



March 13, 2007  
External Review Draft

# **U.S. Environmental Protection Agency DRAFT**

## **Interim Guidance for Microarray-Based Assays: Data Submission, Quality, Analysis, Management, and Training Considerations**

**Prepared for the U.S. Environmental Protection Agency  
by Members of the Genomics Workgroup,  
a Group Tasked by EPA's Science Policy Council**

**Science Policy Council  
U.S. Environmental Protection Agency  
Washington, DC 20460**

### **NOTICE**

**This document is an External Review draft. It has not been formally released by the U.S. Environmental Protection Agency and should not at this stage be construed to represent Agency position**

## **DISCLAIMER**

This draft interim guidance, when finalized, will represent EPA's current thinking on this topic. It does not create or confer any legal rights for or on any person or operate to bind the public. The use of any mandatory language in this document is intended to describe laws of nature, scientific principles, or technical requirements and is not intended to impose any legally enforceable rights or obligations. Alternative approaches may be used if the approach satisfies the requirements of the applicable statutes and regulations. If you would like to discuss an alternative approach (you are not required to do so), you may contact the EPA staff responsible for implementing this guidance. Mention of trade names or commercial products does not constitute endorsement of recommendation for use.

Note: This is an external review draft, and is not approved for final publication.

---

---

**Genomics Microarray Workgroup  
Co-Chairs**

William H. Benson  
Office of Research and  
Development

Kathryn Gallagher  
Office of the  
Science Advisor

J. Thomas McClintock  
Office of Pollution  
Prevention and Toxics

Kerry Dearfield  
Office of the Science Advisor  
(2004 - June 2005)

**Science Policy Council Staff**

Jeremy Johnson  
(2004)

**Subgroup Co-Chairs**

**Performance Approach to  
Quality Assessment**

David Lattier, ORD  
Susan Lundquist, OEI

**Data Management**

Susan Hester, ORD  
Joseph Retzer, OEI

**Data Submission**

Greg Miller, OPEI  
Doug Wolf, ORD

**Training**

Bobbye Smith, Region 9 RSL  
Julian Preston, ORD

**Data Analysis**

David Dix, ORD  
Brenda Groskinsky, Region 7 RSL

**Microbial Source Tracking**

Jorge Santo Domingo, ORD  
Ron Landy Region 3 RSL

**Additional Coordinating Committee Members**

Wafa Harrouk, US FDA  
Lee Hofmann, OSWER  
Robert Kavlock, ORD  
Rita Schoeny, OW

**Genomics Workgroup Lead for the Science Policy Council**

Larry Reiter  
Office of Research and Development

---

## Genomics Microarray Workgroup Members

Gregory Akerman, OPPTS  
Wenjun Bao, ORD  
David Bencic, ORD  
Lynn Bradley, OEI  
Kevin Cavanaugh, ORD  
Barbara Collins, ORD  
Brion Cook, OPPTS  
Don Delker, ORD  
Michelle Embry, OPPTS  
Robin Gonzalez, OEI  
Susan Griffin, Region 8  
Stephanie Harris, Region 10  
Belinda Hawkins, ORD  
Kenneth Haymes, OPPTS  
Michael Hemmer, ORD  
Todd Holdermann, OPPTS  
Gene Hsu, ORD  
Margo Hunt, OEI  
Sid Hunter, ORD  
Channa Keshava, ORD  
Steven Kueberuwa, OW  
Mitch Kostich, ORD  
Richard Leukroth, OPPTS  
Nancy McCarroll, OPPTS  
Jesse Meiller, OPPTS  
Elizabeth Mendez, OPPTS  
Ann Miracle, ORD  
Ines Pagan, ORD  
Santhini Ramasamy, OPPTS  
Ann Richard, ORD  
Mitch Rosen, ORD  
Phil Sayre, OPPTS  
Judy Schmid, ORD  
John Sykes, ORD  
Freshteh Toghrol, OPPTS  
Mark Townsend, OPPTS  
Nancy Wentworth, OEI  
Lori White, ORD  
Witold Winnik, ORD  
Steve Young, OEI

---

## **Additional Genomics Resources**

### **Genomics Training Workgroup Members**

Barbara Abbott, ORD  
Gilberto Alvarez, Region 5  
Michele Burgess, OSWER  
Michelle Embry, OPPTS  
Audrey Galizia, ORD  
Karen Hamernik, OPPTS  
Steven Kueberuwa, OW  
David Lattier, ORD  
David Lee, ORD  
Roseanne Lorenzana, Region 10  
Marian Olsen, Region 2  
Jennifer Seed, OPPTS

### **Microbial Source Tracking Workgroup Members**

Bobbye Smith, Region 9  
Jafrul Hasan, OW  
James Goodrich, ORD  
Rita Schoeny, OW  
Robin Oshiro, OW  
Roland Hemmett, Region 2  
Sally Gutierrez, ORD

---

## Table of Contents

<b>ACRONYMS.....</b>	<b>VIII</b>
<b>EXECUTIVE SUMMARY .....</b>	<b>1</b>
<b>1.0 INTRODUCTION .....</b>	<b>5</b>
1.1 BACKGROUND .....	5
1.2 OVERVIEW OF GENOMIC SCIENCE .....	6
1.3 EMERGING IMPACTS OF GENOMICS TECHNOLOGIES.....	8
1.4 PURPOSE AND INTENT OF THIS DOCUMENT .....	11
<b>2.0 THE PERFORMANCE APPROACH TO QUALITY ASSURANCE FOR MICROARRAYS .....</b>	<b>12</b>
<b>3.0 DATA SUBMISSION GUIDANCE .....</b>	<b>14</b>
3.1 INTRODUCTION.....	14
3.2 ABSTRACT.....	14
3.3 EXPERIMENTAL DESIGN .....	15
3.4 ARRAY DESIGN .....	16
3.5 BIOMATERIALS.....	16
3.6 HYBRIDIZATION .....	17
3.7 MEASUREMENTS .....	17
<b>4.0 DATA ANALYSIS GUIDANCE .....</b>	<b>19</b>
4.1 INTRODUCTION.....	19
4.2 DATA ANALYSIS.....	20
4.3 DATA EVALUATION.....	23
4.4 DATA ANALYSIS CONCLUSIONS .....	24
<b>5.0 DATA MANAGEMENT.....</b>	<b>25</b>
<b>6.0 RECOMMENDATIONS .....</b>	<b>28</b>
6.1 TRAINING NEEDS AND RECOMMENDATIONS .....	28
6.2 COLLABORATIVE DEVELOPMENT OF GENOMIC TOOLS FOR DATA ANALYSIS AND DATA MANAGEMENT .....	31
6.3 APPLYING THIS INTERIM GUIDANCE FOR MICROARRAY-BASED ASSAYS TO CASE STUDIES .....	32
6.4 UPDATING GENOMICS GUIDANCE AS NEEDED .....	32
<b>REFERENCES .....</b>	<b>33</b>
<b>APPENDIX A: EPA QUALITY SYSTEM AND THE PERFORMANCE APPROACH TO QUALITY MEASUREMENT SYSTEMS.....</b>	<b>35</b>
<b>APPENDIX B: MIAME-BASED DATA SUBMISSION TABLES.....</b>	<b>51</b>
TABLE B.1 ABSTRACT .....	51
TABLE B.2 EXPERIMENTAL DESIGN .....	51
TABLE B.3 ARRAY DESIGN.....	53
TABLE B.4 BIOMATERIALS .....	56
TABLE B.5 HYBRIDIZATION.....	61
TABLE B.6 MEASUREMENTS.....	62
<b>APPENDIX C: GENOMICS DATA EVALUATION RECORD (GDER) TEMPLATE.....</b>	<b>64</b>
<b>APPENDIX D: GENOMICS DATA EVALUATION RECORD (GDER) FOR ALACHLOR (SAMPLE) ....</b>	<b>67</b>
<b>APPENDIX E: MIAME GLOSSARY .....</b>	<b>77</b>

**APPENDIX F: ADDITIONAL GLOSSARY FROM GENOMICS WHITE PAPER.....82**

**APPENDIX G: CONTENT AND INSTRUCTIONAL GOALS FOR THE THREE LEVELS OF  
GENOMICS TECHNICAL TRAINING: .....88**

## ACRONYMS

CBI	Confidential Business Information
cDNA	Complementary Deoxyribonucleic Acid
CEBS	Chemical Effects in Biological Systems knowledgebase
cRNA	Complementary Ribonucleic Acid
CWA	Clean Water Act
DER	Data Evaluation Record
DNA	Deoxyribonucleic Acid
DQO	Data Quality Objective
EPA	Environmental Protection Agency
FACS	Fluorescence Activated Cell Sorter
FDA	Food and Drug Administration
FNR	False Negative Rate
FPR	False Positive Rate
gDER	Genomics Data Evaluation Record
HPV	High Production Volume
IRB	Institutional Review Board
IVT	<i>In Vitro</i> Transcription
JPEG	Joint Photographic Experts Group
MAGE	Microarray And Gene Expression
MAGE-OM	Microarray And Gene Expression - Object Model
MGED	Microarray Gene Expression Data
MIAME	Minimal Information About Microarray Experiments
MOA	Mode of Action
MOPS-EDTA	[MOPS] 3-(N-Morpholino) propanesulfonic acid], [EDTA] ethylenediaminetetraacetic acid
MPSS	Massively Parallel Signature Sequencing
mRNA	Messenger RNA
MQO	Measurement Quality Objective
MST	Microbial Source Tracking
NHEERL	National Health and Environmental Effects Research Laboratory
NIEHS	National Institute of Environmental Health Sciences
NPDES	National Pollutant Discharge Elimination System
OEI	Office of Environmental Information
OPPTS	Office of Prevention, Pesticides and Toxic Substances
ORD	Office of Research and Development
OSWER	Office of Solid Waste and Emergency Response
OW	Office of Water
PCR	Polymerase Chain Reaction
PMN	Pre-Manufacture Notification
PMT	Photomultiplier Tube
QA	Quality Assurance
QAARWP	Quality Assurance Annual Report and Work Plan
QC	Quality Control

---



QMP	Quality Management Plan
qPCR	Quantitative Polymerase Chain Reaction
qRT-PCR	Quantitative Reverse Transcriptase PCR
RFU	Relative Fluorescent Unit
RNA	Ribonucleic Acid
RNase	Ribonuclease
RTP	Research Triangle Park
RT-PCR	Reverse-Transcription Polymerase Chain Reaction
SAGE	Serial Analysis of Gene Expression
SNP	Single Nucleotide Polymorphism
SOPs	Standard Operating Procedures
SPC	Science Policy Council
TIFF	Tagged Image File Format
TMDL	Total Maximum Daily Load
U.S. EPA	U.S. Environmental Protection Agency

---

## EXECUTIVE SUMMARY

1  
2  
3       The mapping of diverse animal, plant, and microbial species genomes using molecular  
4 technologies has significantly affected research across all areas of the life sciences. The current  
5 understanding of biological systems is rapidly changing in ways previously unimagined and  
6 novel applications of this technology have already been commercialized. These advances in  
7 genomics will have significant implications for risk assessment policies and regulatory decision  
8 making. In 2002, the U.S. Environmental Protection Agency (EPA or “the Agency”) issued its  
9 Interim Policy on Genomics (U.S. EPA, 2002a) that communicated the Agency’s initial  
10 approach to using genomics information in risk assessment and decision making. The Interim  
11 Policy described genomics as the study of all the genes of a cell or tissue, at the DNA  
12 (genotype), mRNA (transcriptome), or protein (proteome) level. While noting that the  
13 understanding of genomics is far from established, the Agency stated that such data may be  
14 considered in the decision making process, but that these data alone are insufficient as a basis for  
15 decisions.

16  
17       Following the release of the Interim Policy, the Science Policy Council (SPC) created a  
18 cross-EPA Genomics Task Force and charged it with examining the broader implications  
19 genomics is likely to have on Agency programs and policies. The Genomics Task Force  
20 developed a Genomics White Paper entitled “Potential Implications of Genomics for Regulatory  
21 and Risk Assessment Applications at EPA” (U.S. EPA, 2004). That document identified four  
22 areas likely to be influenced by the generation of genomics information within EPA and the  
23 submission of such information to EPA: 1) prioritization of contaminants and contaminated sites,  
24 2) monitoring, 3) reporting provisions; and 4) risk assessment. One critical need in the area of  
25 technical development was identified: the need to establish a framework for analysis and  
26 acceptance criteria for genomics information for scientific and regulatory purposes. The Task  
27 Force recommended that the Agency charge a workgroup to establish such a framework and in  
28 doing so consider the performance of assays across genomic platforms (*e.g.*, reproducibility,  
29 sensitivity, pathway analysis tools) and the criteria for accepting genomics data for use in a risk  
30 assessment (*e.g.*, assay validity, biologically meaningful response).

1 In 2004, the Genomics Technical Framework and Training Workgroup was formed with  
2 the responsibility to ensure that the technical framework and training activities build upon the  
3 Agency's Interim Policy on Genomics while continuing to engage other interested parties.  
4 Information developed by these workgroups will be used by EPA program offices and regions to  
5 determine the applicability of specific genomics information to the evaluation of risks under  
6 various statutes.

7  
8 To this end, the Genomics Technical Workgroup considered all of the "omics"  
9 technologies and applications and decided that an interim guidance document on the use of data  
10 generated by DNA microarray technology would be most beneficial to the Agency and regulated  
11 community at this time. Consequently, this document provides recommendations regarding: 1)  
12 data that should be considered for submission to the Agency for microarray studies, 2) the use of  
13 a performance approach to microarray quality assessment parameters, 3) data analysis  
14 approaches for microarrays, and 4) data management and storage issues for microarray data  
15 submitted to or used by the Agency. The guidance applies to both human health and ecological  
16 DNA microarray data.

17  
18 With respect to experimental performance considerations, the Genomics Workgroup  
19 concluded that quality issues are critical considerations in the application of new technologies  
20 such as genomics. The Genomics Workgroup recommends that the Agency not prescribe  
21 specific methods to be used in microarray experiments at this time, but instead provide general  
22 guidance on the recommended performance of microarray experiments in order to obtain data of  
23 the quality required for a specific use; this guidance is provided herein. Investigators submitting  
24 data to the Agency in support of regulatory decision making, methods development, and  
25 technical transfer, may want to consider, in addition to compliance with MIAME (Minimal  
26 Information About Microarray Experiments) Workgroup standards  
27 (<http://www.mged.org/Workgroups/MIAME/miame.html>), the performance-related experimental  
28 and system factors outlined in this document (Appendix A). Further activities on the part of  
29 investigators to address experimental performance issues will serve to strengthen scientific  
30 arguments and experimental claims.

1 This document also provides information regarding submission of microarray data to  
2 EPA to ensure appropriate review and consistent evaluation of data from multiple sources. In  
3 accordance with accepted practice, it is recommended that submissions include sufficient  
4 information to allow an independent reviewer to reconstruct how the data were collected and  
5 analyzed. This approach allows reviewers to judge the quality of the data and the strength of any  
6 conclusions. Many scientific journal editors grappling with these issues have adopted the  
7 MIAME guidelines as a standard for submission of microarray data as part of a submitted  
8 publication. A slightly modified version of MIAME is proposed as the microarray data  
9 submission template for EPA; this submission template will be subject to change as the  
10 technology evolves.

11  
12 With regard to data analysis, the Genomics Workgroup concluded that a systematic  
13 approach for genomics data evaluation is necessary for the further use of such data in risk  
14 assessments. A genomics Data Evaluation Record template is provided herein as a way to  
15 present and organize data from genomics studies in order to derive information necessary for a  
16 regulatory application (see Appendix C for the Genomics Data Evaluation Record [DER]). A  
17 completed sample DER is also provided in Appendix D to facilitate the use of the template. An  
18 overview of issues to be considered in analyzing microarray data is also provided. The transfer  
19 of these evaluations, and the underlying genomics data, into searchable, electronic databases will  
20 be essential to making the data useful in risk assessments. Furthermore, development of  
21 databases containing gene expression profiles for a wide variety of chemicals should facilitate  
22 creation of statistical/computational methods that will help predict the toxic potential of a  
23 chemical.

24  
25 Due to potentially large volumes of genomic and associated toxicological data, it is  
26 essential that the Agency consider the development of a complete data management solution.  
27 The functional needs of a solution of this magnitude would minimally include items listed in the  
28 section on data management. In addition, this Agency data management solution should address  
29 needs unique to scientifically-based risk assessments, confidential and proprietary data security,  
30 public access, and other aspects of regulatory application. It should be noted that consistency,  
31 scientific and operational robustness, common access, and availability in a scalable environment

---

1 are data management needs for an Agency data management solution. While the Agency has  
2 begun to utilize bioinformatics research approaches, both intramurally (e.g., the National Center  
3 for Computational Toxicology in EPA's Office of Research and Development [ORD]) and  
4 extramurally (Environmental Bioinformatics Centers in North Carolina and New Jersey funded  
5 by EPA's Science to Achieve Results (STAR) Program), an Agency-wide data management  
6 solution integrating genomics, toxicological, and other key data required for regulatory  
7 applications is now necessary.

8  
9 The document concludes with the Genomics Workgroup's recommendations to the  
10 Agency for follow-up activities to this interim guidance including: 1) further development of the  
11 outlined training materials and modules, to be offered throughout the Agency to risk assessors  
12 and decision makers who will be faced with the challenge of interpreting and applying genomics  
13 information, 2) continued collaboration of EPA personnel with staff from other federal agencies  
14 and stakeholders in the development of tools for the analysis of genomics data, 3) application of  
15 this guidance to a series of case studies to evaluate its utility in risk assessment and regulatory  
16 applications; and 4) the updating of this guidance as needed as the technology evolves.

17  
18 This document is intended to provide information to the regulated community and other  
19 interested parties regarding submitting microarray data to the Agency and to provide guidance  
20 for EPA reviewers in evaluating such data and/or information. This interim guidance can be  
21 used by EPA program offices to determine the applicability of specific genomics information to  
22 the evaluation of chemical risks.

23

## 1 **1.0 Introduction**

### 3 **1.1 Background**

5 The mapping of diverse animal, plant, and microbial species genomes using molecular  
6 technologies has significantly affected research across all areas of the life sciences. The current  
7 understanding of biological systems is rapidly changing in ways previously unimagined and  
8 novel applications of this technology have already been commercialized. These scientific and  
9 technological advances have spurred many federal agencies to consider the far-reaching  
10 implications for policy, regulation, and society as a whole.

12 In 2002, EPA released the Interim Policy on Genomics (U.S. EPA, 2002a)  
13 communicating its initial approach to using genomics information in risk assessment and  
14 decision making (<http://www.epa.gov/osa/spc/genomics.htm>). This policy describes genomics as  
15 the study of all the genes of a cell or tissue, at the DNA (genotype), mRNA (transcriptome), or  
16 protein (proteome) level. The Interim Policy notes that while genomics offers the opportunity to  
17 understand how an organism responds at the gene expression level to stressors in the  
18 environment, understanding such molecular events with respect to adverse ecological and/or  
19 human health outcomes is far from established. This policy states that while genomics data may  
20 be considered in the decision making process at this time, these data alone are insufficient as a  
21 basis for decisions. Consequently, currently EPA will only consider genomics information for  
22 assessment purposes on a case-by-case basis.

24 Following the release of the Interim Policy, the Science Policy Council (SPC) created a  
25 cross-EPA Genomics Task Force and charged it with examining the broader implications  
26 genomics is likely to have on Agency programs and policies. To that end, the Genomics Task  
27 Force developed a Genomics White Paper entitled “Potential Implications of Genomics for  
28 Regulatory and Risk Assessment Applications at EPA” (USEPA, 2004,  
29 [www.epa.gov/osa/genomics.htm](http://www.epa.gov/osa/genomics.htm)). The Task Force identified scenarios to describe various  
30 circumstances under which EPA might receive these data. Four areas were identified as those

---

1 likely to be influenced by the generation of genomics information within EPA and the  
2 submission of such information to EPA: 1) prioritization of contaminants and contaminated sites,  
3 2) monitoring, 3) reporting provisions; and 4) risk assessment. The Task Force also identified  
4 several challenges and/or critical needs that included research, technical development, and  
5 capacity (*i.e.*, strategic hiring practices and training).

6  
7 The Genomics Task Force recommended that the Agency charge a workgroup with  
8 developing a technical framework for analysis and acceptance criteria for genomics information  
9 for scientific and regulatory purposes. The Genomics White Paper identified issues that need to  
10 be considered in developing such a framework including the performance of assays across  
11 genomic platforms (*e.g.*, reproducibility, sensitivity, pathway analysis tools) and the criteria for  
12 accepting genomics data for use in a risk assessment (*e.g.*, assay validity, biologically  
13 meaningful response).

14  
15 In June, 2004, the Genomics Technical Framework and Training Workgroup was  
16 established with representatives from ORD, numerous program offices (OPPTS, OSWER, OW,  
17 OEI, OPEI) and regional offices (2, 3, 5, 7, 8, and 9). The Genomics Workgroup was comprised  
18 of a Coordinating Committee, several technical genomics guidance workgroups (Performance  
19 Approach Quality Assurance Workgroup, Data Submission Workgroup, Data Analysis  
20 Workgroup, and a Data Management and Storage Workgroup), a Training Workgroup, and a  
21 Microbial Source Tracking Workgroup. The Genomics Workgroup's responsibility was to  
22 ensure that the technical framework and training activities build upon the Agency's Interim  
23 Policy on Genomics while continuing to engage other interested parties. This document will be  
24 used by EPA program offices and regions to determine the applicability of specific genomics  
25 information to the evaluation of risks under various statutes.

