Psychiatry
Centre for Genomic Sciences
Department of Biochemistry
The University of Hong Kong
Pokfulam, Hong Kong SAR, China

# Content

---

**Hints for large GWAS dataset** (around or over 2.5 million SNPs)

Set or change large memory for KGG3 say, 2000MB, by *Tools->Set System Memory*.

---

# 1. Introduction and general pipeline

KGG (Knowledge-based mining system for Genome-wide Genetic studies) is a software tool to perform knowledge-based analysis for genome-wide association studies (GWAS). At present, the version 3 has been equipped with main functions to conduct multivariate/univariate gene-based association tests using SNP p-values from GWAS[1,2,3] and to carry out advanced univariate biological module-based association analysis (pathway enrichment and protein-protein interaction (PPI) network association) by a set-based test [2]. In addition, KGG has provided direct hyperlinks to several useful bioinformatics annotation databases on sequence variants

([http://jjwanglab.org/gwasrap](http://jjwanglab.org/gwasrap)), genes (GeneCards, [http://www.genecards.org/](http://www.genecards.org/)) and pathways (MsigDB, [http://www.broadinstitute.org/gsea/msigdb](http://www.broadinstitute.org/gsea/msigdb)). A number of functions to model emerging epigenomic regulatory data for prioritizing association signals are still under development.
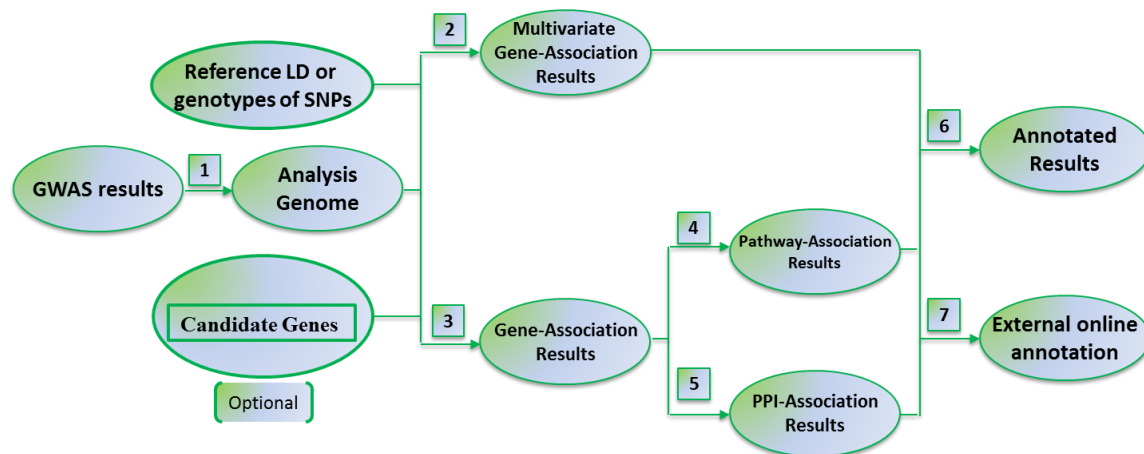


*Figure 1.1 Pipeline chart of KGG analysis* (version 3)

Notes: Circle nodes stand for data and files (input, output); single directional arrows stand for analytical procedures involved.

Main steps involved:

1) Build an **Analysis Genome**: generate an intermediate dataset which integrates original GWAS p-values, SNP annotation and gene annotation, and LD between SNPs WITIN genes together. It is a unified dataset which will be used for all kinds of analyses on KGG.
2) Conduct a **multivariate gene-based association test**: calculate gene-based p-values of multiple phenotypes by a method [3].
3) Conduct **gene-based association test**: calculate gene-based p-values of a single phenotype by GATES[1] or HYST[2].
4) Explore significantly **associated pathways** by HYST [2] and enriched with susceptibility genes by hypergeometric distribution test. One can use either the integrated pathways (gene sets) from MsigDB ([http://www.broadinstitute.org/gsea/msigdb](http://www.broadinstitute.org/gsea/msigdb)) or his or her self-customized pathways on KGG.
5) Explore statistically significant **associated PPI pairs** by HYST [2] which may work together to contribute to the development of the disease or traits. Again, one can use either the integrated PPI pairs from the STRING PPI ([http://string-db.org/](http://string-db.org/)) or his or her self-customized PPI pairs on KGG.
6) Annotate and export significant SNPs, genes, pathways and PPIs.
7) View **external bioinformatics annotation results** of statistically significant SNPs, genes and pathways.

Other plug-in:

1) SPS: a simulation tool for **calculating power of set-based** genetic association tests.

# References

1. Li MX, Gui HS, Kwan JS, Sham PC.  GATES: A rapid and powerful gene-based association test using extended Simes procedure. Am J Hum Genet. 2011 Mar 11;88(3):283-293.
2. Li MX*, Kwan JS*, Sham PC. HYST: A hybrid set-based test for genome-wide association

studies, with application to protein-protein interaction-based association analysis. Am J Hum Genet. 2012 Sep 7;91(3):478-88.
3. Sluis et al. MGAS: a powerful tool for multivariate gene-based genome-wide association analysis. Bioinformatics (In press)

# 2.  Installation

### 2.1 Installation of Java Runtime Environment (JRE)

The Java Runtime Environment (JRE) v1.7 (or higher version) is required to run KGG3 on any operating systems (OS). It can be downloaded from http://java.sun.com/javase/downloads/index.jsp for free. Installing the JRE is very easy in Windows OS and Mac OS X.

