

DNA Array Analysis in a Microsoft® Windows® Environment

BioTechniques 32:110-119 (January 2002)

T. Conway, B. Kraus, D.L. Tucker, D.J. Smalley, A.F. Dorman, and L. McKibben
The University of Oklahoma,
Norman, OK, USA

ABSTRACT

Microsoft® Windows®-based computers have evolved to the point that they provide sufficient computational and visualization power for robust analysis of DNA array data. In fact, smaller laboratories might prefer to carry out some or all of their analyses and visualization in a Windows environment, rather than alternative platforms such as UNIX. We have developed a series of manually executed macros written in Visual Basic for Microsoft Excel® spreadsheets, that allows for rapid and comprehensive gene expression data analysis. The first macro assigns gene names to spots on the DNA array and normalizes individual hybridizations by expressing the signal intensity for each gene as a percentage of the sum of all gene intensities. The second macro streamlines statistical consideration of the confidence in individual gene measurements for sets of experimental replicates by calculating probability values with the Student's *t* test. The third macro introduces a threshold value, calculates expression ratios between experimental conditions, and calculates the standard deviation of the mean of the log ratio values. Selected columns of data are copied by a fourth macro to create a processed data set suitable for entry into a Microsoft Access® database. An Access database structure is described that allows simple queries across multiple experiments and export of data into third-party data visualization soft-

ware packages. These analysis tools can be used in their present form by others working with commercial *E. coli* membrane arrays, or they may be adapted for use with other systems. The Excel spreadsheets with embedded Visual Basic macros and detailed instructions for their use are available at <http://www.ou.edu/microarray>.

INTRODUCTION

With the publication of the first genome-wide expression data six years ago (13), the power of the DNA array was immediately evident. Since then, reviewers have predicted a revolution in the way biologists conduct their research (3–5,12). Today the trend toward widespread use of DNA arrays continues. Vast amounts of DNA array data are accumulating, and the need for standardized array annotation and data representation is being addressed (<http://www.ebi.ac.uk/microarray/MGED/>). Relational databases are being developed to handle microarray data storage in a format that facilitates data processing and visualization, allowing researchers to analyze and interpret their experiments and disseminate the data (16). Generally, the database management software used is Oracle, Web interfaces are written in JAVA or XML, and scripts used for data processing, retrieval, etc. are written in Perl, C, or other programming languages. However, for the common microarray user, the implementation of such a system may be beyond their budget or exceed their actual needs. For these users, a central archive is an attractive alternative for data storage. An example of this, the EcoReg Consortium (<http://gobi.lbl>).

[gov/~ecoreg/index.html](http://ecoreg/index.html)), is being established as a public database for storage and manipulation of *E. coli* microarray and proteome data. Still, it will be necessary for individual users of consortium databases to process, analyze, and format the data for submission.

DNA arrays have been used to examine the genetics and physiology of the comprehensive biological model, *E. coli*. Some research groups have used commercial membrane-based DNA arrays (2,17,18), while others have employed DNA microarrays (6,11,19,20) or oligonucleotide arrays (15). Several practical issues regarding the use of whole-genome arrays have been addressed, and the power of this technology as a means for deducing the physiological state of the bacterial cell is now well established.

In our laboratory, we routinely use membrane-based DNA arrays for *E. coli* gene expression profiling. Multiple replicates of each experimental condition are processed and analyzed with manually executed macros written in Visual Basic and run in Microsoft® Excel® spreadsheets on Microsoft Windows®-based computers. These macros, statistical analysis, and data processing protocols are described here and are freely available to the scientific community (<http://www.ou.edu/microarray>).

MATERIALS AND METHODS

Software and System Requirements

The analysis tools described here are written for Microsoft Excel and Microsoft Access®: Microsoft Office® 97 or higher is required. Optimal system

requirements include a Windows-compatible computer with a PII processor and 128 Mb RAM or higher. Macintosh® computers running Microsoft Office 98 for Macintosh can also be used to run the macros in Excel.

Example *E. coli* Data Set

In this study, we compare a sample data set of *E. coli* MG1655 grown at pH 7.4 and pH 5.5 under otherwise identical conditions. Cultures were grown aerobically in 50 mL MOPS (pH 7.4) or MES (pH 5.5) minimal glucose (0.2%) medium (8) in 250-mL fleakers (Corning, Acton, MA, USA) at 37°C with 300 rpm agitation and harvested in mid-logarithmic growth phase. RNA isolation, radioactive labeling during cDNA synthesis, and hybridization to DNA array membranes were described previously (17).

