

# CGI: Java Software for Mapping and Visualizing Data from Array-based Comparative Genomic Hybridization and Expression Profiling

Joyce Xiuweu-Xu Gu<sup>1</sup>, Michael Yang Wei<sup>1</sup>, Pulivarthi H. Rao<sup>1</sup>, Ching C. Lau<sup>1</sup>, Sanjiv Behl<sup>2</sup> and Tsz-Kwong Man<sup>1</sup>

<sup>1</sup>Department of Pediatrics, Baylor College of Medicine and Texas Children's Cancer Center, Houston, Texas. <sup>2</sup>Department of Computer Science, University of Houston, Victoria, Texas.

**Abstract:** With the increasing application of various genomic technologies in biomedical research, there is a need to integrate these data to correlate candidate genes/regions that are identified by different genomic platforms. Although there are tools that can analyze data from individual platforms, essential software for integration of genomic data is still lacking. Here, we present a novel Java-based program called CGI (Cytogenetics-Genomics Integrator) that matches the BAC clones from array-based comparative genomic hybridization (aCGH) to genes from RNA expression profiling datasets. The matching is computed via a fast, backend MySQL database containing UCSC Genome Browser annotations. This program also provides an easy-to-use graphical user interface for visualizing and summarizing the correlation of DNA copy number changes and RNA expression patterns from a set of experiments. In addition, CGI uses a Java applet to display the copy number values of a specific BAC clone in aCGH experiments side by side with the expression levels of genes that are mapped back to that BAC clone from the microarray experiments. The CGI program is built on top of extensible, reusable graphic components specifically designed for biologists. It is cross-platform compatible and the source code is freely available under the General Public License.

**Keywords:** aCGH, expression profiling, visualization, correlation, and data integration

## Introduction

With the advent of genomic technologies, DNA and RNA-based microarrays are becoming more accessible to biomedical researchers. One of the common DNA platforms is array-based Comparative Genomic Hybridization (aCGH), which can identify DNA copy number aberrations in the genome (Pinkel, 1998; Man, 2004). There are many software tools that have been developed to analyze aCGH data (Jong 2004; Margolin, 2005; Chen, 2005; Cheung, 2005; Price, 2005; Kim, 2005) and expression microarray data (Sykacek, 2005; Shamir, 2005; Saraiya, 2005; Li, 2001; Vaquerizas, 2005; Bumm, 2002; Saeed, 2003); however, no tool is currently available for the biologist to integrate these two types of data. One of the main challenges is that once the significant BAC clones or genes are identified, it is very difficult to correlate the DNA copy number and RNA expression results. This is because the significant genes may not lie within the corresponding BAC clones even though they are located in the same chromosomal region. Therefore, a more precise method of matching is needed in order to properly correlate these two types of data.

A typical way to perform the matching is to manually search the UCSC Genome Browser (<http://genome.ucsc.edu/>) to make sure the significant genes lie within the significant BAC clones. However, this type of manual search is very laborious and error prone if the numbers of BAC clones and genes are large. Thus, it is important to develop a user friendly and flexible tool that can match, correlate and display the aCGH and expression profiling data. Since it is common to identify hundreds to thousands of significant genes by either expression profiling or aCGH experiments, our program can further assist researchers to select genes that are found to be significant by both types of experiments, or genes that may not be identified by using either type of technique alone.

To address this issue, we developed a Java-based, stand-alone program that uses MySQL database (<http://www.mysql.com>) as a backend to store the BAC clones and gene information downloaded from UCSC database. This information is used to match the user-provided BAC clones in aCGH experiments

**Correspondence:** Tsz-Kwong Man, Texas Children's Cancer Center, 6621 Fannin St. MC 3-3320, Houston, TX 77030-2399, Tel: (832) 824-4682; Email: [ctman@txccc.org](mailto:ctman@txccc.org)



Copyright in this article, its metadata, and any supplementary data is held by its author or authors. It is published under the Creative Commons Attribution By licence. For further information go to: <http://creativecommons.org/licenses/by/3.0/>.

and genes in expression profiling experiments. After that, the correlation coefficients and p-values of the matched BAC clone-gene pairs will then be computed and displayed in various formats for data visualization and comparison.