In Linux, you have more work to do. Details of the installation can be found at http://www.java.com/en/download/help/linux_install.xml.

In Ubuntu, if you have an error message like: "Exception in thread "AWT-EventQueue-0" java.awt.HeadlessException …", then please installs the Sun Java Running Environment (JRE) first.

> To install the Sun JRE on Ubuntu(10.04), please use the following commands:
> sudo add-apt-repository "deb http://archive.canonical.com/ lucid partner"
> sudo apt-get update
> sudo apt-get install sun-java7-jre sun-java7-plugin sun-java7-fonts
>
> Detailed explanation of above commands can be found at
> http://www.ubuntugeek.com/how-install-sun-java-runtime-environment-jre-in-ubuntu-10-04-lucid-lynx.html.

> Note: After completing Java installation please make sure that not only the java is executable but the extracted jre/bin directory is added to the PATH, otherwise KGG3 would not start properly. This is easily achievable by executing the following command on the terminal:
>
>  echo 'export PATH=/path/to/installed/jre/bin:$PATH' >> ~/.bashrc && source ~/.bashrc
>
> Thanks **Attila Pulay** for the suggestion!

### 2.2 Installation of KGG

To simplify the installation, we still keep KGG as a green tool (i.e., no formal installation procedure guided by an installation wizard). After decompressing the kgg3.zip file, you will see a "bin" folder where there are 3 script files to initiate KGG3. On Microsoft Windows, please double click kgg3.exe or kgg364.exe file.  On Linux, Mac OS X and Solaris, please type the kgg3 in a Command-line Terminal.

# 3. Interface and functions

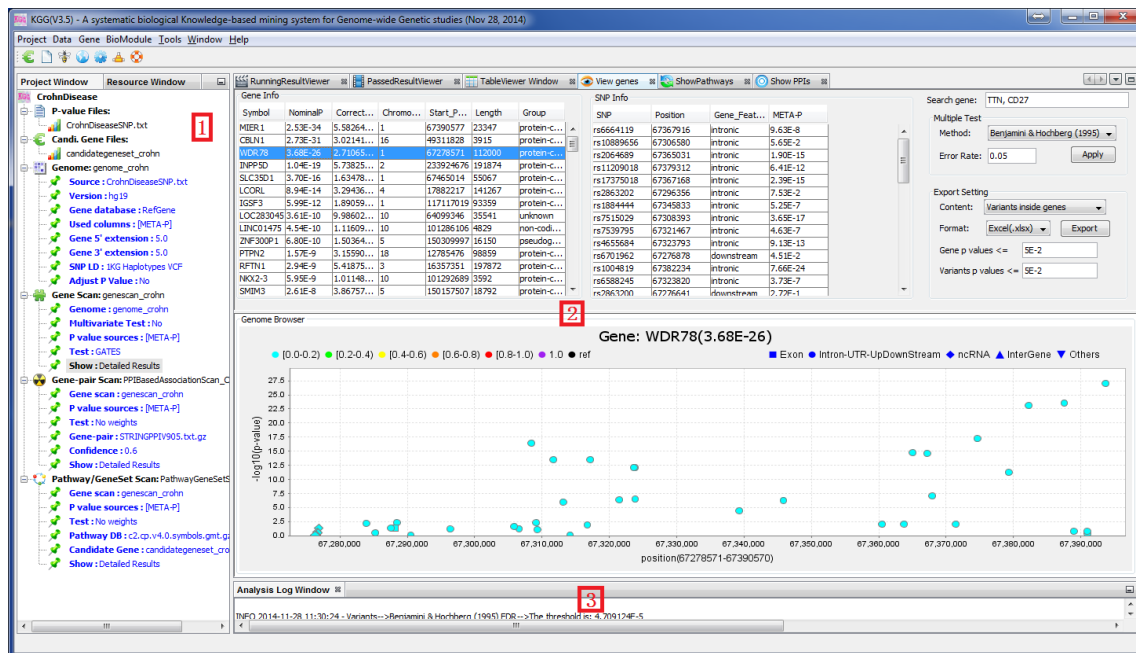Figure 3.1 shows a typical interface of KGG with an active project.

*Figure 3.1 A typical KGG interface*

Illustration:

**Frame 1:** tree-structured branches to manage input data and analysis results of a KGG project;

**Frame 2**: view of input data or output results;

**Frame 3**: running log of KGG analysis results;

The graphic dialogs of KGGs are self-explaining. Therefore, we will not elaborate the function of each buttons.

## 3.1 Project
➢ **Create project**: create a new KGG project.
➢ **Open project:** open an existing KGG project.
➢ **Close project:** close the current project.
➢ **Exit:** exit the KGG application.

## 3.2 Data
➢ **Load P value file:** import your association summary results (e.g., the plink output).
➢ **Define seed genes:** tell KGG the known causal genes of the disease you are studying.
➢ **Build analysis genome:** build an analysis genome in which KGG maps all SNPs to their gene features and calculates the r-square or genotypic correlation of SNPs within genes.