Description of Raw Data

We routinely use Panorama *E. coli* Gene Arrays™ (Sigma-Genosys, The Woodlands, TX, USA) for gene expression profiling. Phosphorimaging of a hybridized membrane array produces a TIFF image file that must be further processed for data analysis. The image analysis software (ArrayVision™ version 5.1; Imaging Research, St. Catharines, Ontario, Canada) makes use of a customized template to accommodate three grid layers (3 × 1; 16 × 24; 4 × 4) according to the design of the Panorama *E. coli* gene arrays. The spot labeling protocol was edited such that each spot is named by its unique array coordinate, allowing the spot intensity measurements to be easily associated with the correct gene identifiers in subsequent processing steps. The customized ArrayVision template file for analysis of Panorama *E. coli* gene arrays, and detailed instructions for its use, are available at <http://www.ou.edu/microarray>. The spot intensities are represented in a row-column format and are exported into Excel spreadsheets for further analysis.

Data Processing

The macros and sample analyses can be downloaded from <http://www.ou.edu/microarray>.

Follow the links to “Macroarray”, “Data Analysis”, and then “Spot-Finding and Image Quantitation” or “*E. coli* Data Analysis (software downloads)”; alternatively, the macroarray section of the site can be accessed directly at <http://www.ou.edu/microarray/macroarray.htm>. Raw DNA array data, exported from ArrayVision, are processed in a series of three Excel workbooks that are used to manually execute four macros written in Visual Basic (Table 1). These macros are designed to filter the data and calculate statistics to allow for further data analysis and interpretation. Detailed, step-by-step instructions for use of these analysis tools are provided on our Web site.

RESULTS AND DISCUSSION

The workbooks, macros, and subroutines used for DNA array data analysis are outlined in Table 1. The subroutines can be run in order individually, or the macro containing all relevant subroutines can be run once to execute all subroutines. In the following section, we provide an overview of important statistical considerations, the specific processing steps, and outcomes.

Statistical Significance

Various approaches for attaching significance to DNA array data have been published, including a simple “rule-of-thumb” criterion for the value of the expression ratio (2,20). Some researchers have used the standard deviation from the mean of the expression ratios as an indicator of confidence (6,17). Arfin et al. (1) applied the Student’s *t* test to experimental replicates and considered the *P* value to be the most important indicator of significance. Richmond et al. (11) combined a confidence interval obtained with the Student’s *t* test and a rule-of-thumb criterion for the expression ratio. Others prefer to consider the significance of a single experimental condition based on the coefficient of variation (18). A precedent for statistical analysis of array data has yet to be firmly established, and a standard is clearly needed. Whatever statistical approach is adopted, it is essential that DNA array exper-

iments are properly replicated and the uncertainty that lies behind individual gene measurements be considered to attach significance to data sets.

We advocate the use of at least two replicates of each experimental condition. Membrane arrays typically have duplicate spots for each gene, and each spot is considered to be a separate determination, providing a total of four determinations for the two replicates. Because membrane arrays are hybridized with a single labeled target mRNA and normalized independently (in effect a one-color experiment), the variation in the measurement is at the level of the raw data, not the measurement of the ratio (unlike the statistical approach that is popular with two-color microarrays and involves internal normalization of the measurements). The uncertainty that lies behind individual gene measurements can be variously calculated as the standard deviation of the determinations or the coefficient of variation. Since we are usually interested in the statistical significance of differences between an experimental condition and a control, we prefer the Student’s *t* test as a means for calculating this probability, based on the uncertainty of the replicate measurements in both conditions. The Student’s *t* test is best applied to natural log-transformed normalized data sets (7). Generally, a *P* value of less than 0.05 is chosen to indicate a 95% probability that the difference in gene expression between conditions is significant. However, it has been pointed out that with very large data sets (e.g., a bacterium with 5000 genes) choosing a value of *P* < 0.05 means that there could be up to 250 false positives in the data set (1). Thus, the researcher is left with two choices: to lower the *P* value to a level where no false positives are expected (*P* < 0.0002) or to consider a second statistical metric that, when combined with a reasonable *P* value (*P* < 0.05), is an excellent indicator of significance of a ratio value.

We use the standard deviation of the mean of the log ratios—within the context of the *P* value—to indicate significant up- or down-regulation of gene expression. This approach is meaningful where the expression level of the majority of genes does not change significantly between conditions and where the researcher is interested in genes that show



