

Profiling Gene Expression Using Onto-Express

Purvesh Khatri,¹ Sorin Draghici,^{1,2} G. Charles Ostermeier,³ and Stephen A. Krawetz^{2,3,*}

¹Department of Computer Science, ²Institute for Scientific Computing, and ³Department of Obstetrics and Gynecology, Molecular Medicine and Genetics, Wayne State University School of Medicine, Detroit, MI 48201, USA

*To whom correspondence and reprint requests should be addressed. Fax: (313) 577-8554. E-mail: steve@compbio.med.wayne.edu.

Gene expression profiles obtained through microarray or data mining analyses often exist as vast data strings. To interpret the biology of these genetic profiles, investigators must analyze this data in the context of other information such as the biological, biochemical, or molecular function of the translated proteins. This is particularly challenging for a human analyst because large quantities of less than relevant data often bury such information. To address this need we implemented an automated routine, called Onto-Express (<http://vortex.cs.wayne.edu:8080>), to systematically translate genetic fingerprints into functional profiles. Using strings of accession or cluster identification numbers, Onto-Express searches the public databases and returns tables that correlate expression profiles with the cytogenetic locations, biochemical and molecular functions, biological processes, cellular components, and cellular roles of the translated proteins. The profiles created by Onto-Express fundamentally increase the value of gene expression analyses by facilitating the translation of quantitative value sets to records that contain biological implications.

INTRODUCTION

The development of microarray technologies permits researchers to monitor mRNA transcript levels for thousands of genes in a single experiment [1-5]. The sheer power and potential for discovery using this whole transcriptome analysis [6,7] move well beyond traditional strategies that rely on the examination of a single gene or the assessment of a group of genes one at a time. The resulting patterns, that is expression profiles, depict subsets of transcripts that qualitatively reflect gene activity at a given instant or in response to a particular state or condition [8,9]. These profiles represent a genetic "fingerprint" that characterizes the cell or tissue being studied (G.C.O. *et al.*, manuscript submitted) [10] and provide a foundation from which to begin an investigation of the underlying biology. Biological functionality must be associated with each element of the genetic fingerprint to appreciate the intricacies of how specific tissues or cells function in a normal state or a disease state, or in response to a given treatment. A description of functionality should include the structural, regulatory, or enzymatic roles of the corresponding proteins; the biological "objectives" to which the proteins contribute; the places within a cell where the proteins are active; and the major biological processes in which these proteins play a role as well as their biochemical activities.

Databases are being constructed to store expression data [11] and algorithms are being implemented to support access to this information [12,13]. Even though these efforts are underway, the interoperability required to convert such data

into functional biological profiles is lacking. This is a daunting task when one considers that genetic fingerprints can contain tens of thousands of records (G.C.O. *et al.*, manuscript submitted) [5]. To this end, we developed the Java-based program called Onto-Express (<http://vortex.cs.wayne.edu:8080>). This algorithm uses UniGene accession or cluster identification numbers generated from microarray investigations to search the public databases to return the functional profiles for each genetic fingerprint. Here we demonstrate the utility of Onto-Express.

RESULTS AND DISCUSSION

Microarray-based genomic surveys and other high-throughput strategies are becoming increasingly important in biology. These methods typically produce data sets containing massive lists of expressed or repressed genes. This information provides only a starting point from which to begin the investigation of specific biological processes. To begin understanding how cells function within a tissue in a normal or diseased state or respond to specific stimuli, these data must be decoded. To facilitate interpretation and stratify the data, information concerning the biological, biochemical, or molecular functions of the translated proteins must be provided.

The stratification of *in silico* and biological data can be an undertaking of epic proportions when investigators use the typical biological paradigm of examining one element at a time. To overcome this impasse, we created Onto-Express.