## 3.4 Gene
➢ **Gene-based association scan:** conduce the gene-based association scans.
➢ **View genes:** view and export gene-based association results.

## 3.5 Module
➢ **PPI-based association scan:** conduct PPI based association scan.
➢ **View PPIs:** view significant PPI pairs.
➢ **Pathway-based association scan:** conduct pathway based association scan.
➢ **View pathways:** view significant pathways.

## 3.5 Tools
➢ **Set system memory**: set the memory of KGG.
➢ **Power calculator**: SPS-a simulation tool for calculating power of set-based genetic association tests.

## 3.6 Window
➢ **AnalysisOutput: show the results when performing multiple tests.**
➢ **Project: depict the structure of the project.**
➢ **Resource: show the resource that KGG contains.**
➢ **ResultViewer: give the real-time results when performing the concrete analysis.**
➢ **RunningResultViewer: record the parameters using in each analysis.**
➢ **TableViewer: display the content of some files.**
➢ **Output: show all the IDE output.**

# 4. Input files

## 4.1 Input file 1 (GWAS results)
  KGG focuses on the downstream analysis of GWA studies, where statistical association p-values (or chi-square values) at SNPs have been generated by conventional statistical genetic methods (such as PLINK). Therefore, the association p-values are the major input of our KGG. KGG flexible supports a user-customized format for the association p-values.

   Once three columns of information, chromosome number and SNP IDs (or physical position) and p-values are available in a file, you can define the column order by yourselves on KGG. The input file can include more than one p-value column. The following is an example.

*Example input format (with rsID) of KGG:*

| CHR | SNP | P-value1 | P-value2 | P-value3 | … |
|-----|-----|----------|----------|----------|---|
| 4 | rs1513559 | 0.02301 | 0.8815 | 0.007688 | … |
| 4 | rs294755 | 0.4384 | 0.9575 | 0.006112 | … |
| 4 | rs835316 | 0.002688 | 0.007688 | 0.4893 | … |
| 4 | rs1841043 | 0.01115 | 0.006112 | 0.119 | … |
| 4 | rs11726946 | 0.005892 | 0.4893 | 0 | … |
| … | … | … | … | … | … |

*Example input format (with only position) of KGG:*

| CHR | SNPID | SNPPOS | P-value1 | P-value2 | P-value3 | … |
|-----|-------|--------|----------|----------|----------|---|
| 4 | Snp1 | 100001 | 0.02301 | 0.8815 | 0.007688 | … |
| 4 | Snp2 | 110011 | 0.4384 | 0.9575 | 0.006112 | … |

| 4 | Snp3 | 120001 | 0.002688 | 0.007688 | 0.4893 | ... |
| 4 | Snp4 | 130011 | 0.01115 | 0.006112 | 0.119 | ... |
| 4 | Snp5 | 140001 | 0.005892 | 0.4893 | 0 | ... |
| ... | ... | ... | ... | ... | ... | ... |

Moreover, a p-value column could include values of different models. KGG will recognize this format if you select the input format as "multiple tests per column" when building the analysis genome.

*Example a more complex input format of KGG:*

| CHR | SNP | P-value1 | Test-Mode | P-value2 | ... |
|-----|-----|----------|-----------|----------|-----|
| 4 | rs1513559 | 0.02301 | additive | 0.007688 | ... |
| 4 | rs1513559 | 0.4384 | recessive | 0.006112 | ... |
| 4 | rs1513559 | 0.002688 | dominant | 0.4893 | ... |
| 4 | rs1841043 | 0.01115 | additive | 0.119 | ... |
| 4 | rs1841043 | 0.005892 | recessive | 0 | ... |
| ... | ... | ... | ... | ... | ... |

### 4.2 Input file 2 (Candidate Gene list)

Candidate genes could be loaded one by one or imported from a TXT file. The input file has only one column without header, while one row contains one gene (symbol or ID).

## 5. Set-based association analysis tutorial

**Step 1**: create a new project, named 'CrohnDisease', and set the project path at C:\KGG\Tutorial (or other path defined by user).



*Figure 5.1 Create project*

**Step 2**: select the menu Data>Load P Value File and choose 'CrohnDiseaseSNP.txt' file which contains the whole-genome association p-values for Crohn disease at SNP-level. This dataset was downloaded from a public domain released by (Barrett, et al., 2008). It includes 7 columns, as SNP, CHR, POS, RISK, NONRISK, META-Z and META-P.

*Figure 5.2 Input GWAS original result file*

**Step 3**: import file 'CrohnCandidateGeneSet.txt' as input of candidate gene; define ATG16L1, CARD9, IBD5, IL23R, NOD2 and TNFSF15 as seed genes. Then, save it as candidategeneset_crohn.