## 26 27 **1.2 Overview of Genomic Science**

28  
29 As a means of introduction to genomics and its potential impact on regulatory decision  
30 making, it is important to understand the basic principles behind genomic technology. Only  
31 about 1-2% of the human DNA actually codes for RNA that can be translated into proteins. This

1 1-2% is considered to be the theoretical functional genome. Any particular cell type (*i.e.*, from  
2 various organs or species) will have its own practical functional genome, which is a subset of the  
3 entire functional genome that encodes for functional proteins in that cell. The functional genome  
4 for any cell type can be assessed by determining the messenger RNA (mRNA) profile of the cell,  
5 tissue, or organ. The mRNA copies the necessary portion of the cell's DNA code and transports  
6 this information to the ribosomes where protein synthesis occurs. Thus, the assessment of  
7 mRNA profiles is called functional genomics. Such profiles are constructed using microarrays  
8 that contain all (or a sampling) of a cell's functional genome. Hybridization of a DNA copy  
9 (cDNA) of the mRNA that is being actively produced by the cell to these microarrays  
10 demonstrates which genes are currently active in that cell. Within the 98-99% of DNA not  
11 coding for RNA message is information that affects the activity of the functional genome by  
12 influencing where and when genes are active in an organism. Thus both coding and noncoding  
13 DNA are important in organismal function and response to perturbations.

14  
15 The study of a cell's protein composition is called proteomics. Currently, it is possible to  
16 analyze only a fraction of a cell's proteins, but rapid advances in this field will allow more  
17 complete profiling in the near future. Another discipline of biology analyzes biofluids and  
18 tissues to determine the profiles of endogenous metabolites present under normal conditions or  
19 when the organism has been affected by factors such as exposure to environmental chemicals.  
20 This type of whole cell analysis is called metabolomics (or metabolic profiling). In order to  
21 understand how a cell functions under normal or stressed circumstances, it is necessary to  
22 characterize the proteins that are manufactured by the cell, as well as endogenous metabolites.  
23 This facilitates an understanding of global metabolism and how proteins interact along  
24 biochemical pathways. This approach describes the area of systems biology, in which the cell,  
25 tissue, or organism is considered as a complete, albeit complex, system.

26  
27 Broadly defined, genomics tools provide the means to examine changes in gene  
28 expression, protein, and metabolite profiles within the cells and tissues, in contrast to current risk  
29 assessment methods which are restricted to whole organism effects or changes in single  
30 biochemical pathways. Genomics tools have the potential to provide detailed data about the  
31 underlying biochemical mechanisms of disease or toxicity (*i.e.*, disease etiology, biochemical

---



1 pathways), sensitive measures of exposures to chemicals, new approaches to detecting effects of  
2 such exposures, and methods for predicting genetic predispositions that may possibly lead to  
3 disease or higher sensitivity to particular stressors in the environment.

4  
5 Another type of application is chemical identification. By utilizing genomic expression  
6 profiles it is possible to identify and classify environmental contaminants. For example,  
7 Hamadeh *et al.* (2002a,b) found chemical-specific gene expression profiles in liver tissue of  
8 exposed rats. The authors demonstrated that 24-hour exposure to compounds from the same  
9 chemical class (peroxisome proliferators) resulted in gene expression profiles that were unique  
10 but more similar to each other than to patterns corresponding to exposure to a chemical of a  
11 different class (enzyme inducers). These gene expression profiles were associated with  
12 differences in histopathology between the different chemical classes following longer durations  
13 These and other published works indicate the utility of genomic approaches in chemical  
14 identification and in investigations of mode-of-action of chemical hazards.

### 17 **1.3 Emerging Impacts of Genomics Technologies**

18  
19 Toxicology has been moving from observation of changes in tissue histology,  
20 physiology, and chemistry to a mechanistic understanding through assessment of large scale  
21 changes of gene activity within those tissues. Identification of changes in gene expression using  
22 microarrays is becoming an important tool for informing our understanding of toxicological  
23 processes as well as informing the hazard identification process and mode of action analysis as  
24 part of safety and risk assessment. As the price of conducting microarray experiments declines  
25 and an appreciation of their value increases, their use for basic research and as part of the  
26 environmental regulatory process is likely to increase.

27  
28 The use of data generated by microarray technology in peer reviewed scientific  
29 publications has grown exponentially over the last few years. Microarray technology allows  
30 monitoring of changes in gene expression across thousands of genes, or even entire genomes or  
31 proteomes in response to experimentally manipulated or natural conditions. We are now

1 beginning to understand several important toxicological processes in terms of changes in the  
2 activity of single genes or ensembles of genes acting in concert. The identification of these  
3 changes is increasingly the product of the use of microarray technology. As a result of these  
4 research trends, EPA anticipates receiving increasing volumes of microarray data from  
5 environmental researchers, and as a part of the regulatory process. In order to ensure optimal  
6 utilization of these data, EPA has developed this guidance to address the quality, submission,  
7 analysis, and storage of microarray data.

8  
9 While many new genomic technologies do exist, most are not as yet ready for application  
10 in risk/safety assessment and decision making. Therefore, it is important for the Agency to  
11 consider how these genomic technologies might be incorporated into existing programs. It  
12 should be noted that genomics will not fundamentally alter the risk assessment process, but is  
13 expected to serve as a powerful tool for evaluating the exposure to and effects of environmental  
14 stressors and will offer a means to simultaneously examine a number of response pathways.  
15 EPA and other regulatory agencies are beginning to address the use of genomics data for various  
16 risk assessment applications, including the need to establish a link between genomic alterations  
17 and adverse outcomes of regulatory concern. Given the rapidly evolving nature of genomics  
18 technologies, care should be taken to develop an acceptable scheme to simplify and refine the  
19 risk-related information and to distinguish it from the large amount of complex scientific and  
20 statistical data available. This strategy should remain dynamic and fluid in anticipation of  
21 continuing technical evolution at the molecular levels (*e.g.*, DNA, RNA, and protein levels).  
22 Furthermore, bioinformatic approaches for data acquisition and analysis, including technologies  
23 designed to store and analyze the profusion of data generated from microarray analyses, should  
24 be considered in parallel with the data generating methods. Finally, many scientific, policy,  
25 ethical, and legal concerns developing along with the emergence of this science will need to be  
26 addressed.

27  
28 The Interim Policy on Genomics provides guidance concerning how and when genomics  
29 information should be used to assess the risks of environmental contaminants under the various  
30 regulatory programs implemented by the Agency at the present time. The standardization of  
31 experimental design, the selection of informative biomarkers, and data analysis for genomics is

---

1 important for the utility of genomics information in future risk assessment and regulatory  
2 decisions. Such standardization will enhance the reproducibility of results obtained and the  
3 reliability of conclusions drawn from microarray data. Furthermore, EPA is considering the  
4 development of data quality standards based on performance of microarrays, as well as other  
5 genomics technologies (*e.g.*, functional genomics). This in turn will help to ensure the integrity  
6 of EPA's approach to assessing the genomics information submitted to the Agency.

7  
8 Genomics issues have already arisen in environmental decision-making. For example, a  
9 pesticide registrant has cited a published genomic article (Genter et al., 2002) as part of the data  
10 package submission for product registration to EPA's Office of Pesticide Programs. The data  
11 were submitted in support of an alternative mode of action that would affect human health  
12 assessment conclusions. Similar submissions are quite likely to be made by other pesticide  
13 registrants.

14  
15 Although this document focuses on the use of microarrays for toxicological studies as  
16 they pertain to macroorganisms, it should be noted that the impact of microarray technologies  
17 goes beyond the exploration of toxicological effects in eukaryotic systems. For example, the use  
18 of microarray techniques in environmental and clinical microbiology has increased significantly  
19 in the last few years. Microarrays can also be used to screen for host specific markers that can be  
20 used in microbial source tracking (MST). As an example of the application of genomics to MST,  
21 a research consortium including State of California regulatory agencies, public utilities, and EPA  
22 recently participated in a study comparing the performance of various genomics-based methods  
23 designed to identify the source of fecal material in ambient waters in an MST approach (Griffith  
24 *et al.*, 2003). Moreover, genomics methods are being evaluated to assist dischargers in  
25 complying with Clean Water Act (CWA) requirements to develop Total Maximum Daily Loads  
26 (TMDLs) for water bodies that are listed as impaired due to the presence of fecal coliforms. This  
27 MST work will also address the issue of beach closures; current microbial methods require  
28 several days to complete and do not distinguish between bacteria from humans and other sources  
29 such as sea gulls or seals. Further details on these MST efforts are described in *Microbial*  
30 *Source Tracking Guide Document* (available at:  
31 <http://www.epa.gov/ORD/NRMRL/pubs/600r05064/600r05064.htm>; U.S. EPA, 2005) .

1  
2           These examples indicate the need to make proactive policy decisions and to develop  
3 processes to address how genomics data will be used in Agency decision-making.  
4

#### 5 **1.4 Purpose and Intent of this Document**

6

7           As a result of research trends, EPA anticipates receiving increasing volumes of  
8 microarray data from environmental researchers, and as a part of the regulatory process. The  
9 Genomics Technical Workgroup considered all of the “omics” technologies and applications and  
10 decided that a guidance document on the use of data generated by DNA microarray analysis  
11 would be most beneficial to the Agency and regulated community at this time. This guidance  
12 applies to microarray data relevant to human health and ecological risk assessment and decision  
13 making. This guidance is provided in order to facilitate appropriate submission, consistent  
14 review, and optimal utilization of these data. Consequently, this document provides  
15 recommendations regarding: 1) data that should be considered for submission to the Agency for  
16 microarray studies, 2) the use of a performance approach to microarray quality assessment  
17 parameters, 3) data analysis approaches for microarrays, and 4) data management and storage  
18 issues for microarray data submitted to or used by the Agency.  
19

20           The purpose of this document is to provide information to the regulated community and  
21 other interested parties regarding submitting microarray data to the Agency and to provide  
22 guidance for reviewers in evaluating and utilizing such data and/or information. This interim  
23 guidance can be used by EPA program offices to determine the applicability of specific  
24 genomics information to the evaluation of chemical risks. It is important to note that microarray  
25 technology is rapidly changing, such that methodologies for generating such data and ensuring  
26 its quality will likely change; however the need to ensure consistency and quality in generating,  
27 analyzing and using the data will not. As the state of the science develops, EPA plans to revisit  
28 the guidance as necessary.

---

---

## 2.0 The Performance Approach to Quality Assurance for Microarrays

Quality issues are critical considerations in the application of new technologies or approaches, such as genomics. The Workgroup recommends that the Agency not prescribe specific methods to be used in microarray experiments at this time. This section instead provides general guidance on the recommended performance of microarray experiments in order to obtain data of the quality needed for a specific use.

The Agency acknowledges that continued advancement of tools and platforms for describing biological phenomena will be pivotal in supporting claims for regulatory decision making. It is also noted that at this time there exist numerous approaches, investigator fabricated and commercially available platforms, hardware and other peripheral equipment by which to measure biologic trends and changes at the level of tissues and cells. The following technical statements relate primarily to “*expression*” measurements (up- and down-regulation of macromolecules) and certain other multiplex technologies used to generate and collect quantitative and qualitative data about changing biologic conditions. This guidance is also relevant to the evolving nature of “*expression*” measures, particularly as recommendations for standardization in experimental performance put forth by the combined efforts of academic, industry and government scientists, become universally accepted and applied.

Although there are currently numerous means by which to observe and acquire biological expression measurements, such as *Massively Parallel Signature Sequencing* (MPSS) and *Serial Analysis of Gene Expression* (SAGE), the most frequently used experimental approach to collecting expression data is microarray-based studies. This technology, which has expanded well beyond the sphere of human health, is exploited to describe changing transcriptional profiles in genes of countless species that are important to numerous areas of biological sciences. Unfortunately, many of these investigations are undertaken without the benefit of explicit consensus for quality assurance and quality control and there has yet to be firmly established criteria for intra-experimental and cross-platform performance evaluation.

---

1  
2           Investigators submitting data to the Agency in support of regulatory decision-making,  
3 methods development, and technical transfer, should also consider at a minimum the  
4 performance-related experimental and system factors outlined in Appendix A, in addition to  
5 compliance with MIAME (Minimal Information About Microarray Experiments) Workgroup  
6 standards (<http://www.mged.org/Workgroups/MIAME/miame.html>) discussed in Section 3  
7 below. Further activities on the part of investigators to address experimental performance issues  
8 will serve to strengthen scientific arguments and experimental claims.

9  
10           Each EPA program, regional, or research and development office's Quality System  
11 should be defined and documented in their Quality Management Plan (QMP). A summary of  
12 their individual office's Quality System activities is detailed in a Quality Assurance Annual  
13 Report and Work Plan (QAARWP), which also includes information on their annual internal  
14 assessment of their Quality System.

15  
16           Additional detailed discussion of the EPA Quality System and the performance approach  
17 to quality assurance for microarrays is provided in Appendix A.

---

---

## 1 **3.0 Data Submission Guidance**

### 3 **3.1 Introduction**

4  
5 EPA developed the following information regarding submission of microarray data to  
6 facilitate appropriate review and consistent evaluation of data from multiple sources. The text  
7 that follows was written as a preliminary template guiding the submission of microarray data to  
8 the EPA. As the state of the science develops, EPA plans to revisit this submission format as  
9 necessary. In accordance with accepted practice, it is useful if submissions include sufficient  
10 information to allow an independent reviewer to reconstruct how the data were collected and  
11 analyzed. This approach allows reviewers to judge the quality of the data and the strength of any  
12 conclusions. It is also useful if the submission includes enough information in a format that  
13 facilitates comparison or integration with similar data from other experiments.

14  
15 Microarray technology is rapidly evolving with many competing platforms, native data  
16 formats, and analysis tools. As a result, a data submission standard should not be so specific as  
17 to stifle flexibility or innovation. Similarly, standards should not be burdensome, discouraging  
18 submission or slowing scientific progress. Many scientific journal editors grappling with these  
19 issues have adopted the Minimal Information About Microarray Experiments (MIAME)  
20 guidelines as a standard for submission of microarray data as part of a submitted publication  
21 (<http://www.mged.org/Workgroups/MIAME/miame.html>). A slightly modified version of  
22 MIAME, described below in Sections 3.2 through 3.7 and Appendix B, is proposed as the  
23 recommended microarray data submission template for EPA, which will be subject to change as  
24 the technology evolves. As genomics science and the associated technologies evolve, it can be  
25 expected that the MIAME guidance will concomitantly evolve. If the MIAME guidance in this  
26 document conflicts with the most recent changes to the MIAME guidance, the reader is directed  
27 to consider the MIAME guidance as the most recent, correct version.

### 29 **3.2 Abstract**

30

---

1 An abstract or executive summary of the source and type of data as well as the type of  
2 data evaluation and its final interpretation would provide a useful introduction to the data  
3 submission. Such a summary would not need to be exhaustive but would optimally provide the  
4 key highlights so that the reader will know the source of the data and how it was interpreted.  
5 The abstract might be written in a similar manner as for the submission to a scientific meeting or  
6 a journal article. It is advantageous if the reader is able to extract the important features of the  
7 submission and its interpretation from the abstract, although it is understood that a thorough  
8 evaluation of the substance of the data will involve a review of all the submitted material.  
9

### 10 **3.3 Experimental Design**

11  
12 It would be beneficial if voluntary submissions of genomics data to EPA included a  
13 sufficient description of the experimental design necessary to understand the source and nature  
14 of the data as well as the materials used to conduct the research. The following discussion is not  
15 an exhaustive listing or meant to be complete but indicates the spectrum of information on the  
16 experimental design that might be submitted for review. The submitter should consider  
17 providing the standard information one would include in the materials and methods section of  
18 any scientific article including a list of all the endpoints examined in the study. Such  
19 information would include information about the biological model system, treatment methods  
20 and doses, husbandry of animals, and cell culture information for *in vitro* systems. If whole  
21 animal models were employed, then submission of information regarding the exposure system,  
22 exposure doses, time points, details on euthanasia, length of time between harvesting of tissues  
23 and freezing or other processing, numbers of samples utilized for DNA array analysis, methods  
24 of RNA processing, and RNA quantification should be considered. The submitter should  
25 consider providing information on the methods employed for hybridization and incorporation of  
26 label and the numbers of hybridizations. When relevant, the submission of additional  
27 information necessary for interpretation of the data should be considered. Such information  
28 might include reference sample information, sample amplification, or any additional information  
29 unique to the study. The submitter should also consider providing information regarding any  
30 problems that arose during the study that could have an impact on interpretation.  
31

---



---

### 1   **3.4   Array Design**

2  
3           The inclusion of a complete description the platform used for transcriptional expression  
4 analysis such that the reviewer can assess the appropriateness of the analysis should be  
5 considered. The platform might be a commercially available platform (*e.g.*, Affymetrix, Agilent,  
6 Clontech) such that reference may be made to the specific type of chip used and the locations  
7 (weblink) of the source of the proprietary information so that the reviewer may access this  
8 information to aid in the review of the data analysis. If the transcriptional expression analysis  
9 was derived from a custom array designed for or by the submitter, then a inclusion of complete  
10 description of the production of the array would be useful. This information would likely  
11 include but certainly not be limited to the source of the nucleotide sequences used on the array,  
12 how the arrays were prepared, equipment used to prepare the arrays, description of the slides or  
13 membranes on which the arrays were spotted, gene lists, and any supportive data which confirms  
14 the specificity of the sequences used. A more complete listing of the types of data that would be  
15 useful in supporting the submission of custom arrays can be found in Appendix B.

### 16 17   **3.5   Biomaterials**

18  
19           It is advantageous if the submitted data package presents the physical characteristics of  
20 the studied biomaterials as these will likely vary between experiments. Such characteristics  
21 might include age, sex, cell type/line, and/or genetic variation. When applicable, this  
22 information would address the biological material from which nucleic acids (or proteins) have  
23 been extracted for subsequent labelling and hybridization. It is also recommended that submitted  
24 information on biomaterials detail the source properties, treatment, extract preparation, and  
25 labelling of the sample. Any pertinent information about sample controls would also be useful in  
26 analyzing submitted data.

27  
28           The exposure conditions applied to each test organism or tissue are important parameters  
29 influencing the experimental response. As a result, it is useful to document the incubation and  
30 treatment conditions applied to the studied biomaterial. Other key submission information might

---

1 include the method of chemical or physical exposure using the appropriate dosing units.  
2 Furthermore, any processing of samples taking place after exposure would be of interest.

3  
4 Information on the hybridization extract preparation protocol might include such details  
5 as the nucleic acid type and amplification method used. It would also be useful to record and  
6 submit the labeling materials and technique used in the experiment. Finally, the data submitter  
7 should consider outlining the type and position on the array of any external controls that may  
8 have been added to the hybridization extract(s). Please see Table B.4 in Appendix B for further  
9 information.

### 11 **3.6 Hybridization**

12  
13 It would be useful to submit a concise description of the procedures adopted for each  
14 hybridization. If a commercially available platform is utilized, reference may be made to the  
15 specific type of hybridization procedures and parameters adopted in the experiment. Web or  
16 literature citations describing the source of the hybridization protocol and materials are useful.  
17 Furthermore, information regarding the relationship between the labelled sample extracts and  
18 their corresponding arrays (design, batch and serial number) would be useful for understanding  
19 the experiment. Documentation of the steps taken in the hybridization including information  
20 regarding the solution, blocking agent and concentration used, wash procedure, quantity of  
21 labelled target used, time, concentration, volume, temperature, and a description of the  
22 hybridization instruments is encouraged.

### 24 **3.7 Measurements**

25  
26 The submitter should consider completely describing the methods used to acquire the  
27 image of the array, the nature of the image (*e.g.*, TIFF), the nature of the extraction of image data  
28 into quantified image data, and the nature of the spreadsheets used to house the quantified data.  
29 Submission of the original TIFF images is encouraged as is the submission of the initial  
30 quantization matrix. The description of the spreadsheet normalization of the TIFF data and any  
31 subsequent data analysis is also of value in a submission. In addition, features of the data used

---

1 for analysis such as background correction, normalization methods, methods used to test  
2 usability of the raw data, and types of analytical approaches would be useful information for the  
3 reviewer. Analytical approaches might include statistical models, graphical models, image based  
4 displays of data, and various analytical software packages. Information about the software may  
5 include weblink, proprietary information from instruction manual, or specific description of  
6 custom analytic methods. More complete description of information that should be considered  
7 for a submission for review may be found in Appendix B.

1

## 2 **4.0 Data Analysis Guidance**

3

4 This section provides information that will assist in regulatory and risk assessment efforts  
5 when considering the use of genomics data. Genomics data can be used to aid in reducing the  
6 level of uncertainty in the decision making process and provide a means to further evaluate  
7 exposure and effects. This guidance effort is also an attempt to highlight the need for developing  
8 genomics data analysis tool criteria, and the standardization of methods for the use of these tools.

9

### 10 **4.1 Introduction**

11

12 Evaluation of qualified genomics data, which have been properly analyzed and submitted  
13 (see Sections 2.0 and 3.0), has the potential to dramatically improve the mechanistic  
14 understanding of toxicities and their relevance to human health and ecological hazard  
15 identification and risk assessments. For example, DNA microarrays may be used to identify  
16 gene expression profiles associated with exposure to particular compounds, or characteristic of  
17 certain modes of action or mechanisms of toxicity. When a correlation has been established  
18 between a gene expression profile and a toxic mechanism, then these genomic data provide  
19 supportive evidence for that mechanism. Even when the mechanism for a particular compound  
20 is unknown, genomic data can help identify plausible toxicity pathways that may be involved in  
21 the biological process under study (Crosby *et al.*, 2000) for the purposes of prioritization or  
22 screening.

23

24 Genomic technologies generate vast amounts of data (gigabytes) quickly (during a single  
25 analytical session), especially when using DNA microarrays for gene expression profiling. This  
26 wealth of data increases the importance of careful documentation of experimental and analytical  
27 methods while working towards data interpretation and evaluation. The Minimal Information for  
28 the Analysis of Microarray Experiments (MIAME) guidelines have helped to standardize DNA  
29 microarray experiment documentation. Extension of the MIAME guidelines into  
30 toxicogenomics has provided even more applicable prerequisites for analysis

---

1 (<http://www.mged.org/MIAME1.1-DenverDraft.DOC>; Fostel et al., 2005). Also critical to  
2 analysis of genomics, and particularly microarray data, is access to the raw data from published  
3 or submitted experiments, and accompanying documentation of experimental and analysis  
4 details. Establishment of public genomic databases such as the Gene Expression Omnibus  
5 (GEO, <http://www.ncbi.nlm.nih.gov/geo/>) provides limited access to microarray data, but these  
6 are not compatible with all monitoring or regulatory applications.