Table 1. Useful Purpose of Workbooks, Visual Basic Macros, and Subroutines Used in This Study

| Workbook | Macro | Subroutine | Purpose |
|---------------------|------------------|---------------------|--|
| Image Data Cruncher | AllDataCrunched6 | ArvAllSort1 | associates array coordinate with spot number |
| | | Nameall2 | associates spot number with unique identifier for gene |
| | | CalcPercentage3 | normalizes data by expressing each spot as percentage of sum of all spot intensities |
| | | Cleanup4 | reorganizes data and calculates average values for duplicate spots |
| | | Statistics5 | calculates averages of genomic DNA controls and blank spots |
| 2-Replicate-Stats | AllAnalysis8 | OrganizebySpotNo1 | sorts each of the four data sets individually by spot number |
| | | CalculateAverages2 | calculates averages of volumes and percentage values for the control and test replicates |
| | | copyvaluesintoPRow3 | copies and pastes percentage values into a separate spreadsheet for calculation of <i>P</i> values |
| | | CalculateLn4 | copies percentage values into a separate spreadsheet and natural log transforms data |
| | | CalculatePRow5 | calculates the <i>P</i> value for the raw data by application of the Student's <i>t</i> test |
| | | CalculatePLn6 | calculates the <i>P</i> value for the log transformed data by application of the Student's <i>t</i> test |
| | | CopyAllValues7 | copies and pastes data used for ratio calculations into separate spreadsheet |
| Data Analysis | AllAnalysis6 | SpotSort1 | sorts control and test data sets by spot number |
| | | DataSort2 | copies data set into spreadsheet used for ratio calculations, sorts by total percentage value |
| | | ThresholdRatios3 | calculates ratio of Test/Control using threshold of total percentage value for 500th lowest gene |
| | | Cleanup4 | reorganizes data and calculates log (10) of ratio |
| | | Stats5 | calculates standard deviation of log ratio values and correlation between Test and Control |
| | (Manual Step) | | copy and paste special values Prow and PLn values from 2-Replicate-Stats to Data Analysis |
| | MakeDBsheet | | reorganizes data and copies into separate spreadsheet for entry into Access database |

substantially different expression. The standard deviation for the log ratios is calculated, and only those genes that differ by more than three standard deviations (99.9% confidence in each tail) from the mean of the log ratio (usually zero, or no change) are considered. In practice, emphasis is placed on those genes that have expression ratios greater than three standard deviations from the mean and have a reasonable probability of being significantly different between the conditions, based on a *P* value of less than 0.05. Where there are four or

more determinations for each gene, the *P* value can be lowered to less than 0.005 with little change seen in the number of genes that are considered to vary significantly between conditions.

There may also be situations when the researcher is interested in changes in gene expression that are not a full three standard deviations from the mean but are still significant (i.e., where the differences in gene expression between conditions are subtle yet meaningful). In this case the Student's *t* test can be used as the sole measure of

significance, but the *P* value must be adjusted to ensure that false positives are avoided. One approach for this is to apply the Bonferroni correction that describes a *P* value for significance in a large data set (7). Various strategies for implementing this correction factor have been described, and its proper use is somewhat controversial (14). The Bonferroni correction effectively lowers the *P* value to a point where false positives are avoided and consideration is given only to those genes for which there is a high degree of confidence in

BioComputing/BioInformatics >>

the ratio value. In this light, the Bonferroni correction seems to be a reasonable statistical tool but may be too stringent for some considerations.

"Image Data Cruncher" Workbook

To begin the data analysis process, the raw data from an experimental replicate is copied and pasted into a blank spreadsheet named "arvdata" in an Excel workbook named "Image Data Cruncher" containing the macro named "AllDataCrunched6" (Figure 1). The first subroutine in the macro, "ArvAllSort1", copies and pastes subsets from the "arvdata" spreadsheet into a second spreadsheet named "allfields" that contains information provided by the membrane manufacturer for associating the array location with a spot number that is unique to each target on the array. The second subroutine, "Nameall2", copies and pastes the data in the "allfields" spreadsheet into a third spreadsheet that associates the spot number with a unique identifier and associated genome annotation information for each gene.

Differences in spot intensities between replicate experiments arise from normal experimental variation, such as differences in growth conditions, ra-

dioactive nucleotide incorporation efficiency, hybridization conditions, or image acquisition. To compare experimental replicates (separate cultures) or technical replicates (same culture and same RNA sample), the data from each array must be normalized. The third subroutine, "CalcPercentage3", normalizes arrays; if the data are not normalized, then the values for replicate experiments, when plotted, will not pass through zero or be directly proportional (Figure 2). Array experiments can be normalized by expressing each gene-specific spot relative to an internal control, if a suitable set of control spots is present on the arrays. Unfortunately, the intensities of the genomic DNA control spots on the Panorama *E. coli* membranes vary significantly (data not shown) and are not reliable for normalization. In the absence of an internal standard, the preferred approach for normalization is to express each gene-specific spot as a fraction of the sum of all gene spots ($n = 4290$), a strategy that at once considers all variables that lie behind the array image (17,18). The "CalcPercentage3" subroutine normalizes the entire data set, expressing each spot as a percentage of the sum of all spots on the array.