*Figure 5.3 Input candidate gene set for crohn's disease*

**Step 4**: select META-P for building analysis genome; extend gene region to its flanking 5 kb region in both sides; and use LD SNP coefficients from 1000 Genome Project to adjust LD.
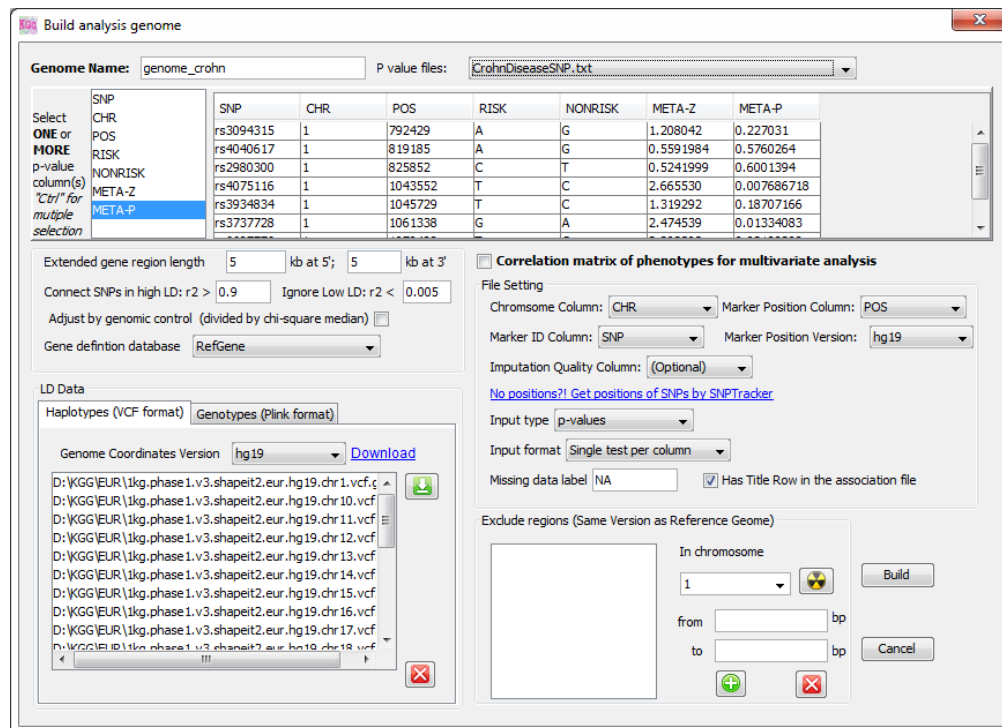
*Figure 5.4.1 Select META-P to build analysis genome and name the genome as genome_crohn*
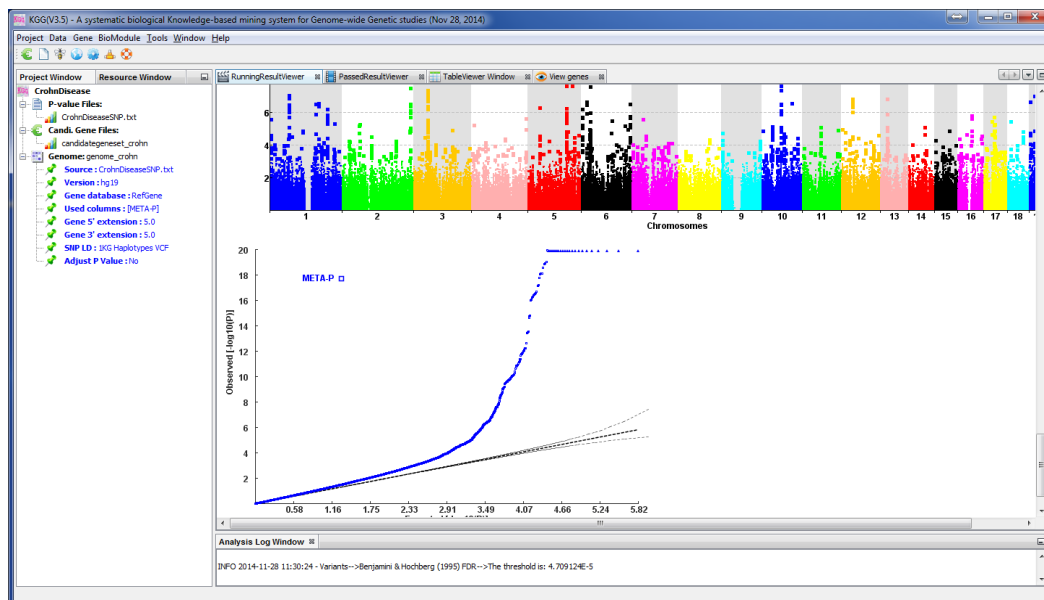