7  
8 In addition to data submission and management activities, computational tools for  
9 genomics data analysis are another critical need for routine application of genomics data.  
10 Although evaluation of many of the currently available computational tools for genomics data  
11 analysis is underway through multiple internal and external Agency research efforts, these tools  
12 have not been examined by the Agency in sufficient detail that would allow for specific final  
13 recommendations to be made. Furthermore, while the variability and complexity of microarray  
14 experiments make prescribing a common, all-encompassing protocol functionally problematic,  
15 general components for the successful analysis and interpretation of all microarray approaches  
16 are discussed. The Agency is currently participating in several projects designed to develop  
17 appropriate protocols and methods for microarray data analysis. These include collaborative  
18 efforts with Food and Drug Administration (FDA) on the Microarray Quality Control project (  
19 <http://www.fda.gov/nctr/science/centers/toxicoinformatics/maq/>) and National Institute of  
20 Environmental Health Sciences (NIEHS) on the Chemical Effects in Biological Systems  
21 knowledgebase (<http://cebs.niehs.nih.gov/>). As an interim solution a genomics Data Evaluation  
22 Record (DER) template (Appendix C) is proposed as a means to outline a framework for  
23 genomics data analysis and documentation.

## 24 25 **4.2 Data Analysis**

26  
27 A few general features of genomic data analysis areas are described below with the intent  
28 to provide a basic but broad overview.

### 29 30 4.2.1 Data Processing and Filtering

31

1 Data processing covers the steps from scanning the array, to obtaining reliable estimates  
2 for the relative abundance of each gene transcript in all of the samples. Generally, these steps  
3 are classified as image analysis, quality control filtering, background correction, transformation  
4 and normalization. Each hybridized array has an associated and unique image file from which  
5 individual values (pixel intensities) can be collected. Data can be filtered to exclude signals that  
6 fail quality criteria. The specifics of data filtering and the threshold levels chosen are dependent  
7 upon the details and goals of the experiment. Standardization of processing and filtering criteria  
8 will be a critical step toward intra- and inter-laboratory agreement. The final output of the initial  
9 processing will be data that can be analyzed further to identify differentially-expressed genes.

#### 10 11 4.2.2 Statistics 12

13 A standard, or common, statistical approach, that would be appropriate for all microarray  
14 experiments, cannot be specified because of unique experimental variables such as differences in  
15 microarray platforms, experimental design (reference versus matched), levels of replication  
16 (technical versus biological), as well as within experiment sources of variation (spot to spot, slide  
17 to slide, etc.). Therefore, the types of methods and tools used for statistical analyses of  
18 microarray results often differ not only from more traditional experimental approaches, but also  
19 from one microarray experiment to another. Sample size strongly affects the statistical method  
20 chosen for analysis. For example, while a relative balance may exist between the number of  
21 samples and data points measured in a standard non-genomic experiment, microarrays, as well as  
22 proteomic and metabonomic technologies, generate hundreds and often thousands of data points  
23 from each sample. Furthermore, a variety of formulae exist to calculate appropriate microarray  
24 sample sizes, depending on experimental design. Nevertheless, the cost of conducting such  
25 experiments prohibits large scale studies with multiple sample sizes. Another constraint is  
26 sample pooling, at times a necessity due to the complex nature and paucity of biological material  
27 (*i.e.*, tissues and/or RNA quantities). It is, nonetheless, important to recognize that sample  
28 pooling may impact microarray experiments at multiple levels, including experimental design  
29 and subsequent analyses. Finally, data replication should be considered. It is important to  
30 distinguish the two types of replication that exist in biological experiments, including  
31 microarrays: technical (repeats of the same sample) and biological (starting material from unique

---

1 sources, such as different animals in a test group). For scientifically sound reasons, the latter  
2 assumes greater significance in most biological assays including microarray experiments.

### 3 4 4.2.3 Interpretation

5  
6 Numerous approaches can be used as a secondary level of analysis to interpret  
7 differentially expressed genes detected using microarray experiments. For example, genes can  
8 be sorted by ontology (gene ontology, GO) and subsequent cluster analyses (principal  
9 component analysis, hierarchical clustering, and  $\kappa$ -means clustering) can be used to better  
10 organize the data and help identify patterns of gene expression.

11  
12 Various bioinformatics (mathematical and statistical) algorithms can be used to integrate  
13 these patterns of expression with common biological pathways and networks of co-regulated  
14 genes. Linking these functional and pathway analyses to concurrent and previously identified  
15 phenotypic characteristics will significantly advance the understanding of the biological  
16 processes involved along the source-to-outcome continuum.

### 17 18 4.2.4 Inference

19  
20 Integration of these various data analyses and interpretation tools can be used to infer  
21 cause and effect relationships from these genomic data (Freeman, 2005). Biological inference  
22 may lead to biomarker development as well as descriptions of dose-response relationships,  
23 mechanisms of action, and predictive toxicity. Biomarkers are recognized as providing data  
24 linking exposure to internal dose and effect. The application of biomarkers to the risk assessment  
25 process that is linked to toxic processes or mechanisms may provide additional information for  
26 risk assessors. Additionally, data generated from microarray studies on model test organisms  
27 could be 1) applied to the identification of susceptible subpopulations, 2) used to develop  
28 surrogate species for toxicity testing, and 3) extrapolated to additional species, once the  
29 biomarkers and mechanism(s) of action are identified.

30

### 4.3 Data Evaluation

The goals of the evaluation of genomics data are directed toward risk assessment for regulatory applications. Currently, however, decisions cannot be made based solely upon gene expression pattern recognition, according to EPA's Interim Genomics Policy; this technology has not yet come to set precedence on its own. Currently, confirmatory studies are useful for potential risk assessment and regulatory use. If the data generated from microarray assays are confirmed using other techniques (*i.e.*, real-time quantitative PCR, functional enzyme assays, protein and metabolite profiles and/or linked to bioassay results), these data will help support links between gene expression, exposure and the resulting adverse effects in organisms. Furthermore, interpretation of microarray data with respect to existing toxicity profiles and endpoints of other perhaps higher level tests (clinical chemistry, immunochemistry, histopathology, and reproductive endpoints) should significantly increase the diagnostic and predictive applications of these technologies in the future.

A genomics Data Evaluation Record is used here as a way to present and organize data from genomics studies in order to derive information necessary for a regulatory application (see Appendix C for the Genomics Data Evaluation Record (DER) Template). For monitoring applications such information and standardization is recommended. The sections of the DER include the general information about a study and a brief executive summary as well as the materials and methods used. The test performance section includes: treatment and sampling times, tissues and cells examined, details of tissue harvest and storage, sample preparation, data analysis, evaluation criteria and statistical analysis. The results, discussions and conclusions are also components of the DER. Sections of the DER are included to provide example information to the risk assessor as a means to document the incorporation of genomics information in the risk assessment process. Genomic data used to support the more conventional data (*e.g.*, limited clastogenesis *in vitro* associated with cytotoxicity, DNA strand breaks, lipid peroxidation) are presented in an example DER for rats exposed to alachlor (see Appendix D: Draft Genomics Data Evaluation Record for Alachlor)



#### 1   **4.4   Data Analysis Conclusions**

2

3           The above considerations demonstrate that a systematic approach for genomics data  
4 evaluation is necessary for further use of genomic data in risk assessment efforts.

5 Documentation methods, like those in the proposed genomics DER (Appendix C) can help  
6 capture some requisite information, but the transfer of these evaluations, and the underlying  
7 genomics data, into searchable, electronic databases will be essential to making the data useful in  
8 risk assessments. Furthermore, development of databases containing gene expression profiles  
9 for a wide variety of chemicals should facilitate creation of statistical/computational methods  
10 that predict the toxic potential of a chemical.

---

## 5.0 Data Management

The goal of this section is to outline recommendations to EPA for an approach to managing genomic data submitted to the Agency or developed internally by EPA scientists. This includes the need to consider an Agency-wide warehouse for storage, retrieval and analysis of information submitted for regulatory or risk assessment purposes.

There are several major types of needs to consider in addressing the issue of an EPA-wide database: broad scientific needs for risk assessment purposes, program-specific regulatory needs, Agency Information Technology (IT) security needs, and public access needs. Although there is an overlap of issues for each of these purposes, it is useful to think of each additional purpose adding another layer of needs.

For scientific risk assessment purposes the key needs include the following items:

- 1) Standardization of data inputs as identified by the Data Submission Workgroup. This includes both microarray data and experiment parameters associated with the toxicogenomics study. It also provides for electronic submission of data.
  - 2) Provision of connectivity to external public biological databases such as Affymetrix, Agilent, and GenBank
  - 3) A quality control mechanism to ensure the fidelity of entered data
  - 4) Capability for importing and exporting data by means of automatic routines
  - 5) Inclusion of a wide range of data analysis and visualization tools such as filtering, clustering, and statistical analysis
  - 6) Sufficient scalability to address large data submissions, many users, and later addition of metabonomics and proteomics data at times in the future
  - 7) Audit trail capability. This would provide a time line and information on who added, changed or deleted specific data. It would also provide versions prior to deletions and changes.
  - 8) Automatic data back up and recovery system
-

1  
2 For security and management purposes, additional key needs include:

- 3  
4 1) Database hosting, administration and management. This includes managing data  
5 submission, database access and privileges, software and hardware updates, back-up and  
6 storage.  
7 2) Physical and electronic security, including user authentication, firewalls, and virus  
8 protection.  
9 3) Governance structure to provide policies and procedures for submissions, access,  
10 security, cost sharing, and priority for development of new features.  
11

12 For regulatory purposes, additional considerations may be necessary:

- 13  
14 1) Electronic signature or other formal identity management capability. If the data are  
15 submitted electronically as part of a regulatory submission, the system needs to ensure  
16 that the submission is linked to the submitter.  
17 2) Capability of partitioning the database to secure Confidential Business Information (CBI)  
18 or other non-public information, if this is part of a regulatory submission.  
19 3) Workflow enabled, so that reviewer can address data in systematic steps needed for  
20 response to submission.  
21

22 For public access purposes, key needs are:

- 23  
24 1) Database is Web enabled, with easy routine for export of data.  
25 2) Clear policies governing the management of the public database as opposed to an internal  
26 or staging database.  
27

28 There may also be staging considerations in building or adopting an Agency-wide  
29 genomic database. The first phase might include genomic data only, and have limited analytic  
30 capability. Eventually the database should provide quality assessment tools, extensive analytical  
31 capability, gene-centric queries, and encompass proteomic, metabonomic, and conventional

1 toxicology assay results. Integrating these diverse types of experimental data will support data  
2 mining as well as the development of predictive toxicology systems.

3  
4 Currently, there is no single database at EPA for managing genomics data each program  
5 or lab is developing its own approach. As the needs are currently identified above, there are  
6 several advantages to creating and maintaining an EPA-wide genomics database:

- 7
- 8 1) **Cost.** All of the scientific and management/security needs identified above should be  
9 addressed by any genomic database used at EPA. Addressing these items once in a  
10 uniform way would avoid duplication of these costs.
  - 11 2) **Data Access.** All users in the Agency would have access to all Agency genomic data  
12 (except CBI data), greatly enhancing our risk assessment capabilities.
  - 13 3) **Quality Control and Consistency.** A quality control mechanism would ensure that all  
14 Agency data passes a consistency test.
  - 15 4) **Availability of a Common Set of Analysis Tools.** As new tools are developed, they  
16 would become available to all users.
  - 17 5) **Scaleable.** While lab or program specific databases may focus on a narrow range of data  
18 or analysis, an EPA database would be built to include a wider range of “omics” data and  
19 a full portfolio of analytical tools enabling Agency scientists to pursue a wider range of  
20 data mining and biological systems-oriented studies.

21

---

## 6.0 Additional Recommendations

The Genomics Workgroup recommends that the Agency undertake a number of follow-up activities to this interim guidance including: 1) further development of the training materials and modules outlined below, to be offered throughout the Agency to risk assessors and decision makers who will be faced with the challenge of interpreting and applying genomics information, 2) continued collaboration of EPA personnel with staff from other federal agencies and stakeholders in the development of tools for the analysis of genomics data, 3) application of this guidance to a series of case studies to evaluate its utility in risk assessment and regulatory applications; and 4) the updating of this guidance as needed as the technology evolves.

### 6.1 Training Needs and Recommendations

The charge to the Genomics Task Force Training Workgroup was to develop an approach and appropriate delivery mechanisms for training Agency risk assessors and managers to understand and interpret genomics data in the context of risk assessment. The need for a better understanding of molecular biology concepts, and ultimately how genomics, proteomics, and other “omics” data may be used to support decision making, is the primary driver for the development of such training for staff and managers.

In designing training genomics, the Training Workgroup considered several issues: 1) the need to develop a modular approach that could build on basic information and change as new information becomes available, 2) the need to vary the level of complexity based on the needs of a particular audience, 3) the importance of considering the target audience, based on the recognition that different staff and managers will have different needs, 4) the need to develop a schedule for production of training materials, recognizing that, by taking advantage of existing public sector resources to build the initial version of the Genomics Training, time and resources may be saved; and 5) identifying internal capacity to provide training, such as ORD scientists and risk assessors to save time and resources.

1           Presented below is a draft outline that describes a modular training course in molecular  
2 techniques, in general, and genomics data interpretation, in particular. The Genomics Training  
3 would consist of three levels of training targeted to specific audiences, each consisting of a series  
4 of modules devoted to a particular group of concepts and/or techniques. Each training level is  
5 outlined below in Table 1, with descriptions of the content and instructional goals. More detail  
6 on the proposed training is provided in Appendix G.

**Table 1. Overview of Genomics Training Plan** (see Appendix G for more detail)

<b>Training Level</b>	<b>Number of Module</b>	<b>Target Audiences</b>	<b>Content</b>	<b>Goal</b>
Level I: Introductory Modules-	8	Non-scientists and/or technical staff without training in biological sciences.	Molecular Biology concepts: cell structure and function, DNA, RNA, proteins, gene arrays, risk assessment concepts, regulatory and risk assessment communication, EPA's current genomics policy.	Provide basic information necessary for understanding assessments of cellular functions at the molecular level and how genomics data may affect risk assessments.
Level II: Intermediate Modules	3	Scientists and/or those likely to use genomics: Intended for staff who need more in-depth understanding of genomics data generation, but do not necessarily generate data.	Background on molecular techniques such as microarrays, DNA amplification techniques, DNA fingerprinting, protein analysis, etc. Modules to be targeted for specific applications. ( <i>e.g.</i> , microbial source tracking, homeland security, field inspectors, etc.)	Provide a general understanding of various applications that may be currently considered by programs throughout EPA. Intended to support human health and ecological risk assessors.
Level III Advanced Modules	Dependent on specific technical needs.	Scientists and those likely to use genomics data to generate risk assessments.	Modules would include statistical, computational and bioinformatics approaches to analyze genomic data, the use of molecular biology in mode-of-action determinations, and using genomics data in hazard/risk assessments. Flexible to account for changes in the field and to meet needs of the different EPA programs. As new technologies/ applications appear, additional modules developed, enhanced and/or revised.	Provide advanced-level knowledge on specific technical needs that scientists performing research or developing hazard/risk assessments associated with chemical registrations and other regulatory activities may face.

## 6.2 Collaborative Development of Genomic Tools for Data Analysis and Data Management

The Agency, in concert with other federal agencies, has begun to investigate and evaluate the currently available computational tools for genomic data analysis. EPA has been testing the toxicogenomic data management and analysis features of the NIEHS Chemical Effects in Biological Systems (CEBS) knowledgebase and FDA National Center for Toxicological Research's ArrayTrack database. EPA has also been collaborating with FDA, National Institutes of Health (NIH), National Institute of Standards and Technology (NIST), and other stakeholders on the microarray quality control (MAQC) project to establish protocols for genomic data analysis. Further, EPA has participated in National Academy of Sciences (NAS) workshops and International Life Sciences Institute (ILSI) projects on the application of genomics to toxicology and risk assessment. Building on these prior efforts, recommendations on the use of genomics tools should be identified recognizing that the goal is the appropriate application of genomic data in risk assessments and regulatory decision making. The Agency should also consider and identify limitations of the currently available tools. Ultimately, the Agency is looking for quantitative and predictive modeling tools, which will likely call for the development of new algorithms and models. These tools will need to provide reliable and repeatable data analyses, and the consistent and necessary information for EPA decision making processes. The scientific, mathematical, and statistical methods that are used for these models and analyses will need to be validated and standardized.

Due to the potentially large volumes of genomic and associated toxicological data, it is essential that the Agency consider the development of a complete data management solution. The functional needs of a solution of this magnitude should minimally include items listed in Section 5.0 Data Management. In addition, this data management solution should address needs unique to scientifically-based risk assessments, confidential and proprietary data security, public access, and other aspects of regulatory application. It should be noted that consistency, scientific and operational robustness, common access, and availability in a scalable environment are important data management needs. While the Agency has begun to develop bioinformatics

---



1 research efforts, both intramurally (e.g., ORD's National Center for Computational Toxicology)  
2 and extramurally (the STAR funded Environmental Bioinformatics Center in NC and NJ), an  
3 Agency-wide data management solution integrating genomics, toxicological, and other key data  
4 for regulatory applications is now needed.

5  
6 **6.3 Applying this Interim Guidance for Microarray-Based Assays to Case Studies to**  
7 **Verify its Utility in Risk Assessment and Regulatory Applications**

8  
9 The EPA's Risk Assessment Forum and other appropriate groups should apply this  
10 interim guidance to several case studies to evaluate its utility in risk assessment and regulatory  
11 applications and to identify potential areas for improvement.

12  
13 **6.4 Updating Genomics Guidance as Needed**

14  
15 This interim guidance should be revised and updated as indicated through its application  
16 to case studies (see section 6.3 above), and as genomics technologies evolve. Additional  
17 genomics guidances (*e.g.*, proteomics, metabonomics) should be developed as needed to ensure  
18 the Agency is prepared to receive and apply such data as the need develops.

19

---

## 1   **References**

- 2  
3   American Public Health Association, American Water Works Association, & Water  
4   Environment Federation. Standard Methods for the Examination of Water and Wastewater.  
5   Revision in process.  
6
- 7   Brooks, A.N., Pennie, W.D. 2001. Transcript profiling of the response to environmental hormone  
8   mimics. *Comments Toxicol* 7:303-315.  
9
- 10   Burczynski, M.E., McMillian, M., Ciervo, J., Li, L., Parker, J.B., Dunn, R.T., Hicken, S., *et al.*  
11   2000. Toxicogenomics-based discrimination of toxic mechanism in HepG2 human hepatoma  
12   cells. *Toxicol Sci* 58:399-415.  
13
- 14   Crosby, L.M., Hyder, K.S., DeAngelo, A.R., Kepler, T.B., Gaskill, R., Benavides, G.R., *et al.*  
15   2000. Morphologic analysis correlates with gene expression changes in cultured F344 rat  
16   mesothelial cells. *Toxicol Appl Pharmacol* 189:205-222.  
17
- 18   Fostel, J., Choi, D., Zwick, C., Morrison, N., Rashid, A., Hasan, A., Bao, W., *et al.* 2005.  
19   Chemical Effects in Biological Systems—Data Dictionary (CEBS-DD): A Compendium of  
20   Terms for the Capture and Integration of Biological Study Design Description, Conventional  
21   Phenotypes, and ‘Omics Data. *Toxicol Sci* 88:585–601.  
22
- 23   Freeman, M.R., Cinar, B., and Lu, M.L. 2005. Membrane rafts as potential sites of nongenomic  
24   hormonal signaling in prostate cancer. *Trends Endocrinol Metab* 16(6):273-9  
25   (Freeman, 2005  
26
- 27   Genter, M.B., Burman, D.M., Soundarapandian, V., Ebert, C.L., Aronow, B.J. 2002. Genomic  
28   analysis of alachlor-induced oncogenesis in rat olfactory mucosa. *Physiol. Genomics* 12:35-45.  
29
- 30   Hamadeh, H.K, Bushel, P.R., Jayadev, S., Martin K., DiSorbo O., Sieber S., *et. al.* 2002. Gene  
31   expression analysis reveals chemical-specific profiles. *Toxicol Sci* 67:219-231.  
32
- 33   Moreau, Y., Aerts, S., De Moor, B., De Strooper, B., Dabrowski, M. 2003. Comparison and  
34   meta-analysis of microarray data: from the bench to the computer desk. *Trends Genetics*  
35   19:570-577.  
36
- 37   U.S. Environmental Protection Agency. 2005. Microbial Source Tracking Guide Document.  
38   Office of Research and Development, Washington, DC EPA-600/R-05/064. 131 pp.  
39   <http://www.epa.gov/ORD/NRMRL/pubs/600r05064/600r05064.htm>  
40
- 41   U.S. Environmental Protection Agency, Science Policy Council. 2004. Potential Implications of  
42   Genomics for Regulatory and Risk Assessment Applications at EPA. EPA 100/B-04/002.  
43   available at: [www.epa.gov/osa/genomics.htm](http://www.epa.gov/osa/genomics.htm)  
44
-

- 1 U.S. Environmental Protection Agency, Science Policy Council. 2002a. Interim Policy on  
2 Genomics  
3
- 4 U.S. Environmental Protection Agency. 2002b. NELAC Constitution, Bylaws and Standards  
5 EPA/600/R-03/049.  
6
- 7 U.S. Environmental Protection Agency. 1997. "Performance Based Measurement System," 62  
8 Federal Register 52098 – 52100, October 6, 1997.  
9
- 10 Waring, J.F., Ciurlionis, R., Jolly, R.A., Heindel, M., Ulrich, R.G. 2001. Microarray analysis of  
11 hepatotoxins in vitro reveals a correlation between gene expression profiles and mechanisms of  
12 toxicity. *Toxicol Lett* 120:359-368.
-

## 1 **Appendix A: EPA Quality System and the Performance Approach** 2 **to Quality Measurement Systems**

3  
4  
5 The best quality data may not be technically available, affordable, or even applicable to  
6 the exact problem at hand. To address a variety of circumstances, EPA has developed a Quality  
7 System by which reasonable quality assurance (QA) guidelines or policies are offered for  
8 assuring, documenting, and assessing data quality. EPA's Quality System is defined in *EPA*  
9 *Order 5360.1 A2, Policy and Program Requirements for the Mandatory Agency-Wide Quality*  
10 *System*, the *EPA Quality Manual for Environmental Programs, EPA Manual 5360 A1*, the  
11 *Contracts Management Manual*, and the Agency's Website ([www.epa.gov/quality](http://www.epa.gov/quality)). The  
12 requirements for EPA-funded organizations and organizations submitting data to EPA under  
13 applicable statutes and regulations are also found in the Code of Federal Regulations (48 CFR  
14 Part 46), also available through [www.epa.gov/quality](http://www.epa.gov/quality). Parties submitting data under applicable  
15 statutes and regulations are expected to document the quality of the data submitted as well as  
16 how it was achieved. Quality System parameters apply to environmental data operations and  
17 measurements or information that describe: (1) environmental processes, (2) location or  
18 conditions, (3) ecological or health effects and consequences; and (4) performance of  
19 environmental technology

### 20 21 **What is a Quality System?**

22  
23 As illustrated in Figure 1, a Quality System is viewed as a tiered organizational approach  
24 for its work processes because it defines how the work is conducted, and provides a scientific  
25 and technical basis for EPA's decision making process. The Quality System is a documented  
26 management structure to ensure the quality of an organization's work processes, products and  
27 services. Adhering to the Quality System helps to ensure that all operations, no matter where  
28 they are performed, occur in a consistent manner and that the processes and outputs in the system  
29 are effective, stable, and consistently followed. Key components in a Quality System are: (1)  
30 Quality management, (2) Quality assurance (QA), and (3) Quality control (QC).