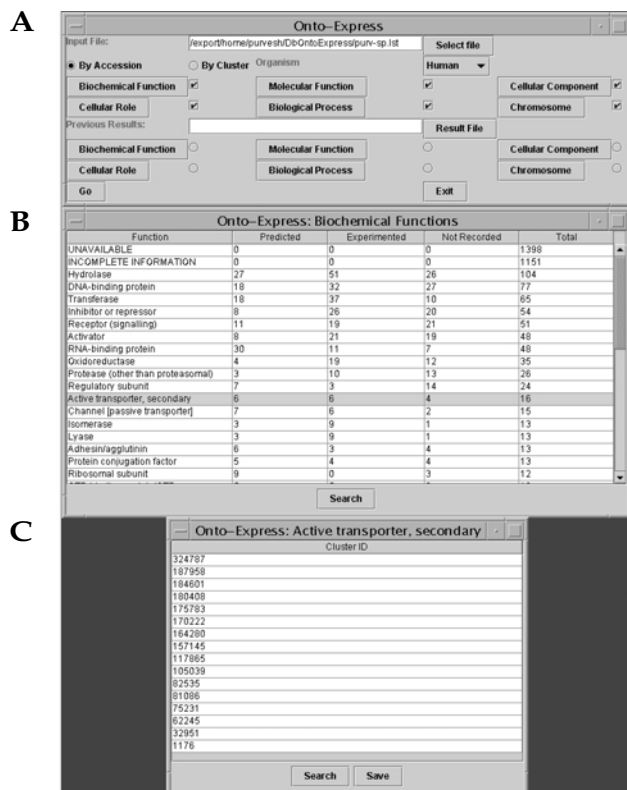


FIG. 1. Program interface windows. Onto-Express uses a graphical user interface. (A) Program window. To initiate the program, the user enters the raw data file path in the appropriate text field. To indicate the format and species contained in the raw data file, the user selects the corresponding radio button (by accession code or by cluster identification code) and organism, for example, human. The checkboxes indicating the desired information are selected and the "Go" button is activated to start data retrieval. (B) Result window. After data retrieval, Onto-Express returns tables that contain the desired information, such as the biochemical functions of the translated proteins. The user can select a function of interest, such as active transporters, secondary, and then click the search button. This will open a window showing the cluster identification numbers that correspond to the function of interest as shown in (C). (C) Search window. The user can highlight one of the cluster identification numbers and click search, which will open the corresponding NCBI web page. Alternatively, the user can click the "save" button and store the generated list in a user-specified file.

This program systematically associates gene expression profiles with the chromosome and cytogenetic gene locations, as well as the biochemical and molecular functions, biological processes, cellular components, and cellular roles of the translated proteins.

Onto-Express starts by reading the input file that exclusively contains single line entries of either accession or cluster identification numbers. The database is built from the `LL_tmpl` file, which is updated weekly (and can be obtained from <ftp://ftp.ncbi.nih.gov/refseq/LocusLink>). After creating the database, Onto-Express queries the NCBI map viewer and retrieves the number of genes on each chromosome. If the input file contains accession numbers, Onto-Express retrieves

the cluster identification number for each accession number and then builds a list of cluster identification numbers. This eliminates redundancy, as multiple accession numbers are often linked to a single cluster identification number. After building the list of cluster identification numbers, Onto-Express retrieves the locus identification number for each cluster. Onto-Express then queries the database using each locus identification number to collect as much information as possible. For each locus this includes biological process, biochemical function, molecular function, cellular role, cellular component, and chromosomal location. The program builds the corresponding list in memory. If the chromosomal location of a cluster is not available in the database created, Onto-Express queries the UniGene database (<http://www.ncbi.nlm.nih.gov/UniGene/>). The software could query other databases such as Ensembl (<http://us.ensembl.org/>) to retrieve additional information or attempt to complete the analysis.

In those cases when an accession number corresponding to a cluster identification number or locus identification number is not available, or for a cluster identification number for which the corresponding locus identification number is not available, it is classified as "Unavailable." If partial information is available for an accession number or for a cluster identification number, the software classifies the observation as "Incomplete."