*Figure 5.4.2 The display after building analysis genome*

**Step 5**: do a gene-based association scan using SNP p-values integrated in the analysis genome named genome_crohn, select 'Extended Simes test(GATES, more powerful for a gene with one or a few independent causal variants' method. Set the parameters as Figure 5.5.1; and name the result as genescan_crohn. Remember that exported Manhattan plots and QQ plots will be shown in "Running Result Viewer Window" (Figure 5.5.2).
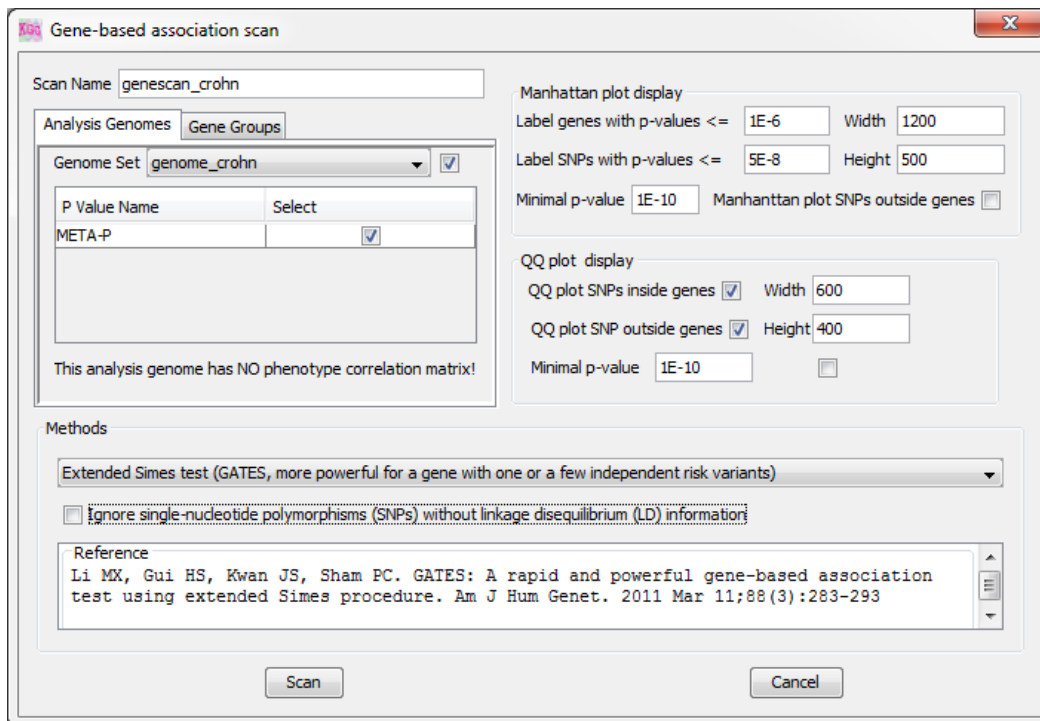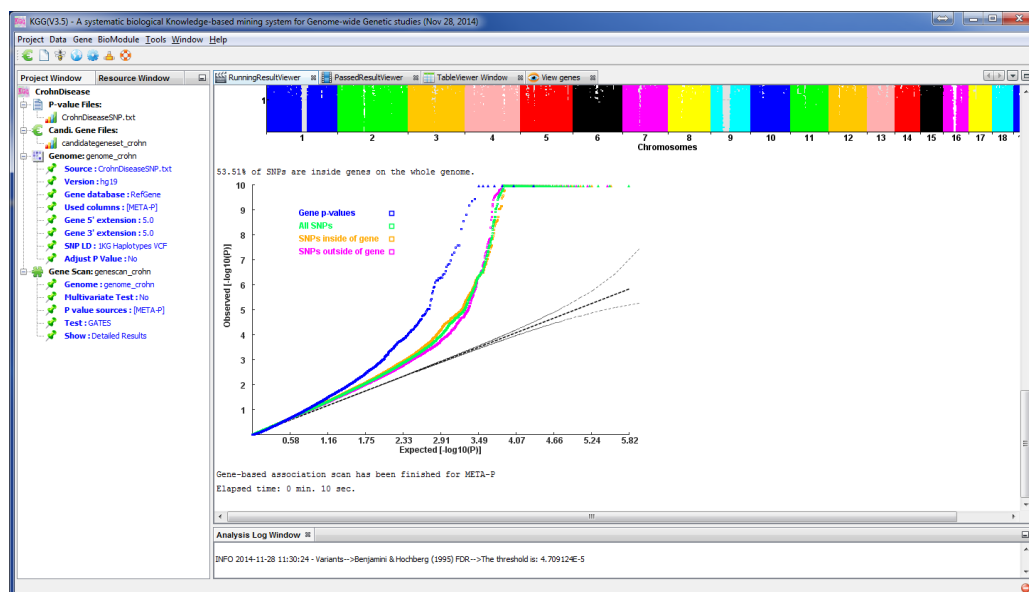
*Figure 5.5.1 Setting for gene-based scans*



*Figure 5.5.2 The display after gene-based scan*

**Step 6:** Click the "Show: Detailed Results" node under "Genome Scan" and a new tab "ShowGenes" will be created to provide you more information about the result (Figure 5.6). You can also export the results you want in this tab.
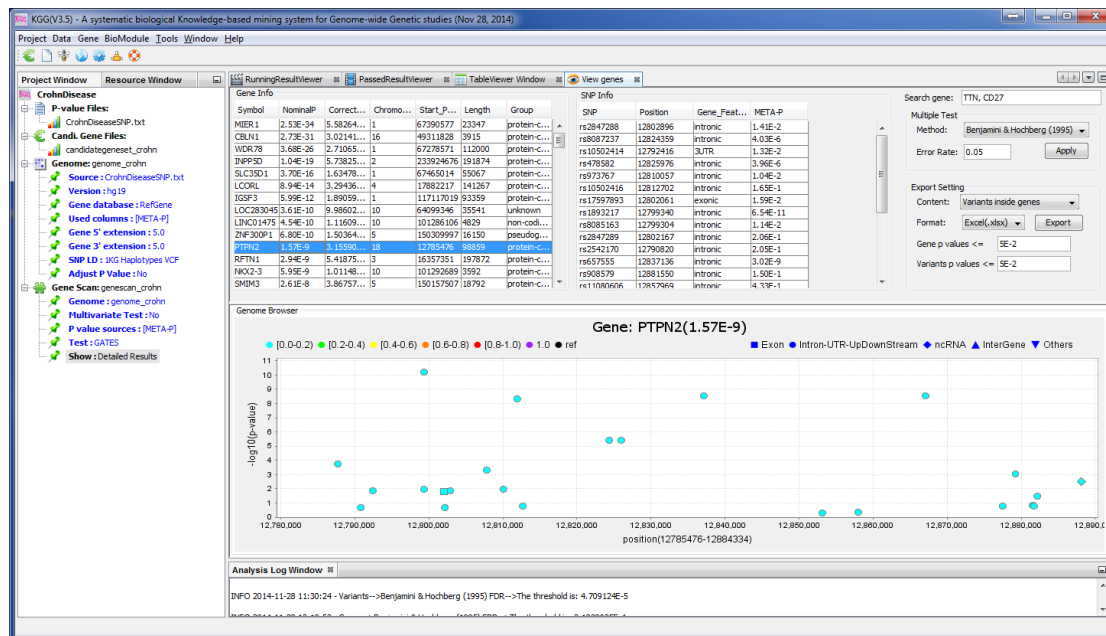
*Figure 5.6 Function of displaying the gene-based association scan result*