The fourth subroutine, "Cleanup4",

reorganizes the data to facilitate subsequent data processing steps. The fifth subroutine, "Statistics5", calculates the averages and statistics for the blank spots (empty spots between genomic DNA control spots) and null spots (empty spots within array). However, we have not found these values to offer a reliable means for establishing background on the array or for empirical determination of a threshold value and therefore do not use them. The "AllDataCrunched6" macro in "Image Data Cruncher" runs all five subroutines at once and results in a file that contains the raw data, normalized data, and associated genome annotation information for each gene on the array.

"2-Replicate-Stats" Workbook

The second workbook, "2-Replicate-Stats", contains a macro named "AllAnalysis8" that is used to calculate the probability that the average of the experimental (test) replicates is significantly different from the average of the control replicates (Figure 3). The four replicate data sets are sequentially copied from the crunched data files, beginning with the first and then the second replicate of the control, followed by the first and then the second replicate of the experimental, and pasted into the "Enter Data (2 Replicates)" worksheet. The first subroutine, "OrganizebySpotNo1", sorts the four data sets by spot number, which aligns the gene-specific data in rows. The second subroutine, "CalculateAverages2", calculates the mean of the normalized (percentage) values for the four spot intensities from each experimental condition (two spots for each gene per membrane). The subroutine, "copyvaluesintoPraw3", copies the percentage values into a separate spreadsheet, and the "CalculateLn4" subroutine transforms the raw percentage values by the natural log while copying them to an additional spreadsheet.

The subroutines "CalculatePraw5" and "CalculatePLn6" are used to calculate the P values for the raw and log transformed data, respectively, by application of the Student's t test to the four determinations for each of the control and experimental conditions. The last subroutine, "CopyAllValues7", reorganizes the data by pasting the data

| SPOT No. | vol a | vol a Pct | SPOT No. | vol b | vol b Pct | avg vol | avg pct | Array Coord | Gene | k# | gene product | Origin |
|----------|---------------|-----------|------------|----------|-----------|-------------|-------------|--------------|-------|-------|---------------------------------|----------|
| 1 | 40883.67 | 0.020257 | 1 | 41059.42 | 0.020657 | 40971.5447 | 0.020457075 | Field3-B2-2 | araD | k0061 | L-rubulose-5-phosphate 4 Carbon | Carbon |
| 2 | 16647.26 | 0.008249 | 2 | 16598.7 | 0.008351 | 16622.98015 | 0.008299618 | Field1-C1-1 | araA | k0062 | L-arabinose isomerase | Carbon |
| 4278 | 4277.6922.287 | 0.00343 | 4277 | 5749.096 | 0.002892 | 6336.63655 | 0.003161121 | Field3-J14-1 | yjiV | b4378 | orf, hypothetical protein | Hypoth |
| 4279 | 4278.21073.11 | 0.010444 | 4278 | 22467.52 | 0.011304 | 21773.31305 | 0.010873876 | Field3-L14-1 | yjiJ | b4380 | orf, hypothetical protein | Hypoth |
| 4280 | 4279.27804.96 | 0.013777 | 4279 | 28848.07 | 0.013507 | 27326.51535 | 0.013642069 | Field3-N14-1 | yjiJ | b4395 | orf, hypothetical protein | Hypoth |
| 4281 | 4280.189045.5 | 0.093668 | 4280 | 163008.7 | 0.082011 | 176027.0892 | 0.087839155 | Field3-P14-1 | smg | b4397 | orf, hypothetical protein | Hypoth |
| 4282 | 4281.46142.38 | 0.022862 | 4281 | 44128.99 | 0.022202 | 45135.6858 | 0.02253202 | Field3-B16-1 | yjiK | b4391 | putative ATP-binding con | Putative |
| 4283 | 4282.8303.315 | 0.004114 | 4282 | 8065.549 | 0.004058 | 8184.43195 | 0.004085963 | Field3-D16-1 | yjiK | b4394 | orf, hypothetical protein | Hypoth |
| 4284 | 4283.12584.51 | 0.006235 | 4283 | 10118.89 | 0.005091 | 11351.70225 | 0.005663105 | Field3-F16-1 | gpmB | b4395 | phosphoglyceromutase 2 | Central |
| 4285 | 4284.8759.552 | 0.00434 | 4284 | 8330.674 | 0.004191 | 8545.11275 | 0.004265689 | Field3-H16-1 | creA | b4397 | orf, hypothetical protein | Putative |
| 4286 | 4285.5978.756 | 0.002962 | 4285 | 5672.92 | 0.002854 | 5825.83795 | 0.002908207 | Field3-J16-1 | yjiY | b4402 | orf, hypothetical protein | Hypoth |
| 4287 | 4286.13865.35 | 0.00687 | 4286 | 12916.84 | 0.006499 | 13391.0961 | 0.00688425 | Field3-L16-1 | lasT | b4403 | orf, hypothetical protein | Hypoth |
| 4288 | 4287.8247.403 | 0.004086 | 4287 | 7055.855 | 0.00355 | 7651.62655 | 0.00381812 | Field3-N16-1 | b0701 | b0701 | rbsC protein in rbs elem | Hypoth |
| 4289 | 4288.10035.64 | 0.005002 | 4288 | 11253.23 | 0.005662 | 10674.43325 | 0.00531186 | Field3-P16-1 | b2086 | b2086 | orf, hypothetical protein | Hypoth |
| 4290 | 4289.3531.696 | 0.00175 | 4289 | 4012.138 | 0.002019 | 3771.9167 | 0.001884202 | Field3-B18-1 | yjiM | b4404 | orf, conceptual translatio | Hypoth |
| 4291 | 4290.17934.93 | 0.008866 | 4290 | 17865.13 | 0.008988 | 17900.03295 | 0.008937199 | Field3-D18-1 | yjiM | b4405 | orf, conceptual translatio | Hypoth |
| 4292 | 2.02E+08 | | totals: | 1.99E+08 | | | | | | | | |
| 4293 | | | | | | | | | | | | |
| 4294 | | | | | | | | | | | | |
| 4295 | | | | | | | | | | | | |
| 4296 | | | | | | | | | | | | |
| 4297 | Blank | 3534.639 | AvgBlank | 4557.739 | | | | | | | | |
| 4298 | Blank | 3564.43 | AvgBlankSt | 2113.05 | | | | | | | | |
| 4299 | Blank | 4444.713 | Avg10ng | 305853.3 | | | | | | | | |
| 4300 | Blank | 2989.546 | Avg10ngStd | 56767.96 | | | | | | | | |
| 4301 | Blank | 2469.33 | Avg5ng | 124537.3 | | | | | | | | |
| 4302 | Blank | 3227.155 | Avg5ngStd | 18307.27 | | | | | | | | |
| 4303 | Blank | 3933.516 | AvgNull | 1809.19 | | | | | | | | |
| 4304 | Blank | 2502.616 | AvgNullStd | 501.2595 | | | | | | | | |
| 4305 | Blank | 3707.056 | PctNull | 0.000896 | | | | | | | | |
| 4306 | Blank | 6623.475 | PctNullStd | 0.000248 | | | | | | | | |
| 4307 | Blank | 4929.158 | PctNull3St | 0.000745 | | | | | | | | |
| 4308 | Blank | 4906.04 | | | | | | | | | | |
| 4309 | Blank | 2173.164 | 10ngnorm | 0.001515 | | | | | | | | |
| 4310 | Blank | 4513.198 | 5ngnorm | 0.001234 | | | | | | | | |