## Software Designs

The CGI software is based on an object-oriented framework designed to conduct searches for features/genes in RNA expression-profiling experiments that mapped back to corresponding BAC clones in aCGH experiments. The program combines bioinformatic data matching from databases with simple correlation analysis. The software is organized into three functional modules (Data, Annotation, and Correlation). The Data module contains DNA copy number and RNA expression data and links them with the Annotation module by interacting with the MySQL database that holds a variety of different types of genomic information including chromosomal localization, Unigene ID, and gene annotation data. Information in the database is used to match the BAC clones and the genes provided by the users. The Correlation module calculates the Pearson correlation coefficients and p-values between the DNA copy numbers and expression values of matched BAC clone-gene

pairs in different experiments. It also displays DNA copy numbers of a specific BAC clone in different aCGH experiments and the associated gene expression values in microarray experiments for easy data visualization and comparison.

## Data Importing

A simple graphical user interface (GUI) prompts users to enter user name, password, database name, and the locations of the aCGH and RNA expression-profiling files (Fig. 1). The aCGH file contains FISH-mapped BAC Clone IDs, cyto-bands, and normalized log ratios representing DNA copy numbers from aCGH experiments. The RNA expression-profiling file contains Unigene IDs, gene symbols, and log-ratios (dual channel arrays) or log intensities (Affymetrix or oligo-based arrays) of gene expressions in a set of experiments involving identical cases as in the aCGH experiments.

## Data Querying and Mining

The program offers two ways to query the data. First, BAC Clone IDs in an aCGH input file are used to query the MySQL database, which stores data downloaded from the UCSC database at

**Figure 1.** Graphic user interface of CGI for data importing and analysis.

URL: <http://genome.ucsc.edu/cgi-bin/hgTables—fishClone> and `uniGene_2` tables. The two tables are first downloaded by the user and imported to the MySQL databases as described in the installation manual (see supplemental information). The Unigene IDs of the genes that reside in each BAC clone in the aCGH input file are retrieved based on chromosome number and their physical locations by SQL commands. Secondly, these Unigene IDs are used to match with the features/genes provided in a RNA expression-profiling input file, so that the matched BAC clone-gene pairs will be identified. The DNA copy numbers and gene expression values of the matched BAC clone-gene pairs will then be extracted from the input files and their Pearson correlation coefficients and p-values are computed by an internal correlation functions. Finally, the correlation coefficients and p-values of the BAC clone-gene pairs will be tabulated together with their BAC Clone IDs, cytobands, Unigene IDs, and gene symbols provided by the input files. If there are multiple genes within a BAC clone, the program will replicate the DNA copy number data of that BAC clone and correlate with the expression data of each of the other genes that are mapped to the BAC clone.

## Data Visualization

CGI uses a correlation table to display a global overview of BAC Clone ID, Cytoband, their corresponding Unigene IDs, gene symbols and Pearson correlation coefficients and p-values of the matched BAC clones and features/genes. The table view is very flexible and the data in the table can be sorted dynamically in an ascending or descending order based on the correlation coefficients, BAC Clone ID, cytoband location, Unigene ID, etc (Fig. 2). It can also change the order of the columns to display different views according to user's preference. Besides the table view, users can also visualize in detail the DNA copy number of a specific BAC Clone and the expression of its associated genes by entering the BAC Clone ID into the text box provided in the GUI (Fig. 1). The CGI program will display two graphic windows if the input BAC Clone ID matches one or more Clone IDs in the correlation table. One window displays three line graphs representing the DNA copy number changes of the queried BAC Clone in aCGH experiments and the expression values of its associated genes in RNA expression-profiling experiments (Fig. 3). The second window displays the DNA copy number data and RNA expression data