Onto-Express stores the information collected as simple text files. If the input file contains accession numbers, the software prefixes the output filenames with "accn." Similarly, if the input file contains cluster identification numbers, Onto-Express prefixes the output filenames with "clst." The first line of the output file states which data are stored followed by a semicolon-delimited listing of the acquired data. Once all of the information is collected, Onto-Express displays the information in tabular form as the user requests.

Figure 1 shows the user interface windows of Onto-Express. The user first enters the input filename that contains the list of either accession or cluster identification numbers. Alternatively, one can browse through directories by pressing the "Select File" button to locate the input file. The user then informs Onto-Express of the input file format by selecting either the "By Accession" or "By Cluster" radio button. Next, the user selects the organism studied from the drop-down list. Subsequently, the user chooses the information they would like to retrieve by selecting or deselecting the check boxes. By default, the software collects all the information. To start data collection, the user activates the "Go" button. The "Go" button then remains unavailable until Onto-Express finishes collecting the data. Version 1.0 processes one input file at a time. While Onto-Express is collecting the information, the upper six buttons, labeled "Biochemical Function," "Molecular Function," "Cellular Component," "Cellular Role," "Biological Process," and "Chromosome," are unavailable (that is, grayed out). They are available only after Onto-Express has collected all of the information. When available, the user can click any of the buttons to display the