**Step 7**: perform pathway enrichment exploration both by gene p-values; settings as Figure 5.7.1 and the output as Figure 5.7.2.



*Figure 5.7.1 Pathway enrichment exploration by gene p-values*

*Figure 5.7.2 The display after pathway-based association scan*

**Step 8**: for more detailed information of the result, you can click the node "Show Detailed Result" (Figure 5.8). You can also change the multiple test methods and export the results you want in this tab.



*Figure 5.8 Function of displaying the results of pathway-based analysis*

**Step9**: search PPIs between significant genes. The significant genes can be picked up according to the gene p-values and SNP p-values; set as Figure 5.9.1; output as Figure 5.9.2

*Figure 5.9.1 PPI association scan by gene-based p-values*



*Figure 5.9.2 The display after running PPI-based association scan*

**Step 10:** Click the node "Show: Detailed Results" and you will get the graph of PPI network. You can also export the results you want in this tab.

*Figure 5.10 Function of displaying the results of PPI-based association scan*
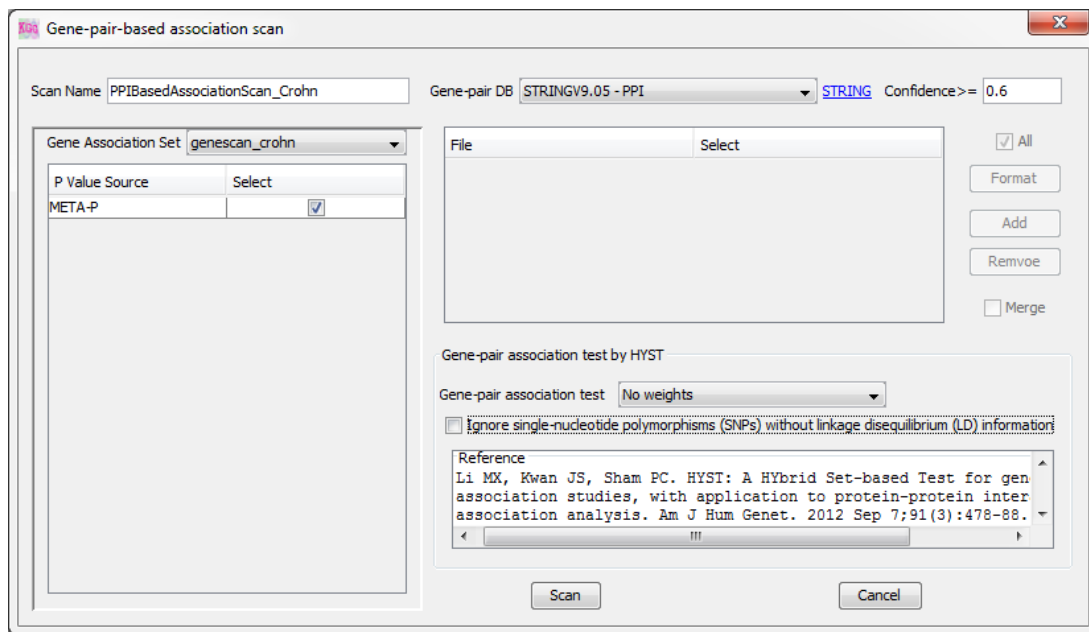
**Step 12:** View results of Crohn's Disease

➢ By text file or Excel file

  Open text or excel file for snp-based or gene-based analysis from local computer

➢ By Graphs

  Check QQ plots and Manhattan plots saved in htmlLog folder

➢ By KGG Interface

  Visualize pathway and PPI network output on KGG interface.

# 6. Power estimation of set-based tests by SPS

**Step 1:** Open the software and enter the main user interface on KGGV3.5 (Tools->Power Estimation). The interface is divided into two parts. The left one is used to set the basic parameters and the right one is to display the results.