Figure 1. "Image Data Cruncher" workbook used to process raw gene expression profile data by normalizing expression values and assigning gene names to array spot coordinates.

columns to be used in subsequent steps into a separate spreadsheet named "All Values". Once the "AllAnalysis8" macro is executed, the "2-Replicate-Stats" workbook serves as an archive for the raw and normalized data from the replicates being compared and contains the *P* values that are associated with the ratio calculations in the "Data Analysis" workbook.

"Data Analysis" Workbook

The third workbook, "Data Analysis", contains two macros. The first macro, "AllAnalysis6", is used to calculate the log ratio of the expression levels in the experimental versus the control condition (Figure 4). The average normalized data from the "All Values" spreadsheet in the "2-Replicate-Stats" workbook are pasted into the "crunched data" spreadsheet in the "Data Analysis" workbook. The first subroutine, "SpotSort1", sorts the data by spot number such that the gene specific data are aligned in rows. The second subroutine, "DataSort2", copies the data and pastes them into the "DataAnalysis" spreadsheet to be used for ratio calculations.

The third subroutine, "ThresholdRatios3", is used to determine a threshold value for ratio calculation and calculates the absolute value of the ratio of the Test/Control such that genes that are more highly expressed in the test condition are given a positive value and genes that are more highly expressed in the control are given a negative value. The threshold value is chosen to represent the limit of detection of an expressed gene (i.e., the signal intensity at which spots are considered to be significantly higher than the array background). Any spot intensity that falls below the threshold value is raised to that value to obtain a reasonable ratio in cases where a gene is expressed below the threshold value in at least one of the two experimental conditions. Ideally, the threshold would be determined independently for each gene, based on the local spot background and the known cross-hybridization to other expressed genes in the sample, but this is not possible because of the dense packing of some membrane arrays and the lack of prior knowledge as to the num-

ber and extent of gene expression in a given growth condition. These factors make determination of the threshold value difficult. We have chosen a conservative approximation of the threshold value corresponding to the 500th lowest expressed gene based on the average of the normalized expression levels in the two conditions. This threshold value is similar to that obtained by visual inspection of array images to determine the faintest of gene-specific spots and is reasonable in light of the predicted number of expressed genes based on the number of mRNA species in the *E. coli* cell and the arrangement of genes in operons (1380 mRNAs \times average 2.5 genes per operon = 3450 expressed genes) (9). Others have calculated threshold values as three standard deviations above the mean of "blank" spots (10), which corresponds to approximately 214 in the example data set used in this study. If desired, researchers can write alternative approaches for threshold determination

into the "ThresholdRatios3" subroutine, or the threshold level can be edited as described on the Web site.