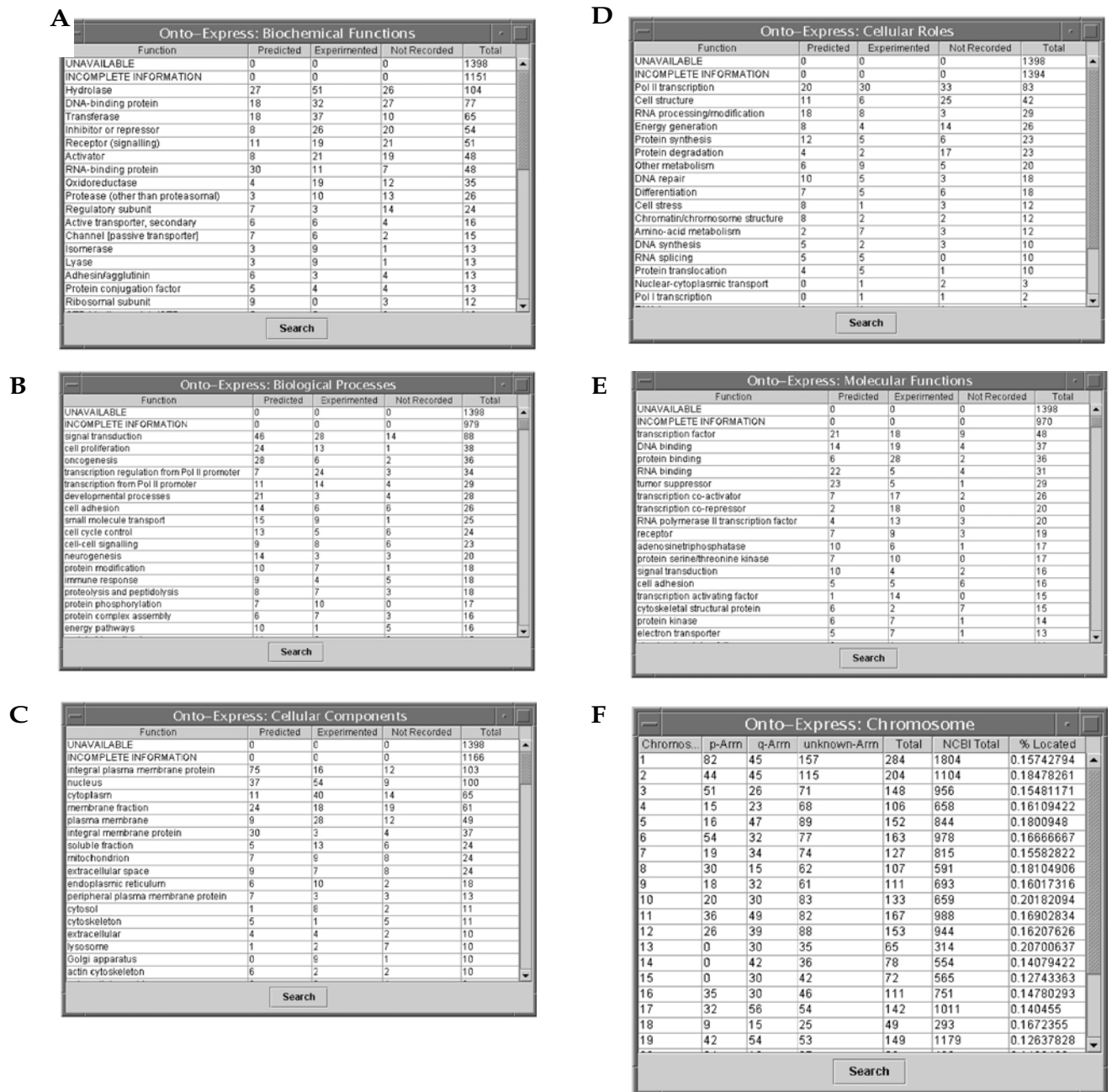


FIG. 2. Result windows. Onto-Express returns six different categories of information. The following panels represent these groupings. (A) Biochemical function. These data depict the principle structural, regulatory, or enzymatic function of the proteins. (B) Biological process. This panel describes the biological “objective” to which the proteins contribute. (C) Cellular component. The information in this panel illustrates the place in a cell where the protein is active. (D) Cellular role. The major biological process involving the protein is demonstrated in this panel. (E) Molecular function. The biochemical activity of the proteins is summarized by these data. The first column in each panel lists specific instances. The second column represents the number of instances that are predicted by similarity or by analogy to other proteins, whereas the third column represents the number of times the specific instances were supported by experimental evidence. The fourth column represents how many times the instance was identified but the type of supporting evidence was not recorded. The fifth column shows the sum of the second, third, and fourth columns. (F) Chromosome. This panel represents data concerning the chromosomal location of the genes assessed. The first column reports the specific chromosome number. The second and third columns represent the number of genes located on p and q arms, respectively. The fourth column explains the number of genes for which the exact location was not reported. The fifth column is the sum of the second, third, and fourth columns. The sixth column represents the total number of genes known to be located on the respective chromosomes and used to standardize the total of genes identified, as shown in the last column.

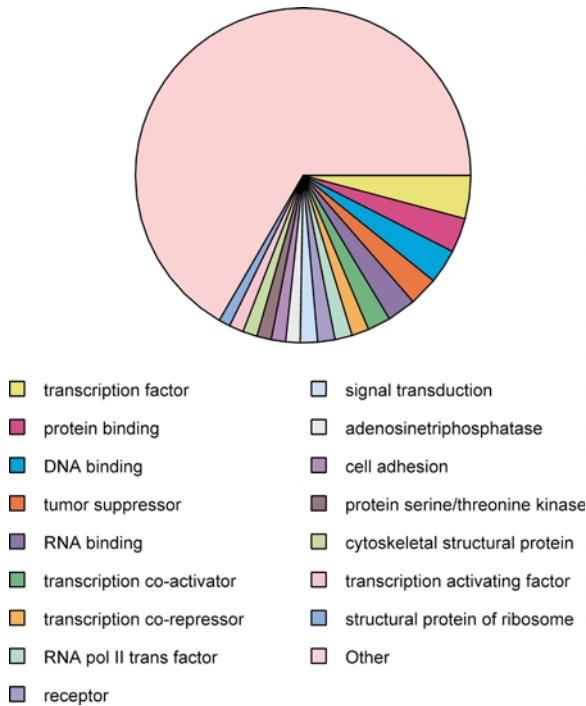


FIG. 3. Ontogeny of spermatozoal RNA. mRNAs were isolated from sperm and used to create labeled cDNAs for microarray analysis. The molecular function of the proteins that represent each expressed sequence tag identified were data mined using Onto-Express. The different sections of the pie chart represent the proportion of proteins identified to have the molecular function that is indicated by color in the legend. "Other" indicates protein groups with fewer than 14 observations.

respective information in a tabular format as shown in Fig. 1B.