Clone ID	Cyto	Unigene ID	Gene Symbol	Correlation	Significance (Z)
RP11-79021	22q12.1	Hs_7370	PITFIB9	0.820927	0.005393
RP3-37420	6q21.1	Hs_237066	ZNF117	0.845792	0.009597
RP11-79014	20p11.21	Hs_7218	N/A	0.843861	0.00959
RP11-383815	19p13.3	Hs_128425	C11orf29	0.839753	0.010219
RP11-91891	17q24.3	Hs_300963	RPL3B	0.839102	0.026819
RP11-91E11	6p21.2	Hs_202331	N/A	0.838573	0.038327
RP11-8095	6p21	Hs_194726	BA3A	0.837596	0.040421
RP11-80106	11q25	Hs_178488	HSPC063	0.828429	0.045904
RP11-660M5	4p15.1	Hs_253305	SLALP	0.810571	0.051903
RP11-88L5	9q22.32	Hs_305059	N/A	0.808906	0.052321
RP11-90822	6q25	Hs_103130	N/A	0.805639	0.05459
RP1-12593	19p38.11	Hs_209983	BTMNI1	0.877339	0.071967
RP11-864	9p21.23	Hs_301404	FBM3	0.875225	0.072416
RP11-7904	17p11.2	Hs_324448	N/A	0.870245	0.076655
RP1-154014	6p15	Hs_349095	N/A	0.86592	0.087108
RP11-7909	22q12.2	Hs_408277	BF3A1	0.84593	0.095959
RP11-8107	17q21.31	Hs_278908	DCO5	0.817128	0.121891
RP11-64L12	19p13.3	Hs_444725	N/A	0.815568	0.123438
RP11-383815	19p13.3	Hs_282177	PPP5K1C	0.808467	0.142467
RP11-80A1	10q23.33	Hs_530374	N/A	0.3912	0.148312
RP11-79F21	6p12	Hs_348921	PHF3	0.38116	0.160991
RP11-80P3	12q24.13	Hs_2714	N/A	0.378476	0.164208
RP11-80P19	3q21.1	Hs_370800	N/A	0.366587	0.194556
RP11-46C8	19q13.2	Hs_22049	N/A	0.359348	0.188352
RP11-79021	22q12.1	Hs_816221	N/A	0.340027	0.194556
RP11-91L1	20q13.2	Hs_32135	N/A	0.353902	0.195625
RP11-81F9	6p21.31	Hs_75243	BRD2	0.348933	0.203164
RP11-81F9	14q11.2	Hs_610026	SC1L2	0.347786	0.204117
RP11-81D7	17q23.3	Hs_279808	DCO5	0.345751	0.208844
RP11-9115	12q13.2	Hs_154057	MMP19	0.343846	0.209524
RP11-80C7	2q11.2	Hs_30992	POLR1A	0.335377	0.221705
RP11-89N1	19q23	Hs_131915	KIAA0883	0.331974	0.226723
RP11-80B9	15q25.3	Hs_33954	HOMER2	0.32769	0.23314
RP3-355C18	22q13.33	Hs_74516	N/A	0.325793	0.236917
RP11-33301	5q24.3	Hs_155418	N/A	0.320538	0.244102
RP11-89C1	2q12	Hs_107845	FLJ10996	0.319311	0.249014
RP11-21616	19p13.2	Hs_171644	PVR	0.314022	0.254261
RP11-563021	19q13.33	Hs_528342	N/A	0.280096	0.283922
RP11-79A12	21q22.11	Hs_331059	NCF1	0.28094	0.292961
RP11-88J10	15q15.3	Hs_259533	CHMT1	0.297844	0.296913
RP11-46C8	19q13.2	Hs_7486	ETH1	0.277544	0.316573
RP11-8106	22q11.21	Hs_300925	N/A	0.276929	0.31769
RP11-80L1	10q25.33	Hs_277497	N/A	0.268663	0.325862
RP11-89D20	10p12.1	Hs_414357	N/A	0.265221	0.339404
RP11-89D23	7q22.1	Hs_349898	N/A	0.264938	0.359341
RP11-89F20	3q14	Hs_77263	SENTAD2	0.252842	0.361678
RP3-329A5	6p21.31	Hs_28522	N/A	0.251149	0.365579
RP11-1293	19p13.11	Hs_249746	HNRPA8	0.246769	0.375271
RP11-90L9	2q26.2	Hs_75426	CC120	0.239896	0.389842
RP11-89P19	8q23.3	Hs_188536	N/A	0.23524	0.396891
RP11-64L12	19p13.3	Hs_18079	PIGO	0.233403	0.402462
RP11-91H15	14q21.1	Hs_523505	N/A	0.232965	0.422324
RP11-01F13	6p21.1	Hs_15250	PEC1	0.219719	0.431383
RP11-89A20	7q11.22	Hs_583	NCF1	0.217048	0.427147
RP11-79A15	15q21.2	Hs_652664	N/A	0.212911	0.448361
RP11-89N1	19q23	Hs_32796	N/A	0.206845	0.459504
RP11-81M8	12q23.2	Hs_235404	PAH	0.202056	0.474012
RP11-79K1	12q13.11	Hs_256864	N/A	0.198162	0.483505
RP11-79J23	6p21.31	Hs_83126	N/A	0.191267	0.494699
RP11-80P3	12q24.13	Hs_5149	FLN2B	0.18893	0.505987
RP11-79L13	2q26.3	Hs_309539	N/A	0.188236	0.508236
RP11-88J6	4q13.3	Hs_164021	CXCL6	0.185907	0.507093
RP11-79E17	1q24.2	Hs_381155	N/A	0.179523	0.524399
RP11-89E17	9q21.31	Hs_84162	N/A	0.17926	0.525981
RP11-01F13	6p21.1	Hs_188891	N/A	0.175238	0.521183
RP11-89D14	9q22.2	Hs_83758	CK2E	0.171577	0.522447
RP3-37891	22q13.2	Hs_35347	N/A	0.162985	0.581903
RP11-64L12	19p13.3	Hs_137572	RHBDL1	0.162549	0.562478
RP11-81H5	4p15.33	Hs_389899	RAO2B	0.159379	0.570464
RP11-90M15	6p12.2	Hs_12463	N/A	0.152878	0.584887
RP11-90M15	13q12.2	Hs_79877	N/A	0.151264	0.590487
RP11-80P5	6p12	Hs_425211	LSM1	0.148684	0.594418