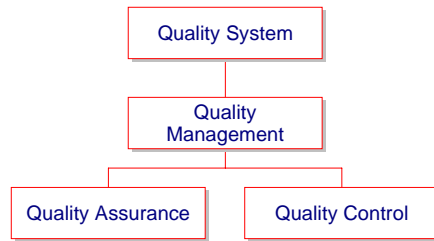


Figure 1. A Generic Quality System

### What Documentation is Needed for Organizations Submitting “Genomics Data??

An organization documents its Quality System in a Quality Management Plan while a laboratory may document its implementation of specific quality policies and practices in a document entitled a Quality Manual or Quality Assurance Plan. However named, the document details the efforts to produce data that are adequate for their intended use and for assuring conformity with regulations and customer requirements for data quality. Examples of a Quality Management Plan are available at [www.epa.gov/quality/qmps.html](http://www.epa.gov/quality/qmps.html).

### What is a Performance Approach?

A Performance Approach conveys “what” needs to be accomplished, but not prescriptively “how” to do it. EPA defines the performance approach as a set of processes wherein the data needs, mandates, or limitations of a program or project are specified, and serve as criteria for selecting appropriate methods to meet those needs in a cost-effective manner. The criteria may be published in regulations, technical guidance documents, permits, work plans, or enforcement orders. Under a performance approach, EPA would specify the questions to be answered, the decisions to be supported by the data, the level of uncertainty acceptable for making decisions, and the documentation to be generated to support this approach (see

1 <http://www.epa.gov/fedrgstr/EPA-WASTE/1997/October/Day-06/f26443.htm>, or 62 FR 52098  
2 for more details about Agency policy regarding the performance approach)

3  
4 Performance approaches can be defined as either: (1) measurement data that are of  
5 specified quality when demonstrating compliance (measurement quality objective (MQO)  
6 approach), or (2) a demonstration of compliance that achieve specified statistical confidence (the  
7 data quality objective (DQO) approach). Any appropriate measurement technology and  
8 sampling frequency/thoroughness may be used as long as MQO or DQO is documented and met.

9  
10 Key components that need to be considered in a performance approach are:

- 11
- 12 a) Sampling procedures and sample acceptance criteria, describing procedures for  
13 collecting, handling (*e.g.*, time and temperature), accepting, and tracking submitted  
14 samples, and procedures for chain-of-custody.
  - 15
  - 16 b) Analytical methods, listing the laboratory's scope for testing and denoting  
17 accreditation/certification status for individual methods, for non-standard methods or new  
18 methods, the laboratory's validation procedures.
  - 19
  - 20 c) Analytical quality control measures, stating the laboratory's requirements for  
21 measurement assurance, *e.g.*, method verification and documentation, error prevention,  
22 and analytical checks such as duplicate analyses, blanks, positive and negative culture  
23 controls, sterility checks, and verification tests.
  - 24
  - 25 d) Documentation control and record keeping specifications, identifying recordkeeping  
26 procedures to ensure data review, acquisition, traceability; accountability noting  
27 procedures to ensure customer confidentiality; and other parameters such as control,  
28 security, storage, retention, and disposal of laboratory records.
  - 29
  - 30 e) Assessments, describing the laboratory's processes to monitor the effectiveness of its  
31 QA program.
-

1  
2 1) Internal audits of laboratory operations, performed on a routine basis,  
3 minimally annually, by the QA officer and supervisor. For a small laboratory, an  
4 outside expert may be needed.

5  
6 2) On-site evaluations by outside experts to ensure that the laboratory and its  
7 personnel are following an acceptable QA program.

8  
9 3) Proficiency test studies, in which the laboratory participates. These  
10 collaborative studies confirm the abilities of a laboratory to generate acceptable  
11 data comparable to those of other laboratories and to identify potential problems.

12  
13 f) Correction and preventative activities, identifying procedures used to determine the  
14 causes of identified problems and to record, correct, and prevent their re-occurrence.

### 15 16 **Systematic Project or Experimental Planning**

17  
18 In general, systematic project planning is essential before any activity begins, whether it's  
19 sampling or analysis. For any project, the scientist needs to develop the experimental study  
20 design by first identifying and documenting what the problem is, why the new information is  
21 needed, and the objectives for the experiment or series of experiments.

22  
23 Once the study objectives are defined, the hypothesis is then developed. In  
24 systematically planning a project, the team or researcher then needs to determine the study  
25 parameters or test variables, both critical and the secondary (if any). The data quality objectives  
26 or performance criteria (*i.e.*, how good the data should be for the intended purpose) should be  
27 defined before the experiment starts along with all the appropriate quality control activities. For  
28 example, how types and numbers of replicates will be followed in the experiment, how is the  
29 specificity/selectivity of the analytical method to the target determined, how will the precision be  
30 determined in terms of repeatability and reproducibility. In the process of determining all these  
31 quality control activities, the experimental design can be optimized and documented.

---

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31

## Parameters of Microarray Platform Performance

Most microarray gene expression experiments fall into three broad classes, depending on outcome, that should have distinct QC reporting needs:

1. The first class of microarray experiments is that for which the investigator concludes that a treatment/exposure causes a biological effect. This is the most common conclusion from published microarray experiments, and the simplest from a QC point of view.
2. The next category of experiments is one wherein the investigator concludes that the treatment has no observable effect. These results are rarely reported in the literature, but might be common in regulatory submissions. These are somewhat more complicated from a QC perspective.
3. The third group offers claims about the magnitude of changes in transcription. Examples of this last class of experiments are rare, and are the most difficult and expensive on which to perform adequate QC. Currently, a cost effective microarray platform on which to perform this last class of experiments is not available. While the minimum QC of the experiment may be unchanged, the extent of documentation needed to verify that an ensuing experimental report is acceptable, may vary based on accompanying results.

Although *negative* and *positive* controls are part of experimental designs and investigative approaches, the requisite controls for “expression” studies – particularly microarray experiments - are not always obvious. Investigators are encouraged to consider not only the biological system under scrutiny, but also the nature of the assertions about the system. The need to have adequate controls should be considered in an experimental scheme, in order to demonstrate that measurements are accurate enough to support scientific claims and assertions. Needs for several straightforward situations are listed below, and can be applied as a guide to more complex scenarios. It is useful if control samples are constructed in a way that ensures the control and experimental samples are as similar to one another as possible (*e.g.*, with regard to

---



1 biologic composition and complexity of RNA) except in such cases where control sample  
2 characteristics are unambiguously presumed to differ.

3  
4 In cases where the investigator proposes that a biologic effect is present (the first case  
5 noted above), the primary QC issues are precision and specificity, and the use of a negative  
6 control is encouraged demonstrate that the measurement system is not likely to produce false  
7 positives. Accuracy is rarely a concern, since claims are not being made regarding the  
8 magnitude of differential expression between experimental and control groups, but only whether  
9 a difference exists. Sensitivity is also not relevant, since no difference or effect would be  
10 observed if sensitivity were too low. It is useful if the negative control and experimental groups  
11 include sufficient replicates, relative to the magnitude of effect and experimental variability, in  
12 order to show that the claimed effect, and no effect cases, can be statistically distinguished with  
13 desired confidence. RNA from untreated samples is usually as adequate and readily available as  
14 a control in this case. Precision is accounted for by statistical procedures (*e.g.*, t-test, chi squared  
15 test, and respective non-parametric analogs) routinely used to determine whether the  
16 experimental and control groups differ significantly. Additional consideration should be given to  
17 demonstrating the specificity of measurement(s) for the effect of interest. Demonstrating that a  
18 variety of probes exhibit binding affinity for discrete regions of a given transcript, can provide  
19 congruent results and is often a sound way to address the issue of platform specificity. The  
20 ability of probes to distinguish similar transcripts, including splice variants, is also useful to  
21 address. The issue of specificity is best addressed by using complementary “expression”  
22 measurement technologies (*e.g.*, quantitative real-time PCR, Northern blot analysis, RNase  
23 protection assays, S1 nuclease protection) to confirm microarray results. This will control for  
24 technique-specific effects, and by using distinct set(s) of amplification primers, help control for  
25 non-specific or unintended hybridization to microarray probes. Alternatively, a different  
26 microarray platform could be used to confirm specificity if the second microarray platform uses  
27 distinct probe sequences for detecting the transcripts of interest. A useful way to control  
28 systematic error is to ensure randomization of both processing order and acquisition of  
29 measurements for control and experimental samples. Using blind samples can be a useful  
30 approach to avoid operator bias.

31

---

1           In cases where the investigator maintains that there is no biologic effect (class two,  
2 above), a positive control is useful to show that the measurement system is capable of detecting  
3 the smallest effect sizes for which the claim is being made. In scenarios such as this, the  
4 additional QC factor of sensitivity comes into play, while specificity becomes less important. It  
5 is advantageous if the positive control and experimental groups both contain a sufficient number  
6 of replicates to show that the two groups can be statistically distinguished with the desired  
7 confidence level. It is useful to avoid absolute claims to the effect of ‘*there is no effect*  
8 *whatsoever,*’ since only effects equal to or larger than those readily observable in the positive  
9 control can realistically be ruled out. Instead, conclusions might take the form of ‘*there is no*  
10 *effect larger than X*’, where ‘X’ represents the smallest magnitude of effect readily detectable in  
11 positive controls. Some validated positive controls, such as samples subjected to a treatment  
12 widely recognized to produce the desired effect, are considered the preferred source for positive  
13 controls. However, there are cases for which no adequate model exists for the effect being  
14 studied. Then, it is useful to construct a positive control using methods such as spiking complex  
15 RNA samples with purified and quantified RNA of interest. Alternatively, investigators might  
16 use mixtures of complex RNA samples in which the RNA of interest is present in varying known  
17 concentrations (see also section on **System Linearity and Calibration**). These controls are  
18 useful for demonstrating that the measurement system can readily detect the effect sizes for  
19 which the negative claim is being asserted.

20  
21           As always, it is beneficial to randomize the order of processing and measurement relative  
22 to the sample group, and using blind samples should be considered. The same statistics, as  
23 applied to the ‘*there is an effect*’ case, are generally used to control for inconsistencies in  
24 precision, but in this case acceptable performance means that the positive and negative controls  
25 may be reliably distinguished from one another, while the experimental sample is statistically  
26 indistinguishable from (appears to come from the same population as) the negative control.

27  
28           When an investigator submits a claim regarding the magnitude of an effect, and not only  
29 the presence or absence of effects, a more complex system of control (i.e., calibration curve)  
30 should be considered (see also section on **System Linearity and Calibration**). In cases such as  
31 this, where quantification of differential expression is critical (e.g., when stating that

---

1 transcription of *gene X* increases 1.8 fold after exposure/treatment), accuracy becomes the  
2 foremost QC factor, and more complex positive controls and statistics should be considered. A  
3 calibration curve typically demonstrates the accuracy of the measurement system across the  
4 range of concentrations being considered. Researchers should consider assembling appropriate  
5 materials for constructing calibration curves in cases where standard reference RNA is not  
6 available. In many cases, investigators may consider methods such as spiking a complex RNA  
7 sample with a known series of concentrations of the RNA species of interest, or using a mixture  
8 series of two complex samples where the concentration of the RNA species of interest differs by  
9 a known amount. In the latter case, combining the two RNA ‘targets’ at different ratios produces  
10 a series of known concentrations of the RNA species of interest. It is useful to adjust the range  
11 and spacing of concentrations on the calibration curve (*e.g.*, log linear scale) and the number of  
12 replicates per concentration based on the level of precision desired and amount of experimental  
13 variability observed. Specificity of signals of interest might be confirmed by showing  
14 congruence with signals produced by probes that hybridize to a different portion of the same  
15 transcript. It is beneficial if conclusions on the magnitude of an effect include confidence  
16 intervals that reflect the performance of the measurement system during calibration curve  
17 construction, as well as variability seen in the experimental samples.

18

### 19 **Overview of Array Technology – The Physical Platform**

20

21 The current method for fabricating DNA microarrays (DNA chips) is to use either cDNA  
22 or oligonucleotides as probes that represent specific genes in the organism of interest, attached to  
23 a suitable solid substrate such as a glass microscope slide. It is acknowledged that the specificity  
24 of these probes is limited by the current understanding of gene sequence, among other things. It  
25 is useful if all sequences are periodically reevaluated based on the newest gene sequence  
26 information to ensure valid assessments.

27

28 Microarrays populated with cDNA probes are created by ‘spotting’ amplified cDNA  
29 fragments in a desired density pattern onto a solid medium such as a glass slide. Arrays using  
30 oligonucleotide probes are either mechanically ‘spotted’ or assembled by chemically  
31 synthesizing short, unique oligonucleotide probes directly onto a glass or silicon surface using

1 covalent chemistry or photolithographic technologies. It has been well established that  
2 numerous possibilities exist for errors to become ‘fixed’ during the manufacture of the arrays;  
3 therefore, the fidelity of the DNA fragments immobilized to microarray surfaces may be  
4 compromised by several different kinds of experimental and manufacturing inconsistencies.  
5 Given the QA/QC challenges in manufacturing gene arrays, a trend has emerged in recent years  
6 towards the use of gene arrays from several large vendors rather than arrays from smaller scale  
7 manufacturers or those prepared “in-house.” While this may limit choice, it may also offer an  
8 advantage to the array community when addressing issues of cross-platform compatibility.  
9 There are a number of sources of technical error which can adversely impact data quality of a  
10 gene array experiment. These include, but are not limited to, poorly functioning probes or probe  
11 sets, cross hybridization of related genomic sequences, scanner settings and function, and  
12 atmospheric ozone. Unfortunately, a set of performance standards by which individual  
13 laboratories may be evaluated are not currently in place, although it is anticipated that such  
14 standards may be developed in the near future.

15  
16 In many array-based studies, the investigators report microarray data for which there is  
17 no corroborating validation for the observed transcriptional measures. For profile data observed  
18 on array platforms regarding novel findings that are not readily supported in the peer-reviewed  
19 scientific literature, it is useful to include supporting data generated by traditional methods of  
20 evaluating gene expression, such as PCR, Northern blot hybridization analysis and RNase  
21 protection assays. In addition, the quality of probe sequences selected for particular transcribed  
22 regions incorporated onto the array is also a critically important consideration. For example, if  
23 probes are selected primarily from the 3’ end of given genes, splice variants of those genes can  
24 evade identification, if the alternative splicing events occur 5’ of a probe region. Additionally,  
25 by microarray analysis, it is very difficult to distinguish between two expressed genes that share  
26 a high degree of sequence homology. Variation in probe specificity is also a commonly  
27 encountered problem in oligonucleotide arrays. This problem frequently arises in instances  
28 where nucleic acid sequences are practically identical between two coding regions and the  
29 oligonucleotide probes are synthesized from 3’ends of the genes.

30

---

---

## 1 Isolation of Nucleic Acid ‘Targets’

2  
3 Since this biological analyte comprises the molecular species that will be measured, it is  
4 beneficial to ensure efficient isolation as well as post-isolation stability and structural integrity.  
5 Total RNA is generally used for gene expression analysis, although, mRNA is also used. RNA  
6 isolation techniques often involve homogenization of either fresh or frozen samples at high  
7 concentrations of guanidine isothiocyanate followed by phenol extraction and alcohol  
8 precipitation, although other methods can produce RNA of high quality.  
9

10 Methods for determining purity (*i.e.*, absence of contaminating reagents) include nucleic  
11 acid analysis by spectrophotometry at absorbance ratios of 260/280 nm, with expected values  
12 between 1.90 and 2.10 at pH 7.5. Another conventional method for determining structural  
13 integrity is through the use of MOPS-EDTA formaldehyde (or glyoxal) agarose gel  
14 electrophoresis, during which either the integrity of ribosomal RNA or the relative size  
15 distribution of mRNA can be evaluated. Recent advances in microfluidics and analytical  
16 equipment, (*e.g.*, 2100 Bioanalyzer , Agilent, Inc.), allow investigators to evaluate the integrity  
17 of nucleic acids with greater speed and accuracy than possible with agarose gel electrophoresis.  
18 It is anticipated that this technology will soon replace the more frequently used methods.  
19

## 20 Experimental Design

21  
22 The importance of pre-planning the experimental design cannot be overemphasized. Since  
23 the critical outputs from biological “expression” analyses are largely dependent upon  
24 experimental design investigators should consider devoting extensive attention to performing  
25 experiments with the appropriate design parameters. It is advantageous if the chosen  
26 experimental design provides sufficient statistical power to unambiguously test the biologic  
27 argument. The level of analytical power needed to allow for the detection of differentially  
28 expressed transcripts at a ratio greater or equal to ‘X-fold’ should be considered. In addition, it is  
29 useful if such analyses take into account the percentage of false positives that the researcher is  
30 willing to accept. The false positive rate (FPR) and the false negative rate (FNR) are necessarily  
31 dependent upon each other *i.e.*, a decrease in one results in an increase in the other.

---

1  
2           When designing an experiment, adequate consideration should be given to sample sizes,  
3 the use of controls, the use of sample randomization, and blind sample procedures. Specific  
4 needs will depend on a number of factors, including the nature of the conclusion being presented,  
5 the manner by which samples are compared with one another, the range of measured effects, and  
6 experimental precision. To ensure adequate statistical power will be realized to support  
7 scientific arguments and conclusions, it is advantageous to consult a statistician during the  
8 experimental design phase. In order to estimate the projected magnitude of effect and  
9 experimental precision conducting a small scale pilot experiment in advance of the definitive  
10 experiment might be considered. Alternatively, one technology (*e.g.*, DNA microarrays) might  
11 be used as an exploratory tool for hypothesis generation, followed by the use of a secondary  
12 technology (*e.g.*, quantitative RT-PCR, qPCR) to generate adequate numbers of experimental  
13 and control replicates to fulfill hypothesis testing. Some general issues that should be considered  
14 are listed below, but their specific applicability will vary among experiments. Careful  
15 consideration of these issues should provide sufficient information for a reliable estimate of  
16 overall experimental performance, and the statistical strength of conclusions put forth.

17  
18           It is expected that technical variation will be introduced at each critical laboratory step  
19 during expression analysis. In addition, unique sources of variation are likely to be associated  
20 with individual laboratories and/or technicians. It is important, therefore, that this variation be  
21 considered during study design and statistical analysis in order to avoid confounding of these  
22 sources of variation with treatment effects. For those experiments in which data are collected  
23 from array-based studies, there are three design schemes typically used; these are briefly  
24 described below. Although these are certainly not all inclusive, identification of the acceptable  
25 system is left to individual research teams. The three fundamental design alternatives typically  
26 used are 1) the flexible *universal reference design*, which is used for analysis of many  
27 experimental factors of equal importance, or those that will be integral to future meta-analysis, 2)  
28 the efficient *balanced block design*, for use in looking for genes that are upregulated or  
29 downregulated between two samples, and 3) the more integral *loop design*, which when  
30 comparing samples of equal interest and high quality results in half the variance per estimate,  
31 because each sample is included two times, rather than once, at the minimal expense of one

---

---

1 additional chip. There is, however, a rather large experimental cost of this latter design, because  
2 it relies on not even one chip failing to reach the highest quality level.

3  
4 The use of universal reference RNA has appeal when conducting experiments using gene  
5 arrays in a two-color hybridization approach. In such experiments, both the control and treated  
6 samples are labeled separately with a single sulfonated indocyanine fluorescent dye (Cy<sup>TM</sup>; e.g.,  
7 Cy3 or Cy5) and are compared to a reference RNA sample which is labeled with the other of the  
8 two Cy<sup>TM</sup> dyes. Not only does this approach help minimize the potential for dye bias, which is a  
9 significant concern when using the two-dye hybridization approach, but this also allows for  
10 comparison of data across studies that use the same reference RNA. One practical approach may  
11 be to take advantage of commercially available universal reference RNAs for gene expression  
12 profiling which, at this time, are offered for use with arrays representing a limited number of  
13 organisms (human, mouse, rat). Another experimental design often used to address dye bias is  
14 the ‘dye swap’ or ‘dye flip’. In this method a second experiment is conducted by exchanging  
15 labeling reactions such that the treated and control samples are conversely labeled with the  
16 respective Cy<sup>TM</sup> dyes. The approach entails the use of additional arrays; however, because dye  
17 bias has been observed by numerous investigators and noted in the literature, such a scheme  
18 should be considered when designing experiments using two-color array systems.