*Figure 6.1 The Main interface of SPS.*

**Step 2:** Set the parameters of all variants. The number of SNPs, the minor allele frequency (MAF) and LD information should be set here. When these SNP markers are divided into several LD blocks, the markers within the same LD block have the same LD as each other, but the LD is set to 0 when the markers belong to different blocks. All of these markers and their LD pattern can be replicate to make up of a larger marker set. Some of these parameters can vary within a certain region, such as MAF and LD, so that the users can investigate how the power will be affected by changing the critical parameters conveniently. Moreover, these parameters can also read from the real data (Plink binary genotype files and vcf file). In this case, the LD information will be calculated from the input genotypes.



*Figure 6.2-1 Set parameters by users.*

*Figure 6.2-2 Set parameters by plink file.*



*Figure 6.2-3 Set parameters by vcf file.*

**Table to list parameters:**

| Parameter | Description |
|---|---|
| Total Variants | The total number of SNPs tested in a set |
| LD Block | The number of LD blocks. Variants in the same block are in LD and that in different blocks have no LD. |
| Repeat Region | The number of copies of SNPs. The SNP will be copied for several times to form a larger set and so does the LD pattern of the. |
| Minor Allele Frequency | The frequency of the least common allele occurs in the population. The MAF can increase from a initial value to a terminal value according to a step value that set from the GUI. |
| SNP Dependence | The relationship between SNPs. If the SNPs are dependent, the user should set the LD value (r), otherwise 0 is set as default. The LD information can also be read from the real data, where it will be calculated based on the allele frequency. |

| | |
|---|---|
| Linkage Disequilibrium (LD, r) | The r score used to represent LD information. The SNPs in the same block are dependent and keep the same r value, while SNPs in the different blocks are independent with each other and the r value is set as 0. The r value can also increase from an initial value to a final value by a step value. |
| Family File Map File BED File | The path of the Plink files. The valid file path can be input by the button on the right. If the three files have the same file prefix and are stored in the same directory, the other file paths will be filled automatically when one file is set. |
| Consider the first several SNPs | The number of SNP that input from the real data. The real data usually include large size of SNPs, which is unnecessary for our simulation. Hence, we just consider the first several SNPs as our study objects. |
| VCF File | The path of a VCF file. |

**Step 3:** Set parameters of risk variants.



*Figure 6.3 Set parameters about risk variants.*

**Table to list parameters:**

| Parameter | Description |
|---|---|
| Risk SNPs | The number of risk SNPs. This parameter can increase from a smaller to a larger value step by step. |
| Odds Ratio | The value used to quantify the association between risk SNPs and disease. This parameter can increase from a smaller to a larger value step by step. |
| Disease Prevalence | The proportion of a population found to suffer the disease. This will be used in the genetic model. |
| Genetic Model | The genetic model of risk loci. The additive model and multiplicative model are candidates in SPS. |
| Position of Risk Variants | The location information of risk variants within the total variants. The users can click the random button for automatic setting or set by themselves. |

**Step 4:** Set population and sample. The larger population size and number of case and control are recommended, because they make the result more accurate and stable, but it will take more time correspondingly. So the user should keep balance between them.



*Figure 6.4. Set population and sample.*

**Table to list parameters:**

| Parameter | Description |
|---|---|
| Population Size | The number of individuals in a population generated by simulation according to the certain genotype and phenotype. |
| Number of Case | The number of individuals that suffer the disease. |
| Number of Control | The number of individuals that do not suffer the disease. |

**Step 5:** Set simulation and meta-analysis parameters. A number of case-control samples will be randomly drawn with replacement from the population and are subject to calculate the p value of the set-based test. The number of p values that pass the threshold will be counted to calculate the power. In order to speed up the simulation process, the user can set several parallel threads, but more memory resource is needed.

The meta-analysis can be carried out at the variant level or set level. When at variant level, the p values of variants in different studies will be combined using Fisher's Combination Test and these meta-p values in a set will be treated by GATES, ScaChi and HYTS. Alternatively, at set level, the p value of variants in a set should be conducted by GATES, ScaChi and HYTS, and then the set-based p values in different studies are aggregated. SPS can also mimic locus heterogeneity by randomizing risk loci of each study in meta-analysis.



*Figure 6.5-1 Set simulation without meta-analysis.*



*Figure 6.5-2 Set simulation with meta-analysis.*

**Table to list parameters:**

| Parameter | Description |
|---|---|
| Sampling Times | The number samples randomly drawn from the case and control group. For each time, a case-control study is achieved. |
| P Value Threshold | The threshold of type I error that used in the case-control study. For SNP-based test, the bonferroni correction is conducted as default. |
| Parallel Running Number | The number of threads that running concurrently. The multiple threads mechanism is used here to speed up the running of program. However, this may cost a large volume of memory. |

| Meta-analysis | Whether to perform meta-analysis. If performed, the users should choose the meta-analysis at variants level or at set level. |
|---|---|
| Number of Studies | The number of studies considered in the meta-analysis. |
| Randomize risk loci of each study | Whether to consider the genetic heterogeneity. If considered, the position of risk loci of each study will be set randomly to mimic the heterogeneity. |

**STEP 6:** Run the program. Click the Start button and run the program. The user can check the results from tables in the right part immediately. The progress bar can also provide the running information in a real time. If the user wants to stop the running program, just click the "stop" button.