The fourth subroutine, "Cleanup4", reorganizes the data and calculates the log (base 10) of the expression ratio. The fifth subroutine, "Stats5", calculates the standard deviation of the mean of the log ratios. The "AllAnalysis6" subroutine executes all five subroutines and results in a "DataAnalysis" spreadsheet that contains the averaged raw and normalized data, corresponding genome annotation information for each gene on the array, and the ratios. The "AllAnalysis6" macro concludes with creation of two empty columns that are used for manual pasting of the *P* value associated with each ratio calculation from the "2-Replicate-Stats" workbook. Finally, the "MakeDB" macro in the "Data Analysis" workbook can be executed to reorganize and paste key data columns into a separate spreadsheet that can be used for data entry into a suitable database (Figure 5).

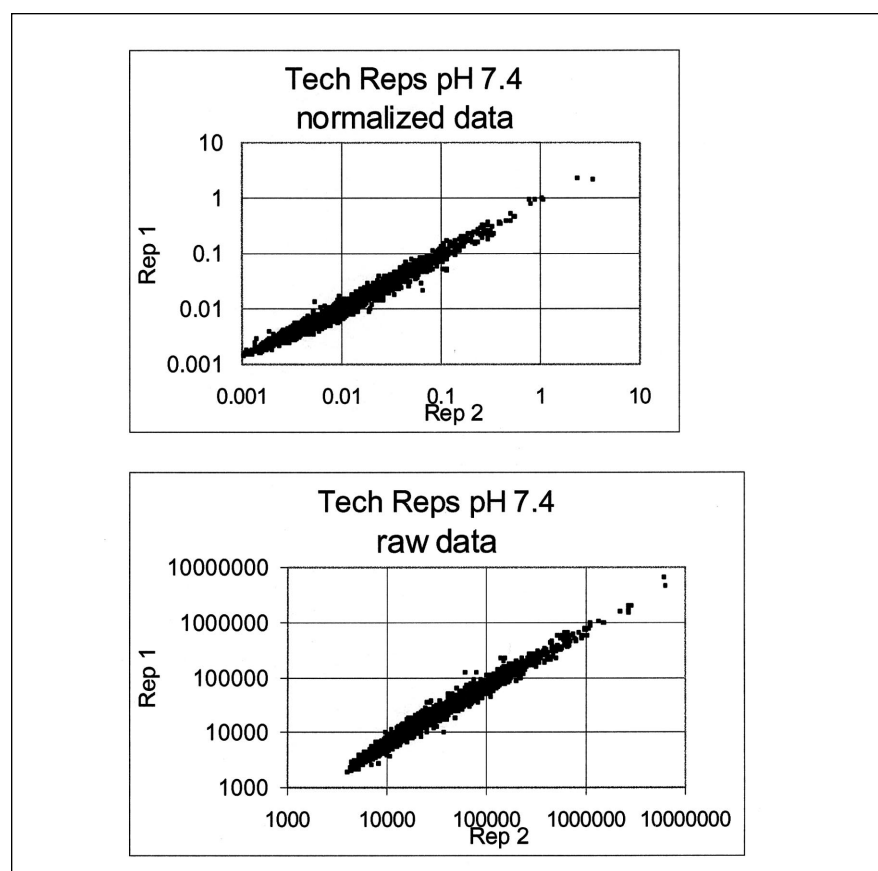


Figure 2. Scatter plot view of normalized (top) and raw (bottom) data showing linearity of experimental replicates for normalized data.



| | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O |
|------|--------|-----------|----------|----------|------------|----------|---|---|---|---|---|---|---|---|---|
| 1 | b# | log ratio | PRAW | PLN | ControlPct | TestPct | | | | | | | | | |
| 2 | b0061 | -0.04279 | 0.471414 | 0.442091 | 0.018735 | 0.016977 | | | | | | | | | |
| 4274 | b4367 | -0.04106 | 0.24586 | 0.253645 | 0.006367 | 0.005793 | | | | | | | | | |
| 4275 | b4363 | 0 | 0.856713 | 0.975312 | 0.002597 | 0.002533 | | | | | | | | | |
| 4276 | b4367 | -0.32545 | 0.00589 | 0.001072 | 0.01022 | 0.00483 | | | | | | | | | |
| 4277 | b4377 | 0.205016 | 0.00408 | 0.00991 | 0.012307 | 0.019732 | | | | | | | | | |
| 4278 | b4378 | 0.170739 | 0.020405 | 0.039649 | 0.003982 | 0.0069 | | | | | | | | | |
| 4279 | b4380 | -0.01075 | 0.646651 | 0.684352 | 0.011744 | 0.011457 | | | | | | | | | |
| 4280 | b4385 | -0.07138 | 0.062612 | 0.063921 | 0.01436 | 0.012893 | | | | | | | | | |
| 4281 | b4387 | -0.31204 | 7.94E-05 | 2.3E-06 | 0.087073 | 0.042446 | | | | | | | | | |
| 4282 | b4391 | -0.0343 | 0.297945 | 0.303097 | 0.024896 | 0.022996 | | | | | | | | | |
| 4283 | b4394 | -0.07738 | 0.166286 | 0.15983 | 0.004821 | 0.004034 | | | | | | | | | |
| 4284 | b4395 | -0.02365 | 0.624003 | 0.62169 | 0.006253 | 0.005921 | | | | | | | | | |
| 4285 | b4397 | 0.039349 | 0.233647 | 0.228714 | 0.004776 | 0.005229 | | | | | | | | | |
| 4286 | b4402 | 0.130621 | 0.000483 | 0.000137 | 0.003106 | 0.005155 | | | | | | | | | |
| 4287 | b4403 | 0.009503 | 0.677438 | 0.638552 | 0.007273 | 0.007434 | | | | | | | | | |
| 4288 | b40701 | -0.01846 | 0.553509 | 0.599018 | 0.004207 | 0.004032 | | | | | | | | | |
| 4289 | b2098 | 0.03718 | 0.369149 | 0.351572 | 0.006154 | 0.006704 | | | | | | | | | |
| 4290 | b4404 | 0 | 0.918079 | 0.974584 | 0.002222 | 0.002196 | | | | | | | | | |
| 4291 | b4405 | -0.0163 | 0.645221 | 0.703073 | 0.010204 | 0.003628 | | | | | | | | | |
| 4292 | | | | | | | | | | | | | | | |
| 4293 | | | | | | | | | | | | | | | |
| 4294 | | | | | | | | | | | | | | | |
| 4295 | | 0.003816 | | | | | | | | | | | | | |
| 4296 | | 0.116571 | | | | | | | | | | | | | |
| 4297 | | 0.902625 | | | | | | | | | | | | | |
| 4298 | | | | | | | | | | | | | | | |
| 4299 | | | | | | | | | | | | | | | |
| 4300 | | | | | | | | | | | | | | | |
| 4301 | | | | | | | | | | | | | | | |
| 4302 | | | | | | | | | | | | | | | |
| 4303 | | | | | | | | | | | | | | | |
| 4304 | | | | | | | | | | | | | | | |
| 4305 | | | | | | | | | | | | | | | |
| 4306 | | | | | | | | | | | | | | | |
| 4307 | | | | | | | | | | | | | | | |
| 4308 | | | | | | | | | | | | | | | |
| 4309 | | | | | | | | | | | | | | | |
| 4310 | | | | | | | | | | | | | | | |

Figure 5. "MakeDB" worksheet containing selected data columns ready for entry into gene expression experiments database.

Benefits of Semi-Automated Data Processing

Our earliest protocol for membrane array data acquisition and analysis required more than 10 man-hours for each experimental replicate (17). Thus, it was important to automate the data analysis process to the greatest extent possible. By taking advantage of improvements in commercially available software and by writing a series of macros in Visual Basic for semi-automated data processing in Microsoft Excel, the time was shortened to 5 min per experimental replicate. In addition, it is critical that the possibility of human error be eliminated when carrying out a large number of manipulations of massive data sets such as those involved in gene expression profiling. The macros described here effectively minimize the level of manual data processing by automating much of the process. The integration of four macros in three workbooks provides for a significant level of human supervision of an otherwise error-free, automated process.

ACKNOWLEDGMENTS

The authors wish to thank Bill Cuevas and Simon Sims for information be-

fore publication and helpful discussions. This work was supported by National Institutes of Health grant nos. to RO1 AI48945-01 and RR-01-005 to T.C.