Figure 2. Table view to display the correlations of the matched BAC clones in aCGH and genes in expression microarray experiments.

as separate bar graphs for better visualization of the individual experiments if the number of matched genes is high (Fig. 4). This function provides a graphical visualization of the correlation between a BAC clone and its matching genes/features.

## Application

To test this program, we have analyzed a previously published dataset that contains data from both aCGH arrays (Man, 2004) and cDNA microarrays (Man, 2005) of a set of pediatric osteosarcoma patients. We found several genes with RNA expressions correlating with the DNA copy numbers in the corresponding BAC clones ( $r > 0.5$ ,  $n = 15$ ,  $p < 0.05$ , Fig. 2). One of the highly correlated genes (ZNF187) is mapped back to the BAC clone RP5-874C20, which is one of the most frequently amplified regions (6p21.1) in osteosarcoma (Man, 2004). ZNF187 or SRE-ZBP is induced by serum response and may regulate oncogene c-fos by binding to its serum response element (Attar, 1992). We have validated matching and correlation results of CGI by manually searching the UCSC genome browser to confirm the match between BAC Clone ID and

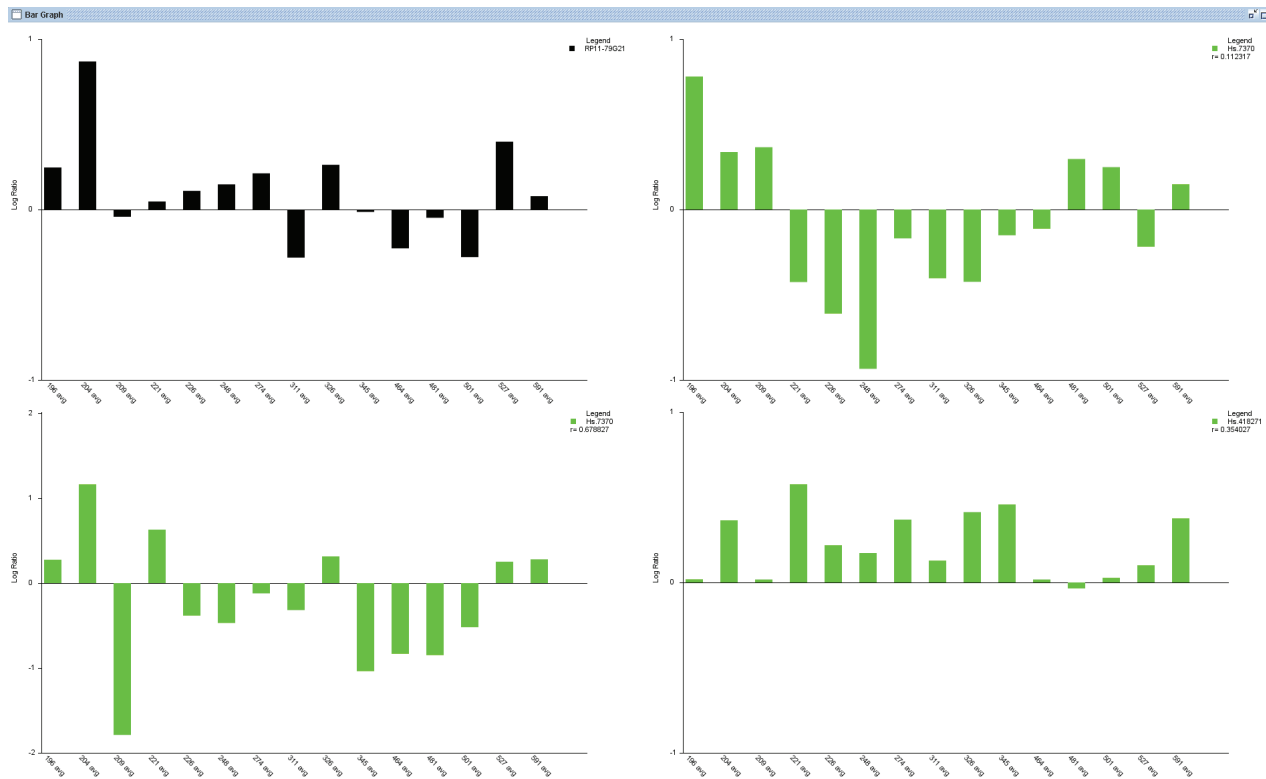
Unigene ID, and recalculated their correlation coefficients using an independent method.