### 19 20 **Experimental Replication**

21  
22 It is not possible to analyze expression data without an estimate of variance. Since  
23 experimental variance has both technical and biological components, replication could be  
24 incorporated at several levels. In the case of a gene array experiment, technical replication could  
25 be in the form of multiple spots per gene on the same array or, perhaps, multiple arrays for a  
26 given sample. While including technical replicates will improve data analysis it is not an  
27 absolute necessity. On the other hand, biological replication is an important consideration.  
28 While it is generally accepted that in a gene array experiment an absolute minimum of three  
29 biological replicates is needed, additional replication is often needed to detect a treatment effect  
30 when less than robust changes in gene expression are observed. Pilot studies could be conducted  
31 to estimate variance and give insight as to what level of replication may be useful. Although not

---

1 comprehensive, additional considerations for determination of replicate numbers are the relative  
2 quality and integrity of samples, the range of expected effects and the method of raw data  
3 analysis. The optimal replicate number is affected other factors such as the type of array  
4 technology and platform (single or dual channel RFU capture), array platform linearity  
5 (precision), feature density (number of representative gene probes), and the selected percentage  
6 value of FPR. Since replication is an asymptotic process, even a small number of replicates will  
7 strengthen any conclusions that can be drawn from the data, irrespective of the technological  
8 approach used to collect these data.

### 9 **Pooling of Samples**

10  
11 From a theoretical perspective, most biological material used in expression studies arises  
12 from ‘pooled’ sources because most tissues used in such investigations contain many distinct cell  
13 types. Pooling of samples is primarily encouraged in those cases where the quantity of nucleic  
14 acid ‘target’ (total RNA) is limiting to the point that this represents the only means by which to  
15 obtain the requisite mass. It is recognized that, in certain studies, pooling of samples across  
16 individuals is a logical approach in order to limit study size. In fact, pooling of samples can help  
17 to minimize biological variation. However, it should be recognized that pooling will not be as  
18 effective in controlling biological variation as increasing the number of biological replicates in a  
19 study. Theoretically, most total RNA samples are pooled, since they are isolated from cells of  
20 related or different types after having been amplified from the original source to produce the test  
21 product. Combining samples does have the advantage of decreasing noise in the system. If  
22 biological variability is not a major concern, a pooled sample could be considered the same as a  
23 single individual when applying an experimental design. If biological variability is important to  
24 the interpretation of the data, and RNA from pooled sources is used in determination of  
25 expression measurements, it is useful to include more than one independent pool of samples for  
26 the purpose of estimating biological variability. Biological replicates are generally regarded as  
27 more critical than are technical replicates to measures of expression in biologic systems unless  
28 otherwise indicated.

### 29 **Specificity and Sensitivity**

30

---



1           **Specificity** and **sensitivity** of assays are affected by sequence-dependent (length and  
2 inclusive base composition) and sequence-independent (relative concentrations of probes and  
3 targets, hybridization time, temperature, etc.) factors. The specificity and sensitivity of assays  
4 have been the subject of numerous cross platform comparison studies recently cited in scientific  
5 literature (Venkatasubbarao, 2004; Enders, 2004; de Longueville *et al.*, 2004). The term  
6 specificity refers to the ability of an expression platform to discriminate or select between  
7 distinct members of the same gene family, whereas sensitivity is the potential to discriminate  
8 transcripts expressed at low level in a complex background. In recent years there has been a  
9 trend in microarray design towards oligonucleotide probe sets to improve the specificity of gene  
10 targeting. Oligonucleotide microarrays (25 to 70 bp) have some advantages over arrays on  
11 which cDNA probes have been affixed. Oligonucleotide probes are designed to be identical with  
12 respect to the number of bases (length) and concentration, with comparable annealing  
13 temperatures of hybridization. These considerations account for enhanced uniformity over the  
14 entire platform. Oligonucleotides are also designed to reduce inadvertent target cross-  
15 hybridization, thereby increasing specificity during hybridization reactions. These combined  
16 properties increase the stability and reproducibility of hybridization signal on each feature on the  
17 array.

18  
19           In addition to the quality of the probe sequences, the specific region of a gene that is  
20 selected as a probe to be incorporated onto the array is also critically important. For example, if  
21 probes are selected primarily from the 3' end of given genes, as is often the case, there is a  
22 distinct possibility that splice variants of those genes will evade identification if the alternative  
23 splicing events occur 5' of a probe region. Additionally, it is very difficult to distinguish  
24 between two expressed genes that share a high degree of sequence homology by microarray  
25 analysis. Irregularity in probe specificity is also a frequently encountered problem in  
26 oligonucleotide arrays. This problem frequently arises in instances where nucleic acid sequences  
27 are practically identical between two coding regions and the oligonucleotide probes are  
28 synthesized from 3' ends of the genes.

29  
30           Decrease in specificity on microarray platforms generally results from the technical  
31 limitations inherent in enzymatic labeling of the RNA target. One of the most widely used

---

1 methods for enzymatic modification of total RNA, for microarray analysis of gene expression,  
2 uses T7 viral RNA polymerase *in vitro* transcription (IVT) to produce complementary RNA  
3 (cRNA) that can be hybridized to gene-specific probes affixed to arrays. Multiple rounds of  
4 amplification are used to label a limited mass of RNA by this IVT method, which has been  
5 shown to inadvertently introduce errors. Because cRNA-DNA sequence mismatches are more  
6 thermo stable than comparable cDNA-DNA mismatches, intensity artifacts have been observed  
7 due to increased non-specific hybridization.

### 8 9 **System Linearity and Calibration**

10  
11 Linearity of signal responses and other measurable output are perhaps among the most  
12 significant aspects of obtaining reliable gene expression data. Regardless of technological  
13 modes (*e.g.*, microarray-based studies, semi-quantitative gel based PCR, ‘real-time’ PCR, or  
14 densitometric scanning of pixel density), usable data collected within the linear region of the  
15 output curve for any chosen system is essential. Given the increased number and overall density  
16 of gene-specific probes present on microarrays, it is particularly useful to demonstrate linearity  
17 of relative fluorescence units (RFUs) for the greatest number of discrete features represented on  
18 the chip. Recent observations from microarray workgroups suggest that specific reference RNA  
19 is the most efficient means by which to accomplish this. In an attempt to measure precision (B.  
20 Aronow, personal communication), it was determined that the greatest coverage of features was  
21 attained by hybridizing 4-5 different ratio mixtures, on as many chips, of species-specific RNA  
22 obtained from different sources. For instance, the study in question mixed RNA prepared from  
23 the colons of 8-week old C57BL/6 8 mice and post partum day one C57BL/6 whole animals in  
24 different proportions. The relative fluorescent unit (RFU) value changes for every gene probe  
25 that yielded a response to the mixture of mouse RNA, were statistically analyzed using least-  
26 square linear regression. This suggested approach permits investigators to ascertain a global  
27 perspective regarding the degree of linear response in a chosen system.

### 28 29 **Randomization of Samples**

30

---

1            Technical variations or differences in “expression” measurements can be introduced at  
2 several junctures in the experimental process including, but not limited to, methods of RNA  
3 labeling, the choice of microarray platform, capture methods for RFU intensity and signal  
4 quantitation, ozone-mediated fluorescent signal degradation, humidity and temperature, and  
5 moreover, those individuals charged with performing the experiments.  
6

7            Many experimental designs suggest that blind randomization of samples is integral to the  
8 analyses. This approach offers the promise of ‘flattening’ both internal and external  
9 experimental sources of variation. Sample randomization should be considered wherever  
10 practical. Numerous confounding variables have been identified that can distort microarray  
11 results. Some of these sources of variation are well known (*e.g.*, RNA degradation during tissue  
12 extraction), and others have been more recently identified (*e.g.*, ozone-mediated bleaching of  
13 some florescent dyes), and some causative factors have yet to be characterized. If adequate care  
14 is taken to randomize the order in which samples are processed, and operators are unaware of the  
15 nature of each sample, known and unanticipated sources of variability are not likely to bias the  
16 outcome of the experiments. However, such sources of variability can nevertheless exert  
17 influence on the observed precision of the system. For instance, if all the experimental group  
18 samples are run on a given day, and all the negative control samples are run on the following  
19 day, it is possible that experimental features can differ on the two days (*e.g.*, operator identity,  
20 photomultiplier drift, and/or ozone concentration). Such differences could systematically bias  
21 results for the experimental samples relative to the control samples, creating the false impression  
22 of a real difference between the two groups. On the other hand, if samples for the two groups are  
23 randomized, with half the samples run on day one, and the other half on day two, factors that  
24 differ between the two days will decrease the precision observed in both groups (a readily  
25 detected and addressed occurrence), without creating the false impression of systematic  
26 differences between the groups in question.

---

## APPENDIX B: MIAME-Based Data Submission Tables

<b>Table B.1 Abstract</b>				
<b>MIAME</b>	<b>Description</b>	<b>When applicable</b>	<b>Notes</b>	<b>Values</b>
<b>B.1 Abstract</b>	Brief summary of the purpose and findings of the experiment.	Always		

<b>Table B.2 Experimental Design</b>				
<b>MIAME</b>	<b>Description</b>	<b>When applicable</b>	<b>Notes</b>	<b>Values</b>
<b>B.2 Experiment design</b>	Design and purpose common to all hybridizations	Always	Related hybridizations interpreted as a single experiment.	
<u>Author, laboratory, and contact</u>	Person(s), organization(s), names and contacts (address, phone, FAX, email, URL).	Always		Contact details
<u>Experiment type(s)</u>	A controlled vocabulary that classifies an experiment.	Always	<u>Experimental Factor(s)</u> .	Time course, dose response, comparison (disease vs normal, treated vs untreated), temperature shock, gene knock out, gene knock in (transgenic), etc.
<u>Experiment Description</u>	Description of the experiment and relevant electronic peer-reviewed journal publication(s)	When additional information is available and an electronic publication exists.	Consistent with experimental design.	Text description, citation, URL. Database entry

<b>Table B.2 Experimental Design</b>				
<b><u>MIAME</u></b>	<b><u>Description</u></b>	<b><u>When applicable</u></b>	<b><u>Notes</u></b>	<b><u>Values</u></b>
<u>Experimental factor(s)</u>	Parameter(s) or condition(s) tested in the experiment.	Always	Experimental factor(s) consistent with <u>Experiment Type(s)</u>	Time, dose, compound, temperature, extraction, hybridization, labelling, scanning
<u>Number of hybridization replicates</u>	Number of hybridization replicates	Always	Consistent with <u>Experiment Type(s)</u>	Single, multiple
<u>Common reference</u>	A hybridization to which all the other hybridizations have been compared.	Always		Yes, no
<u>Quality control steps</u>	Measures to ensure quality: replicates (number and description), dye swap (for two channel platforms) or other	When appropriate		Text description. biological, technical
<u>Qualifier, value, source (may use more than once)</u>	Any further information about the experiment .	Additional useful information		Qualifier= name Value= value Source= database entry or ontology entry

<b>Table B.3 Array Design</b>				
<b>MIAME</b>	<b>Description</b>	<b>When applicable</b>	<b>Notes</b>	<b>Values</b>
<b>B.3 Array design</b>	Array layout. Description of the common features of the array and each array element.	When an array design is novel and cannot refer to manufacturer	Array design should be provided by the array manufacturer.	
<b>B.3.1. Array related information</b>	Overall description of the array.			
<u>Array design name</u>	Unique name, that identifies a specific design	Array is novel and cannot refer to manufacturer	Consistent with the design name given for the array.	Design name, number of features, version (e.g.: EMBL yeast 12K ver1.1)
<u>Platform type</u>	Technology to place biological sequence on array.	Array is novel and cannot refer to manufacturer		in situ synthesized, spotted cDNA, etc.
<u>Surface and coating specification</u>	Surface coating <u>type and name</u>	Array is novel and cannot refer to manufacturer	Consistent with <u>Platform Type</u>	SurfaceType: glass, membrane, coating type
<u>Array dimensions</u>	Dimensions of the array support slide.	Array is novel and cannot refer to manufacturer		width, length
<u>Number of features on the array</u>	The number of features on the array.	Array is novel and cannot refer to manufacturer		number of features
<u>Production protocol</u>	A description of how the array was manufactured.	Array is novel and cannot refer to manufacturer		Protocol description, printing hardware, printing software
<u>Provider</u>	The primary contact (manufacturer) for the information on the array design.	Always		Contact details of manufacturer

<b>Table B.3 Array Design</b>				
<b><u>MIAME</u></b>	<b><u>Description</u></b>	<b><u>When applicable</u></b>	<b><u>Notes</u></b>	<b><u>Values</u></b>
<b>B.3.2 Reporter related information</b>	Information on the nucleotide sequence present in a particular location on the array.			
<b>B.3.2.1 For each reporter type</b>				
<u>Reporter type</u>	Physical nature of the reporter (e.g. PCR product, synthesized oligonucleotide).	Array is novel and cannot refer to manufacturer	Consistent with <u>Platform Type</u>	Types: empty, PCR, synthesized oligonucleotide, plasmid, colony, etc.
<u>Single or double stranded</u>	Reporter sequences are single or double stranded.	Array is novel and cannot refer to manufacturer	Consistent with <u>Platform Type</u>	Single, double
<b>B.3.2.2 For each reporter</b>				
<u>Reporter sequence information</u>	Nucleotide sequence for each reporter: accession number (from DDBJ/EMBL/GenBank), the sequence itself or reference sequences and primers pair information	Array is novel and cannot refer to manufacturer	Consistent with <u>Platform Type</u> and clone	Sequence annotation, accession number, PCR primer pair
<u>Reporter approximate length</u>	The approximate length of the reporter sequence.	When the exact reporter sequence is NOT known		Number of bases
<u>Clone information</u>	For each reporter, identity of the clone, clone provider, date obtained, and availability.	When elements are from clones When an array design is novel and cannot refer to manufacturer	Consistent with <u>Platform Type</u>	Clone ID, provider, date obtained, availability
<u>Reporter generation protocol</u>	A description of how the reporters were generated.	Array is novel and cannot refer to manufacturer		Protocol

<b>Table B.3 Array Design</b>				
<b><u>MIAME</u></b>	<b><u>Description</u></b>	<b><u>When applicable</u></b>	<b><u>Notes</u></b>	<b><u>Values</u></b>
<b>B.3.3 Features related information</b>	Information on the location of the reporters on the array			
<b>B.3.3.1 For each feature type</b>				
<u>Feature dimensions</u>	Dimensions of each feature.	Array is novel and cannot refer to manufacturer	Consistent with array dimensions and number of features	Width, length, height, diameter
<u>Attachment</u>	How the elements (reporters) are physically attached to the array.	Array is novel and cannot refer to manufacturer	Consistent with element generation protocol	Covalent, ionic, hydrophobic, etc.
<b>B.3.3.2 For each feature</b>				
<u>Reporter and location</u>	Arrangement and system used to specify location of each feature	Array is novel and cannot refer to manufacturer	Consistent with array dimensions and number of features	Row, column, x microns, y microns, zone
<b>B.3.4 Composite sequence related information</b>	Information on the set of reporters used collectively to measure an expression of a particular gene.			
<b>B.3.4.1 For each composite sequence</b>				
<u>Composite sequence information</u>	The set of reporters contained in the composite sequence.	When elements are composite array is novel cannot refer to manufacturer	Consistent with element type	Oligonucleotide sequences, number of oligonucleotides, reference sequence
<u>Gene name</u>	The gene represented at each composite sequence	Array is novel and cannot refer to manufacturer	Consistent with clone and composite sequence information	Gene name, accession number, annotation
<u>Qualifier, value, source (may use more than once)</u>	Describe any further information about the array in a structured manner.	When additional information is available that would be useful to base queries on		Qualifier= name Value= value Source= database entry or ontology entry



<b>Table B.3 Array Design</b>				
<b><u>MIAME</u></b>	<b><u>Description</u></b>	<b><u>When applicable</u></b>	<b><u>Notes</u></b>	<b><u>Values</u></b>
<b>B.3.5 Control elements related information</b>	Array elements that have an expected value and/or are used for normalization.			
<u>Control element position</u>	The position of the control features on the array.	When any elements on the array were used as controls	Consistent with Quality Control Description	Row, column, x microns, y microns, zone
<u>Control type</u>	The type of control used for the normalization and their qualifier.	When any elements on the array were used as controls	Consistent with Quality Control Description	Control type (spiking, negative, positive), control qualifier (endogenous, exogenous)

<b>Table B.4 Biomaterials</b>				
<b><u>MIAME</u></b>	<b><u>Description</u></b>	<b><u>When applicable</u></b>	<b><u>Notes</u></b>	<b><u>Values</u></b>
<b>B.4 Biomaterials</b>	The biological material from which the nucleic acids have been extracted for subsequent labelling and hybridization.	Always		
<b>B.4.1 Biosource properties</b>	Information on the source of the sample.			
<u>Organism</u>	The genus and species (and subspecies) of the organism from which the biomaterial is derived.	Always		Genus, species, subspecies from NCBI taxonomy

<b>Table B.4 Biomaterials</b>				
<b>MIAME</b>	<b>Description</b>	<b>When applicable</b>	<b>Notes</b>	<b>Values</b>
<u>Sample Contact details</u>	The resource used to obtain the biomaterial	When biomaterial was prepared or grown outside of the laboratory listed for the author		Biosource provider Type of specimen (tumor biopsy, paraffin section, stool sample}
<u>Cell type</u>	Cell type(s) or organs used in the experiment.	Always	Consistent with organism and targeted cell type	Name of organ tissue cell type (ATCC # ) and source
<u>Sex</u>	Term applied to any organism able to undergo sexual reproduction in order to differentiate the individuals or type involved.	When applicable	Consistent with organism	Mating type alpha, F <sup>+</sup> , F <sup>-</sup> , Hfr, Mating type a, Mixed sex, Unknown sex
<u>Age</u>	The time period elapsed since an identifiable point in the life cycle of an organism.	When applicable	Consistent with organism	Age = combination of real number (measurement) and initial time point e.g.: coitus, birth, planting, beginning of stage
<u>Developmental stage</u>	The developmental stage of the organism's life cycle during which the biomaterial was extracted.	For multicellular species	Consistent with organism	Developmental stage (i.e., embryo, fetus, adult)
<u>Organism part</u>	The part or tissue of the organism's anatomy from which the biomaterial was derived.	For multicellular species	Consistent with organism	Organism part term)
<u>Strain or line</u>	Animals or plants that have an ancestral breeding.	When known	Consistent with organism	Strain or line ( e.g.: Jax mouse strains, Cultivar,NCBI taxonomy)

<b>Table B.4 Biomaterials</b>				
<b>MIAME</b>	<b>Description</b>	<b>When applicable</b>	<b>Notes</b>	<b>Values</b>
<u>Genetic variation</u>	The genetic modification introduced into the organism from which the biomaterial was derived.	When the source organism is genetically modified	Consistent with organism	Examples of genetic variation include specification of a transgene or the gene knocked-out.
<u>Individual number</u>	Identifier or number of the individual organism from which the biomaterial was derived.	When the organism can be distinguished on an individual basis with a unique ID	Consistent with organism	Individual ID. For patients, the identifier should be approved by Institutional Review Boards (IRB, review and monitor biomedical research involving human subjects) or appropriate body.
<u>Individual genetic characteristics</u>	The genotype of the individual organism from which the biomaterial was derived.	When applicable	Consistent with organism	Allele, genotype, haplotype, polymorphisms.
<u>Disease state</u>	The name of the pathology diagnosed in the organism from which the biomaterial was derived.	When applicable	Consistent with organism	"Normal" or . disease state description
<u>Targeted cell type</u>	<u>Cell</u> of primary interest.	Biomaterial is a mixed population of cells	Consistent with organism and cell type Biomaterial may be derived from a mixed population of cells although only one cell type is of interest.	Targeted cell type= term, (Mouse Anatomical Dictionary, FlyBase, CBIL vocabulary)
<u>Cell line</u>	Identifier for the cell line	Biomaterial is derived from an immortalized cell line	Consistent with organism and cell type	Cell line term, source of term ( ATCC # ),e.g.,Hela, Caco-2

<b>Table B.4 Biomaterials</b>				
<b><u>MIAME</u></b>	<b><u>Description</u></b>	<b><u>When applicable</u></b>	<b><u>Notes</u></b>	<b><u>Values</u></b>
<b><u>B.4.2 Biomaterial manipulation</u></b>	Information on the treatment applied to the biomaterial			
<u>Growth conditions</u>	Description of environment used to grow organisms			Culture condition details
<u>In vivo treatment</u>	Manipulation to generate variable(s) under study .	When sample has been treated or manipulated for the study	Consistent with Experiment Type and Experimental Factors	Documentation of the set of steps taken in the treatment
<u>In vitro treatment</u>	Manipulation of cell culture condition for generating variables under study.	When the sample has been treated or manipulated in vitro for the study purpose	Should be consistent (where appropriate) with Experiment Type, Experimental Factors	Documentation of the set of steps taken in the treatment
<u>Treatment type</u>	Manipulation for generating variables under study.	When sample has been treated or manipulated for the study	Consistent with experiment type, experimental factors and treatment	Description of treatment (behavioral stimulus, compound based treatment, infection, modification (genetic, somatic)),
<u>Compound</u>	Drug, solvent, chemical, etc., that can be measured.	When sample has been treated or manipulated with a compound	Consistent with treatment type	Description of compound's physical and chemical characteristics
<u>Separation technique</u>	Technique to separate tissues or cells.	When the cells or tissue are separated from a heterogenous sample		Protocol