REFERENCES

- Arfin, S.M., A.D. Long, E.T. Ito, L. Toller, M.M. Riehle, E.S. Paegle, and G.W. Hatfield. 2000. Global gene expression profiling in *Escherichia coli* K12. The effects of integration host factor. *J. Biol. Chem.* 275:29672-29684.
- Barbosa, T.M. and S.B. Levy. 2000. Differential expression of over 60 chromosomal genes in *Escherichia coli* by constitutive expression of MarA. *J. Bacteriol.* 182:3467-3474.
- Brown, P.O. and D. Botstein. 1999. Exploring the new world of the genome with DNA microarrays. *Nat. Genet.* 21(1 Suppl):33-37.
- DeRisi, J.L. and V.R. Iyer. 1999. Genomics and array technology. *Curr. Opin. Oncol.* 11:76-79.
- Ferea, T.L. and P.O. Brown. 1999. Observing the living genome. *Curr. Opin. Genet. Dev.* 9:715-722.
- Khodursky, A.B., B.J. Peter, N.R. Cozzarelli, D. Botstein, P.O. Brown, and C. Yanofsky. 2000. DNA microarray analysis of gene expression in response to physiological and genetic changes that affect tryptophan metabolism in *Escherichia coli*. *Proc. Natl. Acad. Sci. USA* 97:12170-12175.
- Long, A.D., H.J. Mangalam, B.Y. Chan, L. Toller, G.W. Hatfield, and P. Baldi. 2001. Improved statistical inference from DNA microarray data using analysis of variance and a Bayesian statistical framework. Analysis of global gene expression in *Escherichia coli* K12. *J. Biol. Chem.* 276:19937-19944.
- Neidhardt, F.C., P.L. Bloch, and D.F. Smith. 1974. Culture medium for enterobacteria. *J. Bacteriol.* 119:736-747.
- Neidhardt, F.C., J.L. Ingraham, and M. Schaechter. 1990. Physiology of the Bacterial Cell: A Molecular Approach. Sinauer Associates, Sunderland, MA.
- Pomposiello, P.J., M.H. Bennik, and B. Dimple. 2001. Genome-wide transcriptional profiling of the *Escherichia coli* responses to superoxide stress and sodium salicylate. *J. Bacteriol.* 183:3890-3902.
- Richmond, C.S., J.D. Glasner, R. Mau, H. Jin, and F.R. Blattner. 1999. Genome-wide expression profiling in *Escherichia coli* K-12. *Nucleic Acids Res.* 27:3821-3835.
- Schena, M., R.A. Heller, T.P. Theriault, K. Konrad, E. Lachenmeier, and R.W. Davis. 1998. Microarrays: biotechnology's discovery platform for functional genomics. *Trends Biotechnol.* 16:301-306.
- Schena, M., D. Shalon, R.W. Davis, and P.O. Brown. 1995. Quantitative monitoring of gene expression patterns with a complementary DNA microarray. *Science* 270:467-470.
- Scherf, U., D.T. Ross, M. Waltham, L.H. Smith, J.K. Lee, L. Tanabe, K.W. Kohn, W.C. Reinhold et al. 2000. A gene expression database for the molecular pharmacology of cancer. *Nat. Genet.* 24:236-244.
- Selinger, D.W., K.J. Cheung, R. Mei, E.M. Johansson, C.S. Richmond, F.R. Blattner, D.J. Lockhart, and G.M. Church. 2000. RNA expression analysis using a 30 base pair resolution *Escherichia coli* genome array. *Nat. Biotechnol.* 18:1262-1268.
- Sherlock, G., T. Hernandez-Boussard, A. Kasarskis, G. Binkley, J.C. Matese, S.S. Dwight, M. Kaloper, S. Weng et al. 2001. The Stanford Microarray Database. *Nucleic Acids Res.* 29:152-155.
- Tao, H., C. Bausch, C. Richmond, F.R. Blattner, and T. Conway. 1999. Functional genomics: expression analysis of *Escherichia coli* growing on minimal and rich media. *J. Bacteriol.* 181:6425-6440.
- Tao, H., R. Gonzalez, A. Martinez, M. Rodriguez, L.O. Ingram, J.F. Preston, and K.T. Shanmugam. 2001. Engineering a homothanol pathway in *Escherichia coli*: increased glycolytic flux and levels of expression of glycolytic genes during xylose fermentation. *J. Bacteriol.* 183:2979-2988.
- Wei, Y., J.M. Lee, C. Richmond, F.R. Blattner, J.A. Rafalski, and R.A. LaRossa. 2001. High-density microarray-mediated gene expression profiling of *Escherichia coli*. *J. Bacteriol.* 183:545-556.
- Zimmer, D.P., E. Soupene, H.L. Lee, V.F. Wendisch, A.B. Khodursky, B.J. Peter, R.A. Bender, and S. Kustu. 2000. Nitrogen regulatory protein C-controlled genes of *Escherichia coli*: scavenging as a defense against nitrogen limitation. *Proc. Natl. Acad. Sci. USA* 97:14674-14679.

Received 29 June 2001; accepted 30 August 2001.

Address correspondence to:

Dr. Tyrrell Conway
Department of Botany and Microbiology
The University of Oklahoma
Norman, OK 73069-0245, USA
e-mail: tconway@ou.edu