While Onto-Express is collecting information associated with the input file, the user may view any of the previously collected data. The user can view the current results at a later point in time by activating the lower six buttons. The user can either enter the output file name or browse through directories to select the previously generated output file of interest.

Figure 2 shows the results for biochemical function, cellular component, cellular role, molecular function, biological process, and chromosomal location. The information is presented in five columns (Figs. 2A–2E). The first column denotes the function name. The second row of this column, labeled "Unavailable," contains the number of hits for which the corresponding locus identification number was not available, while the third row, labeled "Incomplete Information," contains the number of hits for which partial information was available. The second and third columns are labeled "Predicted" and "Experimented" and show the number of accession or cluster identification numbers having predicted and experimental proof provided for a given function. The fourth column presents how many times the function was identified for which supporting evidence was "Not Recorded." The fifth column, labeled "Total," is the sum of

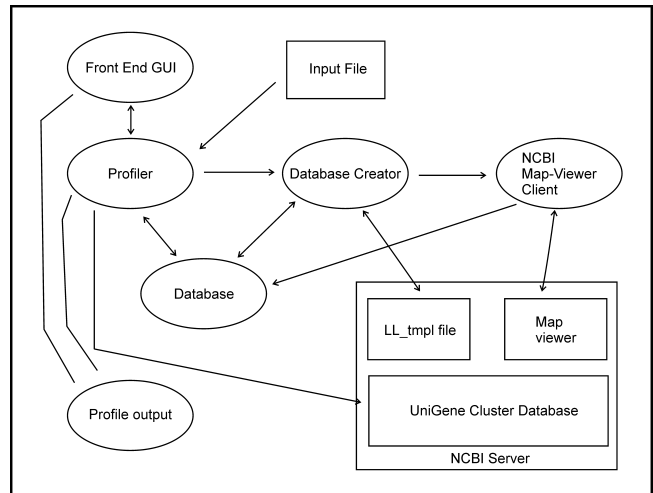


FIG. 4. Software architecture of Onto-Express. Onto-Express uses a GUI interface, a profiler, a database creator, an NCBI map-viewer client, the NCBI server, and database components. The profiler invokes the database creator, which in turn invokes the NCBI map-viewer client. The client queries the NCBI map-viewer and stores the information retrieved in the database. Following the database creator populating the database from the "LL_tmpl" file, the profiler reads the input file. If the input file is comprised of accession numbers, the algorithm retrieves their respective UniGene cluster identification numbers. For each cluster identification number, the algorithm determines the locus identification number to query the database for biochemical function, biological process, cellular component, cellular role, molecular function, and chromosomal location. The profiler stores the information in simple text files, which the front-end GUI displays in tabular format.

second, third, and fourth columns.

Figure 2F shows the related chromosomal information. The first column is the chromosome number, the second column is the total of accession numbers or cluster identification numbers located on the p-Arm, the third is the total of accession numbers or cluster identification numbers located on the q-Arm, followed by the fourth column indicating the number of ESTs identified on that chromosome for which their exact location is not known. The fifth column is the sum of the second, third, and fourth columns. The sixth column, labeled "NCBI Total," shows the number of genes located on that chromosome. The seventh column, labeled "% located," is the ratio of the fifth and the sixth columns expressed as a percentage.