## Discussion

We have developed the CGI program, which provides a simple yet powerful tool for matching, correlating, and visualizing aCGH and gene expression-profiling results simultaneously in multiple experiments. This tool is useful because it correlates the results from DNA profiling with those from RNA expression-profiling experiments in order to identify genes that are important at both DNA and RNA levels. The genes that are significantly altered in both sets of experiments add more confidence to the biological significance of these genes and therefore warrant further investigation. It also alleviates the need for manual matching between BAC clones on the aCGH arrays and the features in gene expression arrays using public databases. For data analysis, it provides a visualization tool and correlation calculations with an interactive and flexible interface. We have also implemented error detection routines to handle the database connection, e.g. user needs to enter



**Figure 3.** Line graph viewer. The top panel is the copy number changes of a BAC clone in a set of aCGH experiments. The bottom panel is the gene expression values of three corresponding genes that matched the BAC clone in expression microarray experiments using the same experimental cases.



**Figure 4.** Bar graphic viewer. The black bar graph shows DNA copy numbers (on the Y-axis) of the BAC clone from 15 different experiments (numbered on X-axis). The relative gene expression values (on the Y-axis) of three corresponding genes mapped to the BAC in expression array experiments were displayed in three separate green bar graphs for clearer visualization.

username, password, and database name for secure connection. The number of experiments in the input files are also checked to ensure comparability of the data. This software was implemented in an object-oriented language, Java, to ensure portability across different operating systems. It is a stand-alone program, which is designed for users to install and run on their own local machine. Therefore, unlike other web-based analytical tools, the users do not need any server support, and are not affected by Internet traffic, server-side problems and downtime. Instead of using flat file data storage, the CGI program also provides fast data access and transfer from MySQL database, which is freely available via the web site (<http://www.MySQL.com>).

Different from some analytical tools, such as BioConductor ([www.bioconductor.org](http://www.bioconductor.org)), which uses command line interface, the CGI program uses an easy-to-use and intuitive graphical interface for bench biologists to perform the analysis without any prior computational background. A detailed description on how to install the databases and program is also provided in the supplementary information. The software framework that we

employ supports the development of more sophisticated visualization and analytical functions in the future through its open API for Java-based plugins. The program is coded in Java reusable object classes, thus promoting a rapid development of future program extensions.

Two similar efforts of comparing aCGH and expression-profiling have been published recently. Kingsley et al have recently developed a web-based system, Magellan, which explores the quantitative relationship between aCGH and mRNA expression data (Kingsley, 2006). Magellan computes the relationship of aCGH and expression based on common annotation values between the two sets of experiments. Shankar et al has also developed a program mainly to visualize aCGH and expression data (Shankar, 2006). In contrast to these other two programs, the CGI program is standalone program, which does not require Internet connection and is not affected by the server-side problem. In addition to visualizing the data, the main strength of CGI is to provide an easy-to-use interface for fast matching and correlation of these two types of genomic data using a relational database. Once these candidate genes are identified, they can be

subjected to additional analyses using other existing analytical tools. Since the software is developed in the object-oriented language Java, it can interact with other programs currently available for aCGH and microarray analysis, such as the BioConductor packages. It is straightforward to include other computational algorithms to extend the analytical capability of the program. The modular design of this program also adds flexibility and extensibility for the development of more functions and plug-ins in the future. In summary, we have developed an easy-to-use program CGI to map, correlate, and visualize aCGH and expression profiling data.