*Figure 6.6 Run the program.*

**STEP 7:** Save the result. The user can review the power from two tables at the SNP level and set level. A line chart is draw to show the variation of power within different odd ratios with given the MAF and LD information. The user can also change the MAF and LD values to update the chart. The users can right-click on the tables and save the results as excel file or txt file. The chart is can also be save by right-click.



*Figure 6.7-1. The output of SPS.*

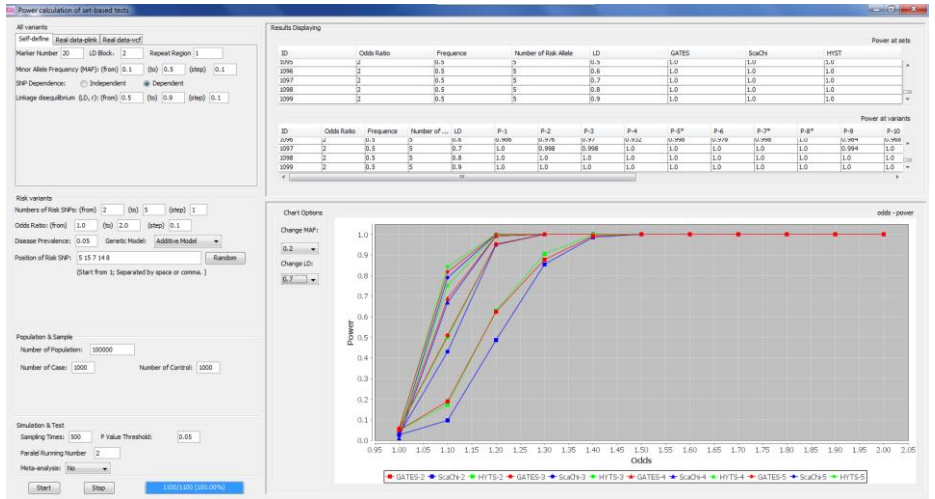| ID | Odds Ratio | Frequence | er of Risk A: | LD | GATES | ScaChi | HYST |
|----|-----------|-----------|---------------|-----|-------|--------|------|
| 0 | 1 | 0.1 | 1 | 0.5 | 0.055 | 0.025 | 0.048 |
| 1 | 1 | 0.1 | 1 | 0.6 | 0.04 | 0.016 | 0.033 |
| 2 | 1 | 0.1 | 2 | 0.5 | 0.042 | 0.025 | 0.041 |
| 3 | 1 | 0.1 | 2 | 0.6 | 0.047 | 0.021 | 0.047 |
| 4 | 1 | 0.1 | 3 | 0.5 | 0.049 | 0.017 | 0.048 |
| 5 | 1 | 0.1 | 3 | 0.6 | 0.049 | 0.029 | 0.044 |
| 6 | 1 | 0.1 | 4 | 0.5 | 0.044 | 0.028 | 0.049 |
| 7 | 1 | 0.1 | 4 | 0.6 | 0.062 | 0.027 | 0.055 |

*Figure 6.7-2 The saved table of set-based power.*

| ID | Odds Ratio | Frequence | er of Risk A: | LD | P-1 | P-2 | P-3* | P-4 | P-5 | P-6 | P-7* |
|----|-----------|-----------|---------------|-----|-------|-------|-------|-------|-------|-------|-------|
| 0 | 1 | 0.1 | 1 | 0.5 | 0.002 | 0.004 | 0.005 | 0.001 | 0.001 | 0.002 | 0.001 |
| 1 | 1 | 0.1 | 1 | 0.6 | 0.003 | 0.001 | 0 | 0.001 | 0 | 0.003 | 0.001 |
| 2 | 1 | 0.1 | 2 | 0.5 | 0.003 | 0.003 | 0.001 | 0.002 | 0.001 | 0.001 | 0.002 |
| 3 | 1 | 0.1 | 2 | 0.6 | 0.003 | 0.006 | 0.005 | 0.002 | 0.002 | 0.001 | 0.004 |
| 4 | 1 | 0.1 | 3 | 0.5 | 0.004 | 0.002 | 0.001 | 0.003 | 0.003 | 0 | 0.001 |
| 5 | 1 | 0.1 | 3 | 0.6 | 0.003 | 0.002 | 0.003 | 0.001 | 0.002 | 0.001 | 0.004 |
| 6 | 1 | 0.1 | 4 | 0.5 | 0.003 | 0.004 | 0.003 | 0.004 | 0.003 | 0.004 | 0.004 |
| 7 | 1 | 0.1 | 4 | 0.6 | 0.001 | 0.002 | 0.001 | 0.003 | 0.004 | 0.002 | 0.001 |
| 8 | 1 | 0.2 | 1 | 0.5 | 0.004 | 0.004 | 0.005 | 0.002 | 0.004 | 0.002 | 0.005 |
| 9 | 1 | 0.2 | 1 | 0.6 | 0.008 | 0.004 | 0.003 | 0.001 | 0.003 | 0 | 0.006 |

*Figure 6.7-3 The saved table of variant-based power.*

# 7. Update from KGG 3.0 to KGG 3.5

Much progress was made from KGG 3.0 to KGG 3.5, mainly including:
1)  Multivariate gene-based association analysis;
2)  Direct link to multiple bioinformatics annotation databases;
3)  Simplified operation; better plotting function;
4)  Integrate regulatory information to prioritize risk genes (under development).
5)  SPS plug-in is included.