Application of Onto-Express

It is well documented that mature spermatozoa contain mRNA [14–20]. However, the function or purpose of the mRNAs have yet to be established. To begin unraveling the mystery of why spermatozoa contain mRNA, we characterized the genetic fingerprint of mature spermatozoa generated from a series of microarray experiments (G.C.O. *et al.*, manuscript submitted), using Onto-Express. Transcription factors, protein binding factors, and DNA binding proteins are the functional groups containing the largest number of identified proteins (Fig. 3). The

identification of transcription factors was somewhat surprising, bearing in mind that spermatozoa are considered transcriptionally quiescent [21-24]. The identification of DNA and protein binding factors is consistent with spermatozoal mRNAs encapsulating spermatogenic gene expression [20,25,26]. These results illustrate that Onto-Express can transform genetic fingerprints having unknown function into profiles that underscore the biology of the system. When interpreting these observations the question often arises: "What are the identities of the various members of each functional group?" Onto-Express provides a graphical search interface to retrieve the accession or cluster identification numbers associated with specific functions (Figs. 1B and 1C). The accession and cluster identification numbers are linked to NCBI (<http://www.ncbi.nlm.nih.gov/>), allowing the user to effortlessly determine which genes were identified in each functional group. In this manner, the user can associate the biological functions of interest to the primary data obtained from the array. For example, the identification of transcription factors in mature spermatozoa was somewhat surprising. The specific factors (for example, UniGene Cluster Identification Numbers Hs.26102, Hs.108106, Hs.6557, Hs.147049, Hs.182528, Hs.28423, Hs.155402, Hs.79058, Hs.2815) are consistent with their functional relevancy in spermatogenesis [27]. It is clear that this tool will find many uses in our emerging studies of various transcriptomes.

MATERIALS AND METHODS

Figure 4 shows the architecture of Onto-Express. The program has been implemented in Java using an object-oriented approach to facilitate platform-independence and maintenance. The basic components of Onto-Express include the graphical user interface (GUI), a profiler, a database creator, and a local database that is updated from the NCBI server. To generate the gene expression profiles, the GUI invokes the profiler, which in turn invokes the database creator. The database creator invokes the NCBI map-viewer client and, in turn, the client repeatedly queries NCBI map-viewer to collect gene/chromosome information. When the NCBI map-viewer client returns, the database creator reads in the LL_tmpl file and creates a local database. This ensures access to the most current information. After the database creator finishes, the profiler starts reading the input file and repeatedly queries the database to collect the information needed. When the profiler has completed reading the input file and collected all of the required information, it outputs the collected information to text files. The front end GUI then reads the profiler output files and displays the assimilated data as the user requests.

ACKNOWLEDGMENTS

This work was supported in part by NIH grant HD36512 awarded to S.A.K. G.C.O. is supported in part by Wayne State University School of Medicine Dean's Post-Doctoral Fellowship.

RECEIVED FOR PUBLICATION OCTOBER 30;
ACCEPTED DECEMBER 12, 2001.