## Acknowledgements

We would like to thank Jaya Visvanthan, and Jianhe Shen for the preparation of the aCGH and microarray data used in this study. We also thank Alexander Yu, Wei-chun Hsu, and Richard Lowry for their help in programming, and Carolyn Pena for her assistance in manuscript preparation. The study is supported by grants from NIH CA81465, the Robert J. Kleberg, Jr. and Helen C. Kleberg Foundation, the Gillson Longenbaugh Foundation, and the Cancer Fighters in Houston (CCL), as well as the Sarcoma Foundation of America and Fleming and Davenport Award (TKM).

## References

Attar, R.M. and Gilman, M.Z. 1992. Expression cloning of a novel zinc finger protein that binds to the c-fos serum response element. *Mol. Cell. Biol.*, 12:2432–43.

Bumm, K and Cheng, M. 2002. CGO: utilizing and integrating gene expression microarray data in clinical research and data management. *Bioinformatics*, 18:327–8.

Chen, W., Erdogan, F., Ropers, H.H. et al. 2005. CGHPRO—a comprehensive data analysis tool for array CGH. *BMC Bioinformatics*, 6:85.

Cheung, S.W., Shaw, C.A., Yu, W. et al. 2005. Development and validation of a CGH microarray for clinical cytogenetic diagnosis. *Genet. Med.*, 7:422–32.

Jong, K., Marchiori, E., Meijer, G. et al. 2004. Breakpoint identification and smoothing of array comparative genomic hybridization data. *Bioinformatics*, 20:3636–7.

Kim, S.Y., Nam, S.W., Lee, S.H. et al. 2005. ArrayCGH: a web application for analysis and visualization of array-CGH data. *Bioinformatics*, 21:2554–5.

Kingsley, C.B., Kuo, W.L., Polikoff, D. et al. 2006. Magellan: A Web Based System for the Integrated Analysis of Heterogeneous Biological Data and Annotations; Application to DNA Copy Number and Expression Data in Ovarian Cancer. *Cancer Informatics*, 1:10–21.

Li, C. and Wong, W.H. 2001. Model-based analysis of oligonucleotide arrays: Expression index computation and outlier detection. *Proc. Natl. Acad. Sci.*, 98:31–36.

Man, T.K., Lu, X.Y., Kim, J. et al. 2004. Genome-wide array comparative genomic hybridization analysis reveals distinct amplifications in osteosarcoma. *BMC Cancer*, 4:45.

Man, T.K., Chintagumpala, M., Visvanathan, J. et al. 2005. Expression profiles of osteosarcoma that can predict response to chemotherapy. *Cancer Research*, 65:8142–50.

Margolin, A.A., Greshock, J., Naylor, T.L. et al. 2005. CGHAnalyzer: a stand-alone software package for cancer genome analysis using array-based DNA copy number data. *Bioinformatics*, 21:3308–11.

Pinkel, D., Segev, R., Sudar, D. et al. 1998. High resolution analysis of DNA copy number variation using comparative genomic hybridization to microarrays. *Nature Genet.*, 20:207–11.

Price, T.S., Regan, R., Mott, R. et al. 2005. SW-ARRAY: a dynamic programming solution for the identification of copy-number changes in genomic DNA using array comparative genome hybridization data. *Nucleic Acids Res.*, 33:3455–64.

Saeed, A.I., Sharov, V., White, J. et al. 2003. TM4: a free, open-source system for microarray data management and analysis. *Biotechniques*, 34:374–8.

Saraiya, P., North, C. and Duca, K. 2005. An insight-based methodology for evaluating bioinformatics visualizations. *IEEE Trans. Vis. Comput. Graph.*, 11:443–56.

Shamir, R., Maron-Katz, A., Tanay, A. et al. 2005. EXPANDER—an integrative program suite for microarray data analysis. *BMC Bioinformatics*, 6:232.

Shankar, G., Rossi, M.R., McQuaid, D.E. et al. 2006. aCGHViewer: A Generic Visualization Tool For aCGH data. *Cancer Informatics*, 2:36–43.

Sykacek, P., Furlong, R.A. and Micklem, G. 2005. A friendly statistics package for microarray analysis. *Bioinformatics*, 21:4069–70.

Vaquerez, J.M., Conde, L., Yankilevich, P. et al. 2005. GEPAS, an experiment-oriented pipeline for the analysis of microarray gene expression data. *Nucleic Acids Res.*, 33:W616–W620.