<b>Table B.4 Biomaterials</b>				
<b>MIAME</b>	<b>Description</b>	<b>When applicable</b>	<b>Notes</b>	<b>Values</b>
<b>B.4.3 Hybridization extract preparation</b>	Information on the extract preparation for each extract prepared from the sample			
<u>Extraction method</u>	The protocol used to extract nucleic acids from the sample.	Always		Protocol
<u>Nucleic acid type</u>	The type of nucleic acid extracted (e.g. total RNA, mRNA).	Always		Polymer type (total RNA, mRNA, DNA)
<u>Amplification method</u>	The method used to amplify the nucleic acid extracted.	When applicable		Protocol
<b>B.4.4 Sample labelling</b>	Information on the labelling preparation for each labelled extract.			
<u>Amount of nucleic acid labelled</u>	Amount of nucleic acid labelled.			Protocol
<u>Label used</u>	Label used.	Always		Label (Cy3, Cy5, etc.)
<u>Label incorporation method</u>	Label incorporation method	Always		Protocol
<b>B.4.5 Spiking control</b>	External controls added to the hybridization extract(s).			
<u>Spiking control feature</u>	Position of the feature(s) on the array expected to hybridize to the spiking control.	When applicable	Consistent with quality control description	row, column, x microns, y microns, zone

<b>Table B.4 Biomaterials</b>				
<b>MIAME</b>	<b>Description</b>	<b>When applicable</b>	<b>Notes</b>	<b>Values</b>
<u>Spike type and qualifier</u>	Type of spike used and its qualifier	When applicable	Consistent with quality control description	Oligonucleotide, plasmid DNA, transcript, concentration, expected ratio, labelling methods
<u>Qualifier, value, source (may use more than once)</u>	Describe any further information about the sample in a structured manner.	When additional information is available that would be useful to base queries on		

<b>Table B.5 Hybridization</b>					
<b>MIAME</b>	<b>Description</b>	<b>When applicable</b>	<b>Notes</b>	<b>Values</b>	<b>Included in DER?</b>
<b>B.5 Hybridization</b>	Procedures and parameters for each hybridization.	Always			
<u>Relationship between samples and arrays</u>	Relationship between the labelled extract	Always	Consistent with technology quality control	Which sample, which extract "array design, batch and serial number, during which hybridization	Yes
<u>Hybridization protocol</u>	Set of steps taken in the hybridization: (solution blocking agent, concentration, wash procedure); quantity of labelled target used; time; concentration; volume, temperature.	Always		Description of the hybridization instruments	Yes

<u>MIAME</u>	<u>Description</u>	<u>When applicable</u>	<u>Notes</u>	<u>Values</u>	<u>Included in DER?</u>
Qualifier, value, source (may use more than once)	Describe any further information about the hybridization in a structured manner.	When additional information is available that would be useful to base queries on			Non-specific

<u>MIAME</u>	<u>Description</u>	<u>When applicable</u>	<u>Notes</u>	<u>Values</u>
<b>B.6.1 Raw data</b>	Each hybridization has at least one image.			
<u>Scanner image file</u>	The image file including header	Always		TIFF, JPEG
<u>Scanning protocol</u>	Steps taken for scanning array and generating an image	Always		Description of the scanning instruments and the parameter settings.
<b>B.6.2 Image analysis and quantitation</b>	Each image has a corresponding image quantitation table, where a row represents an array design element and a column represents different quantitation types			Mean or median pixel intensity.
<u>Image analysis output</u>	The complete image analysis output for each image.	Always.		Spreadsheet or tab-delimited file

<b>Table B.6 Measurements</b>				
(MIAME distinguishes between three levels of data processing: image (raw data), image analysis and quantitation, gene expression data matrix (normalized and summarized data).				
<u>MIAME</u>	<u>Description</u>	<u>When applicable</u>	<u>Notes</u>	<u>Values</u>
<u>Image analysis protocol</u>	Documentation of the set of steps taken to quantify the image	Always.		Image analysis software, the algorithm and all the parameters used
<b>B.6.3 Normalized and summarized data</b>	Several quantitation tables are combined using data processing metrics to obtain the 'final' gene expression measurement table (gene expression data matrix) associated with the experiment.			
<u>Data processing protocol</u>	Documentation of the set of steps taken to process the data.	When normalization has been performed		Normalization strategy and the algorithm used to allow comparison of all data.
<u>Final gene expression table (s)</u>	Derived measurement value summarizing related elements and replicates, providing the type of reliability indicator used.	When a value used for a reliability indicator has been generated	Should be consistent with quality control description and replicate description	Replicates of the elements on the same or different arrays or hybridizations, as well as different elements related to the same entity (e.g., gene). Reliability indicator for each data point (e.g., standard deviation)
<u>Qualifier, value, source (may use more than once)</u>	Describe any further information about the measurements in a structured manner	When additional information is available that would be useful to base queries on		



---

**Appendix C: Genomics Data Evaluation Record (gDER) Template**

<b>Genomics DATA EVALUATION RECORD</b>
--

**STUDY TYPE:****PC CODE:****DP BARCODE:  
SUBMISSION NO.:****TEST MATERIAL (PURITY):****SYNONYMS:****CITATION:****SPONSOR:****EXECUTIVE SUMMARY:****COMPLIANCE:****I. MATERIALS AND METHODS****A. MATERIALS:****1. Test Material:****Description:****Lot/Batch #:****Purity:****CAS # of TGAI:***[Structure]***2. Control Materials:**

Negative control (if not vehicle) :

Final Volume:

Route:

Vehicle:

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42

**3. Test animals:**

- Species:
- Strain:
- Age/weight at study initiation:
- Source:
- No. animals used per dose per duration:
- Properly Maintained?

**4. Compound administration:**

**a. Test material**

- Dose levels:
- Route:
- Method:

**b. Vehicle control:**

**c. Positive control:**

**B. TEST PERFORMANCE**

**1. Treatment and Sampling Times:**

- Duration of dosing:
- Frequency of dosing:
- Total number of doses:
- Timing and frequency of sampling:
- Time elapsed between dosing and sampling:

**2. Tissues and Cells Examined:**

**3. Details of tissue harvest:**

**4. Detail of tissue storage:**

**5. Sample Preparation:**

- a. RNA Isolation, Labelling, Amplification:**
  - b. Histology:**
  - c. Immunochemistry:**
  - d. Western blot analysis:**
  - e. Array analysis:**
-



1 **Appendix D: Genomics Data Evaluation Record (gDER) for**  
2 **Alachlor (Sample)**

3  
4 

<b>Genomics</b> <b>DATA EVALUATION RECORD</b>
--

5  
6 **STUDY TYPE: Mode of Action *In vivo* Genomic Analysis in Rat Olfactory Mucosa**

7  
8 **PC CODE:**

9 **DP BARCODE:**  
10 **SUBMISSION NO.:**

11 **TEST MATERIAL (PURITY): Alachlor (Purity not listed)**

12  
13 **SYNONYMS:**

14  
15 **CITATION: Genter, M.B., Burman, D.M., Vijayakumar, S., Ebert, C.L., Aronow, B.J. (2002).**  
16 **Genomic analysis of alachlor-induced oncogenesis in rat olfactory mucosa.**  
17 **Physiol. Genomics 12:35-45.**

18  
19 **SPONSOR:**

20  
21 **EXECUTIVE SUMMARY:**

22  
23 In an *in vivo* genomic analysis, groups of male Long-Evans rats (1-2 rats/ group) were  
24 administered dietary preparations of an established tumorigenic dose of Alachlor (126  
25 mg/kg/day) or untreated feed 1 day to 18 months. Ethmoid turbinates were removed and frozen  
26 in liquid N<sub>2</sub>. Animals were sacrificed at 3, 4 or 5 months and two separate olfactory mucosal  
27 RNA samples were isolated. Other RNA samples were harvested from single rats treated with  
28 alachlor for 1 or 4 days. After 18 months of treatment, single RNA samples derived from  
29 alachlor-induced tumors were also isolated. Total RNA was extracted from frozen tissue  
30 homogenates by precipitation with ethanol/sodium acetate, screened for quality, labeled with  
31 biotin and hybridized. Histological examinations were performed on additional rats dosed for  
32 the same treatment duration. For the determination of ebrinin (a gene related to the human  
33 tumor suppressor gene, DMBT1) or  $\exists$ -catenin (gene/product associated with the wnt signaling  
34 pathway), immunochemistry was also performed on sections prepared for histology using an  
35 anti-hensin antibody or a commercial antibody; intestinal sections served as the positive control  
36 for antibody staining. CYP2A3 levels were assessed by Western blot analysis. For the array  
37 analysis, total RNA was reverse transcribed followed by second-strand cDNA synthesis. The  
38 resulting cRNA was biotinylated and hybridized to the Affymetrix GeneChip Rat U34A.  
39

---

1 Based on an independent review of qualitative data only (presented as graphs or photographic  
2 copies of tissue sections), it was concluded that alachlor induces olfactory nasal carcinomas  
3 through a nongenotoxic mode of action (*i.e.*, oxidative stress). Support for this conclusion comes  
4 from data showing upregulation ( $\geq 2$ -fold increase over untreated control) of genes correlated  
5 with the following steps in the carcinogenic process: oxidative stress and damage to DNA  
6 (8heme oxygenase, glutathione and metallothionein, GADD 45, apurinic/apurimidinic  
7 endonuclease); progression of adenomas to malignant adenocarcinomas (activation of the wnt  
8 signaling pathway), and transformation to adenocarcinomas (activation of nuclear  $\beta$ -catenin  
9 genes, also associated with the wnt signaling pathway).

10  
11 This study is classified as acceptable (non-guideline) but does not satisfy current regulatory data  
12 requirements for pesticides. Although guidelines do not exist for genomic data, the results  
13 presented in this published article provided critical information that enhances the understanding  
14 of the nongenotoxic mode of action for olfactory mucosal tumors induced by alachlor.

15  
16 **COMPLIANCE:** Not applicable; the publication, however, comes from a reputable, peer-  
17 reviewed scientific journal.

## 18 19 20 **I. MATERIALS AND METHODS**

### 21 22 **A. MATERIALS:**

#### 23 24 **1. Test Material:**

	Alachlor
<b>Description:</b>	Not provided
<b>Lot/Batch #:</b>	Not reported
<b>Purity:</b>	% a.i. Not reported
<b>CAS # of TGAI:</b>	

[Structure]

#### 25 26 **2. Control Materials:**

Negative control

(if not vehicle) :

Vehicle:

**Final Volume:** NA    **Route:** NA

Harlan powder diet

#### 27 28 29 30 **3. Test animals:**

Species: Rat

Strain: Long-Evans

Age/weight at study initiation: Not specified

Source: Harlan, Indianapolis, IN

No. animals used per dose per duration 1-2 males; 0 females

Properly Maintained? Not specified

36  
37

**4. Compound administration:****a. Test material**

Dose levels: 126 mg/kg/day Route Oral – Feeding

**Preliminary: Not performed, referred to a citation (Genter et al., (2000).**

**Main Study:** Dietary administration of 0 or 126 mg/kg/day (tumorigenic dose as per EPA 1985)

**b. Vehicle control:**

Untreated Harlan powdered feed

**c. Positive control:**

See Immunochemistry

**B. TEST PERFORMANCE**

**1. Treatment and Sampling Times:** Male rats were fed dietary preparations of 0 or 126 mg/kg/day for 1 day to 126 mg/kg.

Sampling (after last dose): 1 and 4 days, 3, 4 and 5 month. Tumors were harvested from rats treated for 3, 4, 5, 11 or 18 months.

**2. Tissues and Cells Examined:** Ethmoid turbinates and/or olfactory mucosal tumors

**3. Details of tissue harvest:** Rats were sacrificed by a pentobarbital overdose and decapitated. Ethmoid turbinates were rapidly removed and frozen in liquid N<sub>2</sub> until use. RNA samples of olfactory mucosal tumors derived from two different rats were also harvested.

**4. Detail of tissue storage:** Stored in liquid N<sub>2</sub> until use.

**5. Sample Preparation:**

**a RNA Isolation:** Selected frozen tissues were homogenized and total RNA was precipitated with ethanol/sodium acetate, resuspended in DEPC-treated water and screened for RNA quality using an Agilent Bioanalyzer. Acceptable samples had cut-off ratios of 1.8 for the 28S:18S ribosomal subunits. Duplicate samples were prepared for 2 rats/group at each sampling time (3, 4, 5 months) and one sample was used for single rats at days 1 and 4.

**b Histology:** Progression of the olfactory mucosa tumors was followed by harvesting tissue at 3, 4, 5, 11, and 18 months from dosed rats doses. Tissue was prepared for histological examinations as previously described (Genter *et al.*, 2000)<sup>1</sup> with the exception that decalcified after fixation in multiple changes of cold 0.3 M EDTA prior to embedding in paraffin.

---

<sup>1</sup>Genter, MB, Burman, DM, Dingeldein, MW, Clough, I, Bolon, B (2000). Characterization of cell proliferation and immunochemical markers of alachlor-induced olfactory mucosal tumors in Long-Evans rats. Toxicol Pathol 28:770-781.

---

- 1 **c** **Immunohistochemistry:** Sections (5- $\Phi$ m) prepared for histology were stained for  
2 detection of ebnerin (a gene related to the putative human tumor suppressor gene,  
3 DMBT1) using immunohistochemistry techniques. Ebnerin was localized in the nasal  
4 cavity sections and in tumors with anti-hensin antibody. Localized ebnerin was  
5 reacted with anti-guinea pig horseradish peroxidase (HRP)-conjugated secondary  
6 antibody (1:100) with tyramide signal amplification. Since this antibody also  
7 detects the intestinal crp-ductin, intestinal sections were taken from the rats,  
8 stained and served as the positive control.  $\exists$ -catenin (associated with the  
9 activation of the wnt signaling pathway) was localized with antibody from BD-  
10 Transduction Laboratories, Lexington, KY and visualized using an HRP-  
11 conjugated anti-mouse secondary antibody and TSA amplification as described.  
12
- 13 **d** **Western blot analysis:** Gene expression changes in the olfactory specific  
14 cytochrome P-450 enzyme (CYP2A3), were assessed by Western blot analysis of  
15 5  $\Phi$ g of olfactory mucosal microsomal protein per lane. Visualization was  
16 achieved with HRP-conjugated secondary antibody, enhanced chemiluminescence  
17 and exposure to X-ray film.  
18
- 19 **e** **Array analysis:** Characteristics of the arrays: The total number of probe sets  
20 (genes or expressed sequence tags, ESTs) interrogated was not reported.  
21 However, the study authors provided the following information: ESTs represented  
22 on the U34A GeneChip were derived from the Rat Unigene Build No. 34  
23 assembly. All clones represented on the chip were ESTs on gene lists of interest  
24 and were subjected to re-annotation by use of Unigene and execution of the  
25 National Center for Biotechnology Information (BLASTN) searches  
26 (<http://www.ncbi.nlm.gov/BLAST>) against non-redundant nucleotide databases  
27 during February-April, 2002. Gene category information was based on all  
28 publically available gene ontology information from the Gene and Ontology  
29 Consortium (<http://www.geneontology.org>) as harvested from SWISS-PROT,  
30 GeneCards, Compugen, LocusLinks, and GeneBank as well as exhaustive  
31 Medline literature searches.  
32

33 **Methods:** Total DNA was reverse transcribed with an oligo-dT primer and  
34 second-strand cDNA was synthesized. The resulting T7 RNA polymerase-  
35 mediated cRNA was biotin-labeled and hybridized on the Affymetrix GeneChip  
36 Rat U343A using the recommended protocol provided by Affymetrix  
37

## 38 **7. Data Analysis:**

39  
40 The study authors provided the following information. "MicroArray Suite 5.0 software  
41 was used to scan and quantitate the GeneChip data using a default scan setting; intensity  
42 data were scaled to target intensity of 1,500, and results were analyzed using the  
43 MicroArray Suite 5.0 and GenSpring 4.1.5 software. Data values used for filtering and  
44 clustering were "signal", "signal confidence", "absolute call" (absent or present), and  
45 "change" (increased, decreased, unchanged) as implemented in MicroArray Suite 5.0.

---

1 Data were normalized as follows: the 50<sup>th</sup> percentile of all measurements was used as the  
2 positive control for each array. Each measurement for each gene was divided by this  
3 synthetic positive control, assuming that this was at least 10. The bottom 10<sup>th</sup> percentile  
4 signal level was used as a test for correct background subtraction. The measurement of  
5 each gene in each sample was divided by the corresponding value in untreated samples,  
6 assuming that the value was at least 0.01. Genes regulated across the experimental study  
7 were identified by data filtering for those over- or under- expressed in at least two samples  
8 whose signal strength was greater than 500 in two samples, and were also called  
9 “present” in at least two samples. An additional approach combined those with genes  
10 that could predict length of alachlor exposure or histological responses using Kruskal-  
11 Wallis ANOVA at  $p < 0.001$  and a Benjamini-Hochberg multiple testing correction as  
12 implemented in GeneSpring. K-means analysis were similarly executed in GeneSpring  
13 to organize genes into clusters based on similar expression across the treatment time  
14 course.”

## 15 16 **8. Evaluation Criteria/Statistical Analysis**

17  
18 **Initial Filter Criteria:** The study authors indicated that 4,777 probe set elements (a pool  
19 of genes that fulfill a series of initial filter criteria) were called “present” by the  
20 Affymetrix algorithm on at least one of 26 chips, 998 probe set elements were  
21 overexpressed by 1.8X or more in at least 2 samples, whereas 584 were underexpressed  
22 0.5 X in at least 2 samples. Additionally, significant gene regulation was detected using a  
23 Welch t-test with a cutoff of  $p < 0.001$  (without correction for false discovery rate error).  
24 Using this approach, alachlor-exposed samples could be distinguished from untreated  
25 controls based on differential expression of 644 genes.

26  
27 **Cluster analysis:** 1392 probe sets elements were provided for cluster analysis by  
28 combining the under and over expressed genes along with the alachlor-regulated genes  
29 and then restricting these to only genes that met the ‘present’ criteria. Using the K-means  
30 algorithm, 16 sets were determined to serve as excellent representations of the prominent  
31 patterns in the data set. Clusters with “highly chaotic” patterns were eliminated from  
32 further analysis. Accordingly, 1,265 genes whose variance was well represented by 16  
33 K-means sets were found. These K-means sets were grouped into the following behavior  
34 patterns:

- 35 • Sets that were upregulated acutely
- 36 • Sets that were upregulated only in alachlor-induced tumors
- 37 • Sets that were downregulated following alachlor treatment
- 38 • Sets that were downregulated in alachlor-induced tumors
- 39 • Sets that were persistently upregulated across the treatment intervals.

## 40 41 **II. REPORTED RESULTS**

42  
43 **A. IMMUNOCHEMISTRY:** The study authors stated that ebnerin was highly expressed after  
44 4 months of treatment with 126 mg/kg/day alachlor in nasal respiratory mucosa. Tissue sections  
45 of olfactory mucosal tumors induced by alachlor and control nasal mucosa were provided to

---



1 support the study authors' claim that this gene product was detected in vehicle control nasal  
2 respiratory tissue but was absent in the control olfactory mucosa. Increased gene expression of  
3 this protein was also noted in the progression of alachlor-induced tumors. In the olfactory  
4 tumors, ebnerin was displayed on the surface and in the ductal lumens of the tumor. In addition,  
5 nuclear localization of  $\exists$ -catenin was also confirmed using immunochemical staining. It was  
6 stated that alachlor-induced polyps and early adenomas did not exhibit nuclear localization of  $\exists$ -  
7 catenin but more advanced adenocarcinomas displayed abundant cytoplasmic and nuclear  $\exists$ -  
8 catenin; a tissue section of olfactory mucosal tumors induced by alachlor was presented to  
9 support this claim.

10  
11 **B. WESTERN BLOT ANALYSIS:** K-mean analysis indicated that after exposure of the rats  
12 to alachlor for 2 or 4 days or 1 month, 137 genes were downregulated; included among these  
13 genes was the olfactory predominant cytochrome P-450 enzyme, CYP2A3 and CYP2F1, an  
14 olfactory marker protein. Western blot analysis, depicted in graphs and accompanied by  
15 histological alterations, also indicated that these genes/products returned to background levels in  
16 the presence of foci of respiratory metaplasia (3- and 4- month samples), in the presence of more  
17 pronounced epithelial atypia and small neoplasms present in ~25% of the animals (5-month  
18 samples). CYP2A3 and CYP2F1 were downregulated in the presence of numerous tumors, some  
19 of which were invasive. These results were supported by the composite graphs of the 16 K-  
20 means sets presented in the publication.

21  
22 **C. MICROARRAY ANALYSIS:** One-hundred and forty-eight genes and ESTs that were  
23 upregulated (*i.e.*,  $\exists$ 3-fold increase in the normalized intensity value of 1.0) with acute (1 day to 1  
24 month) exposure to alachlor were identified. These include genes associated with the control of  
25 extracellular matrix such as matrix metalloproteinase-9 (MMP-9, upregulated 9-fold),  
26 carboxypeptidase Z (upregulated 7-fold) and tissue inhibitor of metalloproteinase-1 (upregulated  
27 3-fold); immune system functions; cell proliferation/ cell cycle regulation, including apoptosis-  
28 related genes; calcium homeostasis/signaling; olfactory-related; nervous system-related;  
29 oncogene-related; transporters; and structural machinery. These genes, subgrouped according to  
30 the key functional categories listed above are presented in Table 1 of the article (see Attachment  
31 1). Other genes mentioned by the study authors as being upregulated  $\exists$ 2-fold following acute  
32 exposure included multiple genes which encode proteins associated with oxidative stress; these  
33 included heme oxygenase, glutathione synthase and metallothionein (MT)-1 and MT-2.  
34 Additionally, the GADD 45 gene (associated with mutagenesis possibly caused by oxidative  
35 damage to DNA) was listed as one of the most highly regulated genes by alachlor.

36  
37 An additional, 417 genes and ESTs were identified, based on a  $\exists$ 2-fold upregulated  
38 expression, in alachlor-induced tumors as compared to the untreated mucosa. These genes  
39 included several immune response genes (*i.e.*, neutrophil defensin, mast cell proteases, squamous  
40 cell carcinoma antigens and major histocompatible complex antigens and genes associated with  
41 cell proliferation (*e.g.*, nucleolin, the major nucleolar protein in exponentially-growing  
42 eukaryotic cells). Another set of highly expressed gene were axin2 and frizzled. The study  
43 authors claim that the increased expression of these genes is suggestive of activation of the wnt  
44 signaling pathway. This pattern is consistent with the results of immunochemical staining  
45 confirming nuclear localization of  $\exists$ -catenin late in the carcinogenesis process. Primary

1 normalized data, gene lists and K-means groups can be obtained from <http://genet.chmcc.org> in  
2 the U34A folder listed under Genter *et al.*, 2002.

### 3 4 **III. DISCUSSION and CONCLUSIONS**

5  
6 **A. INVESTIGATORS' CONCLUSIONS:** Based on these analyses, the study authors  
7 concluded that “initiation and progression of alachlor-induced olfactory mucosal tumors is  
8 associated with alterations in extracellular matrix components, induction of oxidative stress,  
9 upregulation of ebnerin, and final transformation to a malignant state by **wnt pathway**  
10 activation.”

11  
12 **B. REVIEWER COMMENTS:** Based on an independent analysis of the genomic data  
13 presented by the study authors, Agency reviewers conclude the following with respect to the  
14 proposed steps in the **alachlor**-mediated carcinogenesis model:

15  
16 • **Initial progression from histologically normal olfactory mucosa to foci of abnormal  
17 mucosa**

18 This step, which is regulated by genes in the acute phase of exposure, is accompanied by  
19 “upregulation” ( $\geq 2$ -fold increase) of genes consistent with a mutagenic response possibly  
20 as a result of oxidative damage to DNA (**8GADD 45, apurinic/apurimidinic  
21 endonuclease**). While the exact role of GADD (growth arrest and DNA-damage  
22 inducible) gene products is not known, this gene group is upregulated in response to  
23 stress to allow cells time to repair macromolecular damage or to lead cells into apoptosis  
24 so that a genetic defect is not propagated. Types of environmental stress that induce  
25 GADD genes include UV irradiation, alkylating agents and glucose starvation (Takahashi  
26 *et al.*, 2001; Jackman *et al.*, 1994). Stokes *et al.* (2002) also demonstrated that GADD 45  
27 gene induction occurs in response to reactive oxygen species (ROS) and quinones and is  
28 abolished in the presence of the antioxidant, ascorbic acid. It is of note that quinones,  
29 which are operationally non-genotoxic (Clayson *et al.*, 1994), are highly redox active  
30 molecules which can redox cycle with their semiquinone radicals, leading to formation of  
31 ROS, including superoxide, hydrogen peroxide, and ultimately the hydroxyl radical.  
32 Production of ROS can cause severe oxidative stress within cells through the formation  
33 of oxidized cellular macromolecules, including lipids, proteins and DNA (Bolton *et al.*,  
34 2000). Supporting the hypothesis of oxidative stress, Genter *et al.*, also observed  
35 upregulation of other genes associated with oxidative stress, [*i.e.*, **heme oxygenase**  
36 (Otterbein *et al.*, 2000), **glutathione synthase and metallothionein** (Andrews 2000)].

37  
38 • **Progression from histologically altered olfactory mucosa to the development of  
39 adenomas**

40 The study authors stated that this step was accompanied by expression of genes  
41 indicating inhibition of apoptosis [**Bid3(AI102299)**] and enhancement of cell  
42 proliferation (**zyxin**). However, no data were provided to support this claim.  
43 Nevertheless, it is of note that Sarafian and Bredesen (1994) state that ROS can serve as  
44 common mediators of apoptosis.  
45

---

---

- 1 • **Progression to a malignant adenocarcinoma phenotype**

2 This phase was indicated by induction of genes (*i.e.*, **axin2 and frizzled**) related to  
3 activation of the **wnt signaling pathway**, which are generally upregulated late in the  
4 carcinogenesis process.

- 7 • **Transformation to adenocarcinomas**

8 In the late stages of tumor progression, the activation of **nuclear  $\beta$ -catenin genes**, which  
9 is critical for tumor formation in other organs and is associated with mutations in the **wnt**  
10 **pathway**.

11  
12 Several other studies support a role for oxidative stress in **Alachlor**-induced toxicity.  
13 Burman *et al.* (2003) show that dietary exposure of Long-Evans rats to 126 mg/kg/day for 1 day  
14 caused an –20% depletion of the olfactory mucosa antioxidant, GSH followed by a significantly  
15 ( $p < 0.001$ ) increased expression of genes associated with increased GSH production after 2 and 4  
16 days of treatment. A return to control values was seen by 10 days of treatment. A pattern  
17 somewhat similar to GSH was observed for ascorbate in the olfactory tissue of 126-mg/kg/day  
18 male rats (*i.e.*, initially, a significant decrease 1 day post-treatment, followed by significant  
19 increases 2 and 4 days after dosing). In contrast to the GSH data, there was a reduction in  
20 ascorbate at 10 days. We noted, however, that the response with either antioxidant was not dose  
21 related. From these results, the investigators concluded that, “Despite the fact that GSH levels  
22 recovered, acute antioxidant perturbations may have been sufficient to trigger other steps in the  
23 carcinogenic process. Therefore, acute depletion of GSH and ascorbate may trigger more  
24 sustained events involved in both the initiation and promotion of the carcinogenic process.”

25  
26 There is also evidence of the ability of **alachlor** to induce oxidative stress in other tissues.  
27 Bagchi *et al.* (1995) evaluated the potential of **alachlor** to induce oxidative stress and oxidative  
28 tissue damage, as measured by production of lipid peroxidation and DNA-single strand breaks  
29 (SSB), in the liver and brain of Sprague-Dawley rats administered two equal oral doses (at 0 and  
30 21 hours) of 300 mg/kg. As noted by Clayson *et al.* (1994), SSB are considered by to be a good  
31 indicator of oxygen damage to DNA. Results from the study of Bagchi *et al.* (2003) show that  
32 **alachlor** induced moderate lipid peroxidation in liver and brain tissues and SSB in brain but not  
33 liver DNA in samples harvested 24 hours after exposure to the first dose. The same authors also  
34 conducted *in vitro* studies of chemiluminescence on liver and brain homogenates, and found that  
35 1nmol/mL **alachlor** induced 3-fold increases in chemiluminescence in both tissues further  
36 suggesting that **alachlor** induced ROS. Finally, the results from *in vitro* studies with cultured  
37 PC-12 neuroactive cells exposed to 100 nM **alachlor** illustrate the sequence of early events  
38 postulated for this MOA (generation of ROS  $\equiv$  DNA damage  $\equiv$  tissue damage) with a 2-fold  
39 increase in DNA-SSB and a 3-fold increase in LDH leakage. Although olfactory nasal tissue  
40 was not examined in this series of assays, the ability of **alachlor** to generate ROS with  
41 subsequent DNA damage and tissue damage both *in vivo* and *in vitro* has been established.  
42 Finally, Bagchi *et al.* cite the work of Akubue and Stohs (1991) showing that the oral  
43 administration of 800 mg/kg **alachlor** to rats caused the increased urinary excretion of the  
44 “oxidative lipid metabolites, malondialdehyde, formaldehyde, acetaldehyde and acetone”.

1 Based on the above considerations, the postulated MOA (generation of ROS  $\equiv$  DNA  
2 damage  $\equiv$  tissue damage  $\equiv$  cell proliferation  $\equiv$  olfactory nasal tumors) in rats is plausible and  
3 coherent. An additional factor favoring this MOA is the evidence of weak and sporadic  
4 mutagenic effects, generally seen only at concentration near or at cytotoxic concentrations.  
5

### 6 7 **C. STUDY DEFICIENCIES:** 8

9 The independent review of the data presented in this publication was limited to the  
10 analysis of qualitative results presented in graphs or photographs copies of tissue sections.  
11 Attempts to access the link for raw data provided in the article failed. Additionally, there were  
12 no data to support the study authors' claim of upregulation of genes associated with apoptosis or  
13 cell proliferation. These data would complete the sequence of key events in the carcinogenic  
14 process for alachlor. Access to the primary microarray data through a functioning, public website  
15 would have been preferable.  
16

17 Based on an independent review of qualitative genomic data (presented as graphs or  
18 photographic copies of tissue sections) in conjunction with the conventional data, it was  
19 concluded that alachlor induces olfactory nasal carcinomas through a nongenotoxic mode of  
20 action (*i.e.*, cytotoxicity manifested through oxidative stress). Partial support for this conclusion  
21 comes from data showing upregulation (2-fold increase over untreated control) of genes  
22 correlated with the following steps in the carcinogenic process: oxidative stress and damage to  
23 DNA progression of adenomas to malignant adenocarcinomas, and transformation to  
24 adenocarcinomas. Although guidelines do not yet exist for genomic data, the results presented in  
25 this DER provided critical information that enhanced the understanding of the nongenotoxic  
26 mode of action for olfactory mucosal tumors induced by alachlor in the rat.  
27

### 28 **REFERENCES** 29

- 30 Andrews, G.K., .2000. Regulation of metallathionein gene expression by oxidative stress and  
31 metal ions. *Biochem. Pharm.* 59: 95-104.  
32
- 33 Bagchi, D., Bagchi, M., Hassoun, E.A., Stohs, S.J. 1995.. In vitro and in vivo generation of  
34 ROS, DNA damage and lactate dehydrogenase leakage by selected pesticides. *Toxico.* 104: 129-  
35 140.  
36
- 37 Bolton, J.L., Trush, M.A., Penning, T.M., Dryhurst, G. Monks, T.J. 2000. Role of quinones in  
38 toxicology. *Chem. Res. Toxicol.* 13:135-160.  
39
- 40 Burman, D.M., Shertzer, H.G., Senft, A.P., Dalton, T., Genter, M.B. 2003. Antioxidant  
41 perturbations in the olfactory mucosa of alachlor-treated rats. *Biochem Pharm* 66:1707-1715.  
42
- 43 Clayson, D.B., Mehta, R., Iverson, F. 1994. Oxidative DNA damage - The effect of certain  
44 genotoxic and operationally non-genotoxic carcinogens. *Mutat. Res.* 317: 25-42.  
45
-

- 
- 1 Kasai, H. 1997. Analysis of a form of oxidative DNA damage, 8-hydroxy-2'-deoxyguanosine, as  
2 a marker of cellular oxidative stress during carcinogenesis. *Mutat. Res.* 387:147-163.  
3
- 4 Jackman, J., Alamo I.Jr., Forance, A.J. Jr. 1994. Genotoxic stress confers preferential and  
5 coordinate messenger RNA stability on the five *gadd* genes. *Cancer Res.* 54:5656-5662.  
6
- 7 Otterbein, L.E., Augustine, M.K.C. 2000. Heme oxygenase: colors of defense against cellular  
8 stress. *Am. J. Physiol. Lung Cell. Mol. Physiol.* 279: 1029-1037.  
9
- 10 Sarafian, T.A. and Bredesen, D.E. 1994. Is apoptosis mediated by ROS? *Free Rad. Res.* 21:1-8.  
11
- 12 Stokes, A.H., Freeman, W.M., Mitchell, S.G., Burnette, T.A., Hellman, G.M., Vrana, K.E. 2002.  
13 Induction of GADD 45 and GADD153 in Neuroblastoma Cells by Dopamine-Induced Toxicity.  
14 *Neuro.Toxicol.* 23:675-684.  
15
- 16 Takahashi, S., Saito, S., Ohtani, N., Sakai, T. 2001. Involvement of the Oct-1 regulatory  
17 Element of the *gadd45* Promoter in the p53-independent Response to Ultraviolet Irradiation.  
18 *Cancer Res.* 61:1187-1195.
-

## Appendix E: MIAME Glossary

**For the most recent version of the MIAME glossary, please see:**  
**[http://www.mged.org/Workgroups/MIAME/miame\\_glossary.html](http://www.mged.org/Workgroups/MIAME/miame_glossary.html)**

**Age:** The time period elapsed since an identifiable point in the life cycle of an organism. (If a developmental stage is specified, the identifiable point would be the beginning of that stage. Otherwise the identifiable point must be specified such as planting) [MGED Ontology Definition]

**Amount of nucleic acid labeled:** The amount of nucleic acid labeled

**Amplification method:** The method used to amplify the nucleic acid extracted

**Array design:** The layout or conceptual description of array that can be implemented as one or more physical arrays. The array design specification consists of the description of the common features of the array as the whole, and the description of each array design elements (*e.g.*, each spot). MIAME distinguishes between three levels of array design elements: feature (the location on the array), reporter (the nucleotide sequence present in a particular location on the array), and composite sequence (a set of reporters used collectively to measure an expression of a particular gene)

**Array design name:** Given name for the array design, that helps to identify a design between others (*e.g.*, EMBL yeast 12K ver1.1)

**Array dimensions:** The physical dimension of the array support (*e.g.* of slide)

**Array related information:** Description of the array as the whole

**Attachment:** How the element (reporter) sequences are physically attached to the array (*e.g.* covalent, ionic)

**Author, laboratory, and contact:** Person(s) and organization (s) names and details (address, phone, FAX, email, URL)

**Biomaterial manipulation:** Information on the treatment applied to the biomaterial

**Bio-source properties:** Information on the source of the sample

**Cell line:** The identifier for the immortalized cell line if one was used to derive the BioMaterial [MGED Ontology Definition]

**Cell type:** Cell type used in the experiment if non mixed. If mixed the targeted cell type should be used [MGED Ontology Definition]

**Clone information:** For each reporter, the identity of the clone along with information on the clone provider, the date obtained, and availability

**Common reference:** A hybridization to which all the other hybridizations have been compared

---

- 
- 1 **Composite sequence information:** The set of reporters contained in the composite sequence.  
2 The nucleotide sequence information for each composite element: number of oligonucleotides,  
3 oligonucleotide sequences (if given), and the reference sequence accession number (from  
4 relevant databases)
- 5 **Composite sequence related information:** Information on the set of reporters used collectively  
6 to measure an expression of a particular gene
- 7 **Compound:** A drug, solvent, chemical, etc., that can be measured [MGED Ontology Definition]
- 8 **Contact details for sample:** The resource (*e.g.*, company, hospital, geographical location) used  
9 to obtain or purchase the BioMaterial and the type of specimen [MGED Ontology Definition]
- 10 **Control elements position:** The position of the control features on the array
- 11 **Control elements related information:** Array elements that have an expected value and/or are  
12 used for normalization
- 13 **Control type:** The type of control used for the normalization and their qualifier
- 14 **Data processing protocol:** Documentation of the set of steps taken to process the data,  
15 including: the normalization strategy and the algorithm used to allow comparison of all data
- 16 **Developmental stage:** The developmental stage of the organism's life cycle during which the  
17 BioMaterial was extracted [MGED Ontology Definition]
- 18 **Disease state:** The name of the pathology diagnosed in the organism from which the  
19 BioMaterial was derived. The disease state is normal if no disease has been diagnosed [MGED  
20 Ontology Definition]
- 21 **Element dimensions:** The physical dimensions of each features
- 22 **Experiment description:** Free text description of the experiment and link to an electronic  
23 publication in a peer-reviewed journal
- 24 **Experiment design:** Experiment is a set of one or more hybridizations that are in some way  
25 related (*e.g.*, related to the same publication MIAME distinguishes between: the experiment  
26 design (the design, purpose common to all hybridizations performed in the experiment), the  
27 sample used (sample characteristics, the extract preparation and the labeling), the hybridization  
28 (procedures and parameters) and the data (measurements and specifications)
- 29 **Experiment type (s):** A controlled vocabulary that classify an experiment
- 30 **Experimental design:** Design and purpose common to all hybridizations performed in the  
31 experiment
- 32 **Experimental factor (s):** Parameter (s) or condition (s) tested in the experiment
- 33 **Extraction method:** The protocol used to extract nucleic acids from the sample
- 34 **Features related information:** Information on the location of the reporters on the array
- 35 **Final gene expression table (s):** Derived measurement value summarizing related elements and  
36 replicates, providing the type of reliability indicator used
-

- 1 **Gene name:** The gene represented at each composite sequence: name and links to appropriate  
2 databases (*e.g.* SWISS-PTOR or organism specific database)
- 3 **Genetic variation:** The genetic modification introduced into the organism from which the  
4 BioMaterial was derived. Examples of genetic variation include specification of a transgene or  
5 the gene knocked-out [MGED Ontology Definition]
- 6 **Growth conditions:** A description of the isolated environment used to grow organisms or parts  
7 of the organism [MGED Ontology Definition]
- 8 **Hybridization protocol:** Documentation of the set of steps taken in the hybridization,  
9 including: solution (*e.g.* concentration of solutes); blocking agent and concentration used; wash  
10 procedure; quantity of labelled target used; time; concentration; volume, temperature, and  
11 description of the hybridization instruments
- 12 **Hybridization extract preparation:** Information on the extract preparation for each extract  
13 prepared from the sample
- 14 **Hybridizations:** Procedures and parameters for each hybridization
- 15 **Image analysis and quantitation:** Each image has a corresponding image quantitation table,  
16 where a row represents an array design element and a column to a different quantitation types  
17 (*e.g.* mean or median pixel intensity)
- 18 **Image analysis output:** The complete image analysis output for each image
- 19 **Image analysis protocol:** Documentation of the set of steps taken to quantify the image  
20 including: the image analysis software, the algorithm and all the parameters used
- 21 **In vitro treatment:** The manipulation of the cell culture condition for the purposes of  
22 generating one of the variables under study and the documentation of the set of steps taken in the  
23 treatment
- 24 **In vivo treatment:** The manipulation of the organism for the purposes of generating one of the  
25 variables under study and the documentation of the set of steps taken in the treatment
- 26 **Individual genetic characteristics:** The genotype of the individual organism from which the  
27 BioMaterial was derived [MGED Ontology Definition]
- 28 **Individual number:** Identifier or number of the individual organism from which the  
29 BioMaterial was derived. For patients, the identifier must be approved by Institutional Review  
30 Boards (IRB, review and monitor biomedical research involving human subjects) or appropriate  
31 body [MGED Ontology Definition]
- 32 **Label incorporation method:** The label incorporation method used
- 33 **Label used:** The name of the label used
- 34 **Measurements:** MIAME distinguishes between three levels of data processing: image (raw  
35 data), image analysis and quantitation, gene expression data matrix (normalized and summarized  
36 data)
-



- 
- 1 **Normalized and summarized data:** Several quantitation tables are combined using data  
2 processing metrics to obtain the ‘final’ gene expression measurement table (gene expression data  
3 matrix) associated with the experiment
- 4 **Nucleic acid type:** The type of nucleic acid extracted (*e.g.* total RNA, mRNA)
- 5 **Number of elements on the array:** The number of features on the array
- 6 **Number of hybridizations:** Number of hybridizations performed in the experiment
- 7 **Organism:** The genus and species (and subspecies) of the organism from which the BioMaterial  
8 is derived [MGED Ontology Definition]
- 9 **Organism part:** The part or tissue of the organism's anatomy from which the BioMaterial was  
10 derived [MGED Ontology Definition]
- 11 **Platform type:** The technology type used to place the biological sequence on the array
- 12 **Production protocol:** A description of how the array was manufactured
- 13 **Provider:** The primary contact (manufacturer) for the information on the array design
- 14 **Qualifier, value, source (may use more than once):** Describe any further information about  
15 the array in a structured manner
- 16 **Quality control steps:** Measures taken to ensure or measure quality: replicates (number and  
17 description), dye swap (for two channel platforms) or others (unspecific binding, low complexity  
18 regions, polyA tails)
- 19 **Raw data:** Each hybridization has at least one image
- 20 **Relationship between samples and arrays:** Relationship between the labelled extract (related  
21 to which sample which extract) and arrays (design, batch and serial number) in the experiment
- 22 **Reporter and location:** The arrangement and the system used to specify the location of each  
23 features on the array (*e.g.* grid, row, column, zone)
- 24 **Reporter approximate length:** The approximate length of the reporter's sequence
- 25 **Reporter generation protocol:** A description of how the reporters were generated
- 26 **Reporter related information:** Information on the nucleotide sequence present in a particular  
27 location on the array
- 28 **Reporter sequence information:** The nucleotide sequence information for reporter: sequence  
29 accession number (from DDBJ/EMBL/GenBank), the sequence itself (if known) or a reference  
30 sequences (*e.g.* for oligonucleotides) and PCR primers pair information (if relevant)
- 31 **Reporter type:** Physical nature of the reporter (*e.g.* PCR product, synthesized oligonucleotide)
- 32 **Sample:** The biological material from which the nucleic acids have been extracted for  
33 subsequent labelling and hybridization. MIAME distinguishes between: source of the sample  
34 (bio-source), its treatment, the extract preparation, and its labeling
- 35 **Sample labeling:** Information on the labeling preparation for each labelled extract
- 36 **Scanner image file:** The TIFF file including header
-

- 1 **Scanning protocol:** Documentation of the set of steps taken for scanning the array and  
2 generating an image including: description of the scanning instruments and the parameter  
3 settings
- 4 **Separation technique:** Technique to separate tissues or cells from a heterogenous sample (e.g.  
5 trimming, microdissection, FACS)
- 6 **Sex:** Term applied to any organism able to undergo sexual reproduction in order to differentiate  
7 the individuals or type involved. Sexual reproduction is defined as the ability to exchange  
8 genetic material with the potential of recombinant progeny [MGED Ontology Definition]
- 9 **Single or double stranded:** Whether the reporter sequences are single or double stranded
- 10 **Spike type and qualifier:** The type of spike used (*e.g.* oligonucleotide, plasmid DNA,  
11 transcript) and its qualifier (e.g. concentration, expected ratio, labeling methods)
- 12 **Spiking control:** External controls added to the hybridization extract (s)
- 13 **Spiking control feature:** Position of the feature (s) on the array expected to hybridize to the  
14 spiking control
- 15 **Strain or line:** Animals or plants that have a single ancestral breeding pair or parent as a result  
16 of brother x sister or parent x offspring matings [MGED Ontology Definition]
- 17 **Surface and coating specification:** Type of surface and name for the type of coating used
- 18 **Targeted cell type:** The targeted cell type is the cell of primary interest. The BioMaterial may  
19 be derived from a mixed population of cells although only one cell type is of interest [MGED  
20 Ontology Definition]
- 21 **Treatment type:** The type of manipulation applied to the BioMaterial for the purposes of  
22 generating one of the variables under study [MGED Ontology Definition]

23  
24

---

---

## Appendix F: Additional Glossary from Genomics White Paper

**Allele:** An alternative form of a gene or any other segment of a chromosome

**Bioinformatics:** The analysis of biological information using computers and statistical techniques; the science of developing and utilizing computer databases and algorithms to accelerate and enhance biological research.

**Biomarker:** A molecular indicator of a specific biological property; a biochemical feature or facet that can be used to measure the progress of disease or the effects of treatment.

**Complementary DNA (cDNA):** DNA made from a messenger RNA (mRNA) template. The single-stranded form of cDNA is often used as a probe in physical mapping.

**Biotechnology:** Set of biological techniques developed through basic research and now applied to research and product development. In particular, biotechnology refers to the use by industry of recombinant DNA, cell fusion, and new bioprocessing techniques.

**Computational Toxicology (Comp Tox):** Word used first in EPA's Interim Policy on Genomics - "Computational Toxicology is defined as the application of models from computational and mathematical biology and computational chemistry for prediction and understanding mechanisms" - Computational Toxicology Framework Document, ORD, April 2003.

**DER:** Data Evaluation Record

**Deoxyribonucleic acid (DNA):** Nucleic acid that constitute the genetic material of all cellular organisms and DNA viruses. The genetic information is used in the synthesis of ribonucleic acids (RNAs) from DNA templates (transcription), and in the synthesis of proteins from messenger RNA (mRNA) templates (translation).

**DNA Microarray:** Microarray is a tool used to sift through and analyze the information contained within a genome. A microarray consists of different deoxyribonucleic acid (DNA) probes that are chemically attached to a substrate, which can be a microchip, a glass slide or a microsphere-sized bead.

**Expressed sequence tag:** A unique stretch of DNA within a coding region of a gene that is useful for identifying full-length genes and serves as a landmark for mapping.

**FACS:** Fluorescence Activated Cell Sorter

**Gene:** The fundamental physical and functional unit of heredity. A gene is an ordered sequence of nucleotides located in a particular position on a particular chromosome that encodes a specific functional product (*i.e.*, a protein or RNA molecule).

---

1  
2 **Gene chip technology:** Development of cDNA microarrays from a large number of genes. Used  
3 to monitor and measure changes in gene expression for each gene represented on the chip.  
4

5 **Gene expression:** Process by which a gene's coded information is converted into the structures  
6 present and operating in the cell. Expressed genes include those that are transcribed into mRNA  
7 and then translated into protein and those that are transcribed into RNA but not translated into  
8 protein (*e.g.*, transfer and ribosomal RNAs).  
9

10 **Genetics:** Study of inheritance patterns of specific traits.  
11

12 **Genetic testing:** Analyzing an individual's genetic material to determine predisposition to a  
13 particular health condition or to confirm a diagnosis of genetic disease.  
14

15 **Genomics:** Comprehensive study of whole sets of genes, gene products and their interaction.  
16

17 **Genome:** All the genetic material in the chromosomes of a particular organism; its size is  
18 generally given as its total number of base pairs.  
19

20 **Genotype:** The genetic composition of an organism or a group of organisms; a group or class of  
21 organisms having the same genetic constitution.  
22

23 **Hazard Assessment:** The process of determining whether exposure to an agent can cause an  
24 increase in the incidence of a particular adverse health effect (*e.g.*, cancer, birth defect) and  
25 whether the adverse health effect is likely to occur in humans.  
26

27 **Hazard Characterization:** A description of the potential adverse health effects attributable to a  
28 specific environmental agent, the mechanisms by which agents exert their toxic effects, and the  
29 associated dose, route, duration, and timing of exposure.  
30

31 **Hazard identification:** The process of determining whether it is scientifically correct to infer  
32 that toxic effects observed in one setting will occur in other settings (*e.g.*, whether substances  
33 found to be carcinogenic or teratogenic in experimental animals are likely to have the same  
34 results in humans).  
35

36 **In Vitro:** A biological study is one which is performed in isolation from a living organism (in  
37 contrast to In Vivo studies).  
38

39 **In Vivo:** A biological study is one which is performed within a living biological organism (as  
40 opposed to an In Vitro study).  
41

42 **Knockout:** Inactivation of specific genes. Knockouts are often created in laboratory organisms  
43 such as yeast or mice so that scientists can study the knockout or null organism as a model for a  
44 particular disease.  
45

---

---

1 **MAGE:** MicroArray and Gene Expression; the group aims to provide a standard for the  
2 representation of microarray expression data that would facilitate the exchange of microarray  
3 information between different data systems.

4  
5 **MAGE-OM:** Microarray Gene Expression: Object Model

6  
7 **MGED:** The Microarray Gene Expression Data (MGED) Society is an international  
8 organization of biologists, computer scientists, and data analysts that aims to facilitate the  
9 sharing of microarray data generated by functional genomics and proteomics experiments.

10  
11 **Mapping:** Charting the location of genes on chromosomes.

12  
13 **Mass spectrometry:** A method used to determine the masses of atoms or molecules in which an  
14 electrical charge is placed on the molecule and the resulting ions are separated by their mass to  
15 charge ratio.

16  
17 **Metabolome:** Entire complement of all the small molecular weight metabolites inside a cell  
18 suspension (or other sample) of interest (Aberystwyth, University of Wales Web site-  
19 <http://dbk.ch.umist.ac.uk/metabol.htm>). This profile is a product of the genome of the organism,  
20 the expression of that genome, and the operation of the metabolism is a particular part of the  
21 organism, in a particular environment.

22  
23 **Metabolomics:** Involves the systematic estimation of metabolomes from a range of organisms,  
24 followed by statistical analyses and other investigations of that large quantity of data.

25  
26 **Metabonomics:** Study of the endogenous composition of biofluids and tissues of an organism in  
27 order to probe the metabolic state in homeostasis, and when under interventional stress. Hector  
28 Keun (Biological Chemistry and Biological Sciences, Imperial College, London); Metabolic  
29 Profiling: Application to Toxicology and Risk Reduction. International Conference, May 14-15,  
30 2003, NIEHS, Research Triangle Park, North Carolina.

31  
32 **MIAME:** Minimum Information About a Microarray Experiment that is needed to enable the  
33 interpretation of the results of the experiment unambiguously and potentially to reproduce the  
34 experiment

35  
36 **Microarray:** A tool used to sift through and analyze the information contained within a  
37 genome. A microarray consists of different nucleic acid probes that are chemically attached to a  
38 substrate, which can be a microchip, a glass slide or a microsphere-sized bead.

39  
40 **Mode of Action:** Key events and processes, starting with the interaction of an agent with a cell,  
41 through functional and anatomical changes observed on the progression to toxicity

42  
43 **MOPS-EDTA:** [MOPS] 3-(N-Morpholino) propanesulfonic acid], [EDTA]  
44 ethylenediaminetetraacetic acid

45

---

1 **Northern blot:** A technique used to separate and identify RNA.

2  
3 **Nucleotide:** A subunit of DNA or RNA. To form a DNA or RNA molecule, thousands of  
4 nucleotides are joined in a long chain.

5  
6 **“Omics”:** Term including genomics, proteomics, metabonomics (some differentiate this term  
7 from metabolomics), transcriptomics, and associated bioinformatics (Environmental Health  
8 Perspectives, 110: 2002, 1047-1050; Meeting Report: Use of Genomics in Toxicology and  
9 Epidemiology: Findings and recommendations of a workshop). Carol J. Henry and Vanessa Vu,  
10 first and last authors, respectively.

11  
12 **Omics Technologies:** A quote often cited describes this phrase“...are based on comprehensive  
13 biochemical and molecular characterizations of an organism, tissue or cell type” Sumner *et. al.*  
14 2003.

15  
16 **Phenotype:** The observable physical or biochemical traits of an organism, as determined by  
17 genetics and the environment; the expression of a given trait based on phenotype; an individual  
18 or group of organisms with a particular phenotype.

19  
20 **PMT:** Photomultiplier tube; used in the capture of raw data

21  
22 **Polymorphism:** The quality or character of genes occurring in several different forms.

23  
24 **Proteome:** All of the proteins produced by a given species, just as the genome is the totality of  
25 the genetic information possessed by that species.

26  
27 **Proteomics:** Study of the function of all expressed proteins (Nature, 422: 2003, 193-197).

28  
29 **Quality policy statement:** Describing the specific objectives and commitment of the laboratory  
30 and its management to quality and data integrity. An ethics statement may be included at this  
31 point.

32  
33 **RNA:** Nucleic acid found in all living cells that plays a role in the transfer of information from  
34 DNA to the protein-forming system of the cell. The base sequence of an RNA is specified by the  
35 base sequence of a section of the DNA (a Gene) which is used as the template for RNA synthesis  
36 (transcription). (Dorland’s Medical Dictionary)

37  
38 **Risk Assessment** (in the context of human health): The evaluation of scientific information on  
39 the hazardous properties of environmental agents (hazard characterization), the dose-response  
40 relationship (dose-response assessment), and the extent of human exposure to those agents  
41 (exposure assessment). The product of the risk assessment is a statement regarding the  
42 probability that populations or individuals so exposed will be harmed and to what degree (risk  
43 characterization).

44

---

---

1 **Signal transduction pathway:** The course by which a signal from outside a cell is converted to  
2 a functional change within the cell.

3  
4 **Single nucleotide polymorphism (SNP):** A change in which a single base in the DNA differs  
5 from the usual base at that position.

6  
7 **Standard operating procedures (SOPs):** listing all routine laboratory operations documented  
8 and signed by management which are available to clients upon request and readily accessible to  
9 staff. Also known as laboratory operating procedures and protocols.

10  
11 **Susceptibility:** Increased likelihood of an adverse effect, often discussed in terms of relationship  
12 to a factor that can be used to describe a human subpopulation (*e.g.* life stage, demographic  
13 feature, or genetic characteristic).

14  
15 **Susceptible Subgroups:** May refer to life stages, for example, children or the elderly, or to  
16 other segments of the population, for example, asthmatics or the immune-compromised, but are  
17 likely to be somewhat chemical-specific and may not be consistently defined in all cases.

18  
19 **Systems Biology:** A holistic approach to the study of biology with the objective of  
20 simultaneously monitor all biological processes operating as an integrated system. Sumner *et.*  
21 *al.*, 2003.

22  
23 **Systems Toxicology:** "...involves the study of perturbation of organisms by chemicals and  
24 stressors, monitoring changes in molecular expression and conventional toxicological  
25 parameters, and iteratively integrating biological response data to describe the functioning  
26 organism".

27  
28 **Throughput:** Output or production, as of a computer program, over a period of time.

29  
30 **Toxicity:** Deleterious or adverse biological effects elicited by a chemical, physical, or biological  
31 agent.

32  
33 **Toxicology:** The study of harmful interactions between chemical, physical, or biological agents  
34 and biological systems.

35  
36 **Toxicogenomics:** The collection, interpretation, and storage of information about gene and  
37 protein activity in order to identify toxic substances in the environment, and to help treat people  
38 at the greatest risk of diseases caused by environmental pollutants or toxicants. Study of the  
39 roles that genes play in the biological responses to environmental toxicants and stressors  
40 (Environmental Health Perspective Toxicogenomics (NIEHS)).

41  
42 **Transgenic:** Having genetic material (DNA) from another species. This term can be applied to  
43 an organism that has genes from another organism.

---

44

**Web-based Glossary Sources**

- 1
  - 2
  - 3 1- National Center for Toxicogenomics (NCT, NIEHS) Glossary
  - 4 <<http://www.niehs.nih.gov/nct/glossary.htm>>
  - 5
  - 6 2- Human Genome Project Information Web Glossary
  - 7 <[http://www.ornl.gov/sci/techresources/Human\\_Genome/glossary/](http://www.ornl.gov/sci/techresources/Human_Genome/glossary/)>
  - 8
  - 9 3- Cambridge Healthtech Institute <<http://www.genomicglossaries.com/CONTENT/omes.asp>>
  - 10
  - 11 4- The Physical and Theoretical Chemistry Laboratory, Oxford University Chemical and Other
  - 12 Safety Information <<http://ptcl.chem.ox.ac.uk/MSDS/>>
  - 13
  - 14 5- NIH Glossary <<http://www.accessexcellence.org/AE/AEPC/NIH/gene27.html>>
  - 15
  - 16 6- Integrated Risk Information Systems (IRIS, EPA) Glossary
  - 17 <<http://www.epa.gov/iris/gloss8.htm>>
-



---

## Appendix G: Content and Instructional Goals for the Three Levels of EPA Genomics Technical Training:

### Level I: Introductory Modules – Molecular Biology Concepts

#### Modules 1-8

Goal: Provide the basic information necessary for understanding the more intricate assessments of cellular functions at the molecular level. Introduce gene arrays and discuss how genomics data may affect risk assessments in the future – this module will tie into EPA’s current Genomics Policy. Issues relating to how to communicate genomics information to risk managers and the public will be addressed.

Target Audience: Non-scientists and/or technical staff without training in biological sciences, such as:  
Managers from Office of Research and Development, Regional and Program Offices  
Regional Risk Managers (e.g., Remedial Project Managers, On Scene Coordinators)  
Attorneys  
Staff from Regional Office Programs (e.g., Air, Water, Waste, Pesticides, Community Involvement, Tribal Program)  
Staff from States and Tribes

Components: Cell structure and function  
DNA structure and replication  
RNA – Types, functions, transcription (gene expression)  
Proteins – General features, formation (translation)  
Gene Arrays – General principles and types  
Risk Assessment Concepts – Cancer and non-cancer risk, how genomics data may affect risk assessments in the future  
Regulatory Framework and Risk communication (different regulatory applications)

### Level II: Intermediate Level Modules – Techniques in Molecular Biology

#### Modules 9-12

Goal: Provide a general understanding of all of the various applications that may be currently considered by programs throughout EPA and is intended to support human health and ecological risk assessors. Specific modules for individual program applications are considered separately (see Level II Modules – Specific Applications for Molecular Tools)

Target Audience: Scientists and/or those likely to use genomics data generated by risk assessors are the audience. Modules are intended for staff who need a more in-depth understanding of how genomics data is generated, but do not necessarily need to generate that data to support decision-

---

1 making. Modules for specific applications will be developed (e.g., microbial source tracking,  
2 homeland security, field inspectors). Examples include:

3  
4 Laboratory Staff

5 Regional Laboratories

6 Office of Research and Development

7 Enforcement/Compliance Staff (e.g., Water programs, TMDLs, FIFRA)

8 Risk Assessors - Human Health and Ecological

9 Regional Offices

10 Office of Research and Development

11 Program Offices

12  
13 Components: Background on molecular techniques, such as:

14 Microarrays

15 DNA amplification using PCR and RT-PCR

16 Isolation kits

17 Restriction enzymes

18 Electrophoresis

19 DNA fingerprinting

20 Protein Analysis

21 Laboratory exercises using various molecular techniques (see above)

22 Techniques for specific applications, such as:

23 Microbial source tracking

24 Homeland security

25 Field inspection

26 Molecular Biology Approaches in Quantitative Risk Assessment

27  
28 **Level II: Intermediate Level Modules – Specific Applications for Molecular Tools**

29  
30 Module 13: Homeland Security

31 Module 14: Microbial and/or Bacterial Source Tracking

32 Module 15: Molecular Techniques to Assess Exposure in Environmental Media

33 Module 16: Molecular Techniques for Genetically Modified Crop Plant Inspectors

34 Goal: Reinforce information and techniques learned in the Level II Modules – Techniques in Molecular  
35 Biology, and to provide more in-depth knowledge and skills in the performance of molecular techniques.  
36 Each of these modules is focused on a separate and specific application of the molecular tools  
37 (introduced in modules 8-11) to support different programs and needs of the Agency and its staff. Each  
38 module is intended to provide technical training to staff to increase the breadth of scientific  
39 understanding that will assist in improving job competencies with respect to science in their particular  
40 program area.

41  
42 Target Audience: Same as Level II Modules – Techniques in Molecular Biology.

43  
44 Components: Technical training in particular program areas, focusing on research and tools currently  
45 under development by or through ORD. For example, Module 13: Microbial and/or Bacterial Source

---

---

1 Tracking will use a newly developed Guide on Tools for Microbial Source Tracking (Jorge  
2 Santodomingo, in preparation), which compares a number of molecular (RT-PCR, DNA finger printing)  
3 and non-molecular (antibiotic resistance) techniques for identifying pathogenic bacteria from in water.  
4 This information may be supplemented by laboratory exercises.  
5

### 6 **Level III: Advanced Modules**

7  
8 Module 17: Data Analysis (1) – Statistical Analysis

9 Module 18: Data Analysis (2) – Bioinformatics Approaches, Computational Toxicology

10 Module 19: Use of Molecular Biology in Mode-Of-Action Determinations

11 Module 20: Using Genomics Data in Chemical Hazard/Risk Assessment  
12

13 Overall Goal: Provide advanced-level knowledge on specific technical needs that scientists  
14 performing research or developing hazard/risk assessments associated with chemical  
15 registrations and other regulatory activities may face. Due to the novel and continually evolving  
16 nature of the genomics field, the advanced training modules will be flexible to account for these  
17 potential dynamic changes. As new technologies and applications appear, additional or existing  
18 training modules will be developed, enhanced and/or revised. Modules will also be flexible to  
19 meet the needs of the different EPA programs.  
20

21 Target Audience: Scientists and those likely to use genomics data to generate Risk Assessments,  
22 such as:

23       ORD Researchers

24       Program Office Risk Assessors  
25

### 26 **Modules 17 & 18 (Data Analysis 1 & 2)**

27  
28 Goal: Provide information to research scientists and program office risk assessors on  
29 computational toxicology, bioinformatics and statistics. The modules will focus on how to  
30 identify and interpret patterns within the large volumes of genomics data and assess data  
31 significance and accuracy, offering insight into the critical evaluation, including pros, cons and  
32 limitations of possible approaches.  
33

#### 34 Components – Module 17:

35 Statistical approaches to microarray data analysis including, but not limited to:

36       Bayesian statistics

37       Correlation

38       Clustering

39       Principle component analysis  
40

#### 41 Components – Module 18:

42 Computational toxicology and bioinformatic approaches and tools used to analyze genomics  
43 data.  
44

---

1 Models and molecular biological applications used to predict effects and understand the cascade  
2 of events leading to an effect and how statistical analyses fits together with other information to  
3 form a bigger picture.

4 Bioinformatics tools (algorithms and statistics) that will be used to discriminate unique signature  
5 and families of signatures indicative of stressors and groups of stressors.

6 Data access (i.e. data mining) and management of data

## 7 8 **Module 19: Use of Molecular Biology in Mode-Of-Action Determinations**

9  
10 *Goal:* Introduce the approaches for and limitations of data interpretation. This module will  
11 provide a link between molecular biology methods and information and the risk assessment  
12 process.

### 13 *Components:*

14 The module will present the general concept that an understanding of the key events  
15 associated with the production of adverse health outcomes at the molecular level could enhance  
16 our ability to predict these outcomes in a qualitative and quantitative sense. In addition,  
17 variability and other uncertainties (e.g., adaptive responses and homeostatic compensation)  
18 surrounding the analysis and interpretation of microarray data for making quantitative  
19 conclusions about effect/response levels will be discussed. The concept of mode-of-action  
20 (MOA including key events) will also be introduced. The different classes of MOA will be  
21 discussed: these will include genotoxicity, mutagenicity, receptor-mediated, cell killing  
22 regenerative cell proliferation, and mitogenic responses. Each of these will be discussed in terms  
23 of the current understanding at the molecular level. For example, what is cell signaling and how  
24 do changes affect cell function; what is apoptosis and how is it induced; what controls the cell  
25 cycle and how can it be abrogated; what is the mechanism for the induction of mutations and  
26 chromosome changes and the role of DNA repair and replication? These molecular  
27 underpinnings will allow for examples of key event pathways to be discussed and how chemicals  
28 might potentially impact the various pathways.

## 29 30 31 **Module 20: Using Genomics Data in Chemical Hazard/Risk Assessment**

32  
33 *Goal:* Provide guidance on the incorporation of genomics (microarray) data in a weight-of-  
34 evidence approach for hazard assessment. Present principles and pitfalls using simple case  
35 studies. Case studies will be flexible to meet the needs of the programs and offices, for example,  
36 case studies may focus on homeland security and microbial source tracking applications.

### 37 *Components:*

38 Case studies such as:

39 Examples where microarray data quality is high

40 Examples that demonstrate data concerns which could lead to erroneous conclusions

41 Case Studies should be developed to support the need of the programs and Regional offices, e.g.,  
42 homeland security, microbial source tracking, ambient water quality monitoring, etc. to support  
43 the use of microarray data or for other molecular-biology-based or “omics” approaches.

44 Examples include, but are not limited to:  
45

---

- 
- 1 • Demonstration of purported evidence that a particular chemical belongs to
  - 2 a particular class of hepatotoxins
  - 3 • Demonstration of purported evidence that chemical has characteristics of a
  - 4 certain class of hormonally active substances
  - 5

6 These Case Studies should include the following elements:

- 7
  - 8 • Purpose
  - 9 • Overall (microarray or other “omics” approach) study design
  - 10 • Purported mode of action and details of how data support proposal, including purported
  - 11 rationale for utility of microarray data; arguments to support conclusions
  - 12 • Conventional mechanistic support: histopathology, clinical chemistry, metabolic profile,
  - 13 time course to appearance of critical elements, dose-response information, special
  - 14 studies, etc.
  - 15 • Microarray data: summary gene expression profile data presentation and necessary
  - 16 supporting raw data, proposed up and down regulated and constituent gene identification,
  - 17 rationale for platform and chip design, demonstration of reproducibility, analysis of
  - 18 variability, positive and negative controls, dose response/temporal elements analysis,
  - 19 statistical analysis; RNA stability
  - 20 • Correlation and comparison: between conventional and microarray data to support
  - 21 argument; phenotypic anchoring
  - 22 • Other Evidence: Structure-Activity Relationship, etc.
  - 23 • Any perceived data gaps
  - 24 • Potential relevance to humans
  - 25 • Weight-Of-Evidence Conclusion
  - 26
-