REFERENCES

- Schena, M., Shalon, D., Davis, R., and Brown, P. O. (1995). Quantitative monitoring of gene expression patterns with a complementary DNA microarray. *Science* **270**: 467-470.
- Eisen, M. B., Spellman, P. T., Brown, P. O., and Botstein, D. (1998). Cluster analysis and display of genome-wide expression patterns. *Proc. Natl. Acad. Sci. USA* **95**: 14863-14868.
- Lockhart, D. J., et al. (1996). DNA expression monitoring by hybridization of high density oligonucleotide arrays. *Nat. Biotechnol.* **14**: 1675-1680.
- Shalon, D., Smith, S. J., and Brown, P. O. (1996). A DNA microarray system for analyzing complex DNA samples using two-color fluorescent probe hybridization. *Genome Res.* **6**: 639-645.
- Ramsay, G. (1998). DNA chips: State-of-the-art. *Nat. Biotechnol.* **16**: 40-44.
- Iyer, V. R., et al. (1999). The transcriptional program in the response of human fibroblasts to serum. *Science* **283**: 83-87.
- Perou, C. M., et al. (2000). Molecular portraits of human breast tumours. *Nature* **406**: 747-752.
- Golub, T. R., et al. (1999). Molecular classification of cancer: class discovery and class prediction by gene expression monitoring. *Science* **286**: 531-537.
- Tamayo, P., et al. (1999). Interpreting patterns of gene expression with self-organizing maps: methods and application to hematopoietic differentiation. *Proc. Natl. Acad. Sci. USA* **96**: 2907-2912.
- Kudoh, K., et al. (2000). Monitoring the expression profiles of doxorubicin-induced and doxorubicin-resistant cancer cells by cDNA microarray. *Cancer Res.* **60**: 4161-4166.
- Baxevas, A. D. (2000). The molecular biology database collection: an online compilation of relevant database resources. *Nucleic Acid Res.* **28**: 1-7.
- Ashburner, M., et al. (2000). Gene ontology: tool for the unification of biology. The Gene Ontology Consortium. *Nat. Genet.* **25**: 25-29.
- Consortium, T. G. (2001). Creating the gene ontology resource: design and implementation. *Genome Res.* **8**: 1425-1433.
- Concha, I., et al. (1993). U1 and U2 snRNA are localized in the sperm nucleus. *Exp. Cell Res.* **204**: 378-338.
- Kumar, G., Patel, D., and Naz, R. K. (1993). c-MYC mRNA is present in human sperm cells. *Cell. Mol. Biol. Res.* **39**: 111-117.
- Chiang, M. H., Steuerwald, N., Lambert, H., Main, E. K., and Steinleitner, A. (1994). Detection of human leukocyte antigen class I messenger ribonucleic acid transcripts in human spermatozoa via reverse transcription-polymerase chain reaction. *Fertil. Steril.* **61**: 276-280.
- Rohwedder, A., Liedigk, O., Schaller, J., Glander, H., and Werchau, H. (1996). Detection of mRNA transcripts of $\beta 1$ integrins in ejaculated human spermatozoa by nested reverse-transcription polymerase chain reaction. *Mol. Hum. Reprod.* **2**: 499-505.
- Miller, D., et al. (1999). A complex population of RNAs in human ejaculate spermatozoa: implications for understanding molecular aspects of spermiogenesis. *Gene* **237**: 385-392.
- Richter, W., Dettmer, D., and Glander, H. J. (1999). Detection of mRNA transcripts of cyclic nucleotide phosphodiesterase subtypes in ejaculated spermatozoa. *Mol. Hum. Reprod.* **5**: 732-736.
- Wykes, S. M., Miller, D., and Krawetz, S. A. (2000). Mammalian spermatozoal mRNAs: tools for the functional analysis of male gametes. *J. Submicrosc. Cytol. Pathol.* **32**: 77-81.
- Kierszenbaum, A. I., and Tres, L. L. (1975). Structural and transcriptional features of the mouse spermatid genome. *J. Cell Biol.* **65**: 258-270.
- Stewart, T. A., Bellve, A. R., and Leder, P. (1984). Transcription and promoter usage of the myc gene in normal somatic and spermatogenic cells. *Science* **226**: 202-210.
- Clermont, Y., and Leblond, C. (1955). Spermiogenesis of man, monkey, ram and other mammals as shown by the periodic acid-schiff technique. *Am. J. Anat.* **96**: 229-253.
- Balhorn, R., et al. (1999). Protamine mediated condensation of DNA in mammalian sperm. In *The Male Gamete: From Basic Science to Clinical Applications* (C. Gagnon, Ed.), pp. 55-70. Cache River Press, Vienna, IL.
- Kramer, J., and Krawetz, S. A. (1997). RNA in sperm: implications within the sperm genome. *J. Biol. Chem.* **271**: 11619-11622.
- Miller, D. (1997). RNA in the ejaculate spermatozoon: a window into molecular events in spermatogenesis and a record of the unusual requirements of haploid gene expression and post-meiotic equilibration. *Mol. Hum. Reprod.* **3**: 669-676.
- Nelson, J. E., and Krawetz, S. A. (1995). Computer assisted promoter analysis of a human sperm specific nucleoprotein gene cluster. *DNA Seq.* **5**: 